



HAL
open science

Extreme value statistics of smooth Gaussian random fields

Stéphane Colombi, Olaf Davis, Julien Devriendt, Simon Prunet, Joe Silk

► **To cite this version:**

Stéphane Colombi, Olaf Davis, Julien Devriendt, Simon Prunet, Joe Silk. Extreme value statistics of smooth Gaussian random fields. *Monthly Notices of the Royal Astronomical Society*, 2011, 414, pp.2436-2445. 10.1111/j.1365-2966.2011.18563.x. insu-03645948

HAL Id: insu-03645948

<https://insu.hal.science/insu-03645948>

Submitted on 21 Apr 2022

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

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

Extreme value statistics of smooth Gaussian random fields

Stéphane Colombi,^{1*} Olaf Davis,² Julien Devriendt,² Simon Prunet¹ and Joe Silk^{1,2}

¹*Institut d'Astrophysique de Paris, CNRS UMR 7095 and UPMC, 98 bis, bd Arago, F-75014 Paris, France*

²*Department of Astrophysics, University of Oxford, Keble Road, Oxford OX1 3RH*

Accepted 2011 February 18. Received 2011 February 17; in original form 2010 December 22

ABSTRACT

We consider the Gumbel or extreme value statistics describing the distribution function $p_G(v_{\max})$ of the maximum values of a random field v within patches of fixed size. We present, for smooth Gaussian random fields in two and three dimensions, an analytical estimate of p_G which is expected to hold in a regime where local maxima of the field are moderately high and weakly clustered.

When the patch size becomes sufficiently large, the negative of the logarithm of the cumulative extreme value distribution is simply equal to the average of the Euler characteristic of the field in the excursion $v \geq v_{\max}$ inside the patches. The Gumbel statistics therefore represents an interesting alternative probe of the genus as a test of non-Gaussianity, e.g. in cosmic microwave background temperature maps or in 3D galaxy catalogues. It can be approximated, except in the remote positive tail, by a negative Weibull-type form, converging slowly to the expected Gumbel-type form for infinitely large patch size. Convergence is facilitated when large-scale correlations are weaker.

We compare the analytic predictions to numerical experiments for the case of a scale-free Gaussian field in two dimensions, achieving impressive agreement between approximate theory and measurements. We also discuss the generalization of our formalism to non-Gaussian fields.

Key words: methods: analytical – methods: statistical – large-scale structure of Universe.

1 INTRODUCTION

Gumbel or extreme value statistics are concerned with the extrema of samples drawn from random distributions (Gumbel 1958). In the case of sample means, the Central Limit Theorem states that the means of many samples of size N drawn from some distribution will be normally distributed in the large- N limit; analogously, it can be shown that in the same limit the cumulative distribution of the sample maximum or minimum v will tend to one of the family of the following functions:

$$G_{\gamma_G}(v) = \exp[-(1 + \gamma_G y)^{-1/\gamma_G}], \quad (1)$$

with

$$y = \frac{v - a}{b}, \quad (2)$$

where a and b are location and scale parameters (see e.g. Coles 2001). The shape parameter γ_G characterizes the distribution: a distribution with $\gamma_G = 0$ is known as having the ‘Gumbel type’,

$$G_0 = \exp[-\exp(-y)], \quad (3)$$

while $\gamma_G < 0$ and $\gamma_G > 0$ correspond, respectively, to forms of the ‘negative Weibull type’ and the ‘Fréchet type’.

Distributions given by equation (1) have seen application to time-series data in many fields such as climate (see e.g. Katz & Brown 1992), hydrology (see e.g. Katz, Parlange & Naveau 2002), seismology (see e.g. Cornell 1968), insurance and finance (see e.g. Embrechts & Schmidli 1994), etc., in predicting the incidence of extreme events from knowledge of past data. Here we consider applications to 2D and 3D random fields relevant to cosmology, but our approach is sufficiently general so that extension to other fields should not prove difficult.

In three dimensions, one is naturally interested in the occurrence of the most massive clusters in galaxy surveys (Bhavsar & Barrow 1985; Cayón, Gordon & Silk 2010; Holz & Perlmutter 2010; Davis et al. 2011) of large-scale mass concentrations (Yamila Yaryura, Baugh & Angulo 2010) such as the Sloan Great Wall (Gott et al. 2005) or the largest voids observed in the spatial distribution of galaxies. In two dimensions, the most obvious application concerns the temperature fluctuations in the cosmic microwave background (CMB; Mikelsohn, Silk & Zuntz 2009), in particular the analysis of the hottest hotspots (Coles 1988) and the coldest cold spots. There are several works that suggest the existence of anomalies in current CMB experiments (see e.g. Larson & Wandelt 2004; Ayaita et al. 2010), the most prominent one being the cold spot discovered

*E-mail: colombi@iap.fr

in the temperature maps measured by the *Wilkinson Microwave Anisotropy Probe* (WMAP; Vielva et al. 2004; Cruz et al. 2005).

In this work, we consider a random field in two or three dimensions, smoothed on some scale ℓ , and take large ‘patches’ of size $L \gg \ell$. The values of the field at all points inside a given patch constitute one sample, and the extreme value of the field in the patch is our statistic of interest (henceforth we restrict ourselves in particular to the maxima, though the case of patch minima is exactly analogous).

Although some of the results present in this paper also apply to the non-Gaussian case, we restrict our practical calculations to Gaussian fields of known power spectrum and ask how to derive analytically the distribution of patch maxima and how to explicitly relate the results to equations (1) and (2). This effort is not new: the calculation of the extreme value distribution of Gaussian fields has been paid a lot of attention by mathematicians, not only for time-series (see e.g. Leadbetter, Lindgren & Rootzén 1983) but also in larger number of dimensions (see e.g. Bickel & Rosenblatt 1973; Rosenblatt 1976; Adler 1981; Aldous 1989; Piterbarg 1996; Adler & Taylor 2007), leading to a number of important results. In particular, in the very large patch limit, convergence to the Gumbel-type distribution (3) was rigorously demonstrated (e.g. Bickel & Rosenblatt 1973). Thus, many of the results derived in this paper can be found in the mathematical literature, but we recall them for the sake of completeness as they are needed to understand how the statistics behave in various regimes.

The method we employ to estimate the extreme value statistics relies on a local maxima approach and was already utilized in a more rigorous mathematical set-up (e.g. Adler & Taylor 2007). The central point is the observation that the probability of the patch maximum being below some specified density threshold is exactly the probability of encountering zero-points *above* that threshold. If we also identify the highest point in the patch with the field highest *peak* there (an assumption which is in fact non-trivial in general; see Section 2.2.1) then the problem is reduced to that of finding the void probability for peaks as a function of threshold.

The distribution and clustering statistics of peaks has achieved a good deal of attention in the astrophysical literature, in part due to their role as the nucleation points of rich clusters of galaxies when the field in question is that of matter overdensities. Combining results obtained previously for the peak abundances and their correlation functions allows us to predict the void probability for both 2D and 3D fields, and hence the extreme value statistics. Once again, although most of the results can be found in the mathematical literature, the key novelty of the present paper lies in the derivation and test of approximate predictions in an intermediate regime where the patch size is large enough compared to the coherence length of the field, but *not so large* that either the asymptotic limit expected for Gaussian random fields (Bickel & Rosenblatt 1973; equation 3) or the Poisson regime (Aldous 1989) has been reached.

This paper is organized as follows. Section 2 sets up the general framework, with a few definitions followed by general results. In particular, the Gumbel statistics is related to the void probability, which is expressed in terms of average number density of peaks above some threshold and their N -point correlation functions. Section 3 focuses on the Gaussian field case, where a general estimate taking into account full clustering of the peaks is performed and explicit asymptotic formulae are derived and related to equation (1). Convergence to equation (3) for Gaussian random fields is recovered when the patch size tends to infinity. An explicit link to the Euler characteristic is also established, in agreement with the literature.

In Section 4, we test the theoretical predictions against numerical experiments in the 2D case. Finally, Section 5 summarizes the results obtained in this paper and discusses their generalization to non-Gaussian fields.

2 THEORY

2.1 Definitions

We consider, in a D -dimensional space with $D = 2$ or 3 , a random field $\delta(x)$ of zero average. We suppose that this field is statistically stationary (invariance of the N -point correlation functions by translation) and isotropic (invariance of the N -point correlation functions by rotation).

2.1.1 Smoothing window

This field is smoothed with a window of size ℓ :

$$\delta_\ell = \delta * W_\ell(x). \quad (4)$$

For instance, the Gaussian smoothing window $W_\ell(x) \propto \exp(-x^2/2\ell^2)$, which we shall choose for all practical calculations, reads in Fourier space

$$W_\ell(k) = \exp[-(k\ell)^2/2]. \quad (5)$$

The top-hat smoothing window will be needed as well, on scales $L \gg \ell$. In three dimensions, it is a sphere of radius L which reads in Fourier space

$$W_L(k) = 3[\sin(kL) - kL \cos(kL)]/(kL)^3. \quad (6)$$

In two dimensions, it is a disc of radius L which reads in Fourier space

$$W_L(k) = \frac{2J_1(kL)}{kL}, \quad (7)$$

where J_1 is the Bessel function of the first kind and of first order:

$$J_1(x) = \frac{1}{\pi} \int_0^\pi \cos[y - x \sin(y)] dy. \quad (8)$$

2.1.2 Gumbel statistics

From now on we measure the height of the field using the density contrast in units of its standard deviation,

$$v \equiv \frac{\delta_\ell}{\sigma_0}, \quad (9)$$

with

$$\sigma_0^2 \equiv \langle \delta_\ell^2 \rangle. \quad (10)$$

We consider a large spherical patch of size L at random position x_0 , $L \gg \ell$, and measure in that patch the maximum value of the smoothed density field:

$$v_{\max} \equiv \max \{v(x); |x - x_0| \leq L\}. \quad (11)$$

The goal is to study the Gumbel statistics, i.e. the probability distribution function $p_G(v_{\max}) dv_{\max}$ of the values of v_{\max} when we choose an infinite number of random realizations of x_0 . This distribution contains a dependence on the choice of smoothing and on the size of the patch,

$$p_G(v_{\max}) = p_G(v_{\max}, \ell, L), \quad (12)$$

which we leave implicit in the remainder of the paper.

2.2 General results

In a sufficiently non-degenerate field δ_ℓ , the set of local maxima – the peaks of the density field – is a discrete ensemble of points of positions p_i and density $\nu_i = \delta_i/\sigma_0$. We shall assume that this property is valid in all the subsequent calculations.

2.2.1 Fundamental assumption

The fundamental assumption we make is that the maximum of the density in a patch can be approximated by the density at the highest peak contained in the patch:

$$\nu_{\max} \simeq \bar{\nu}_{\max} \equiv \max \{ \nu_i, |p_i - x_0| \leq L \}. \quad (13)$$

This assumption is valid only in the regime where $L \gg \ell$. Indeed, there can be local maxima outside the patch but sufficiently close to its edge such that the density measured at a point on the edge of the patch is larger than the maximum density measured in the set of peaks contained in the patch. In other words, if we consider the population of local maxima of the density field defined inside the $(D - 1)$ -dimensional manifold given by the border of the patch of densities $\hat{\nu}_j$, then we have, in reality,¹

$$\nu_{\max} = \max(\bar{\nu}_{\max}, \hat{\nu}_j) \geq \bar{\nu}_{\max}. \quad (14)$$

Obviously, one expects ν_{\max} to approach $\bar{\nu}_{\max}$ as the ratio L/ℓ increases and the ratio of the patch volume to area near its edge decreases.

2.2.2 General expression of the Gumbel statistics

Let us define the cumulative Gumbel distribution by

$$P_G(\nu) \equiv \text{Prob.}(\nu_{\max} \leq \nu) \equiv \int_{-\infty}^{\nu} p_G(\nu_{\max}) d\nu_{\max}. \quad (15)$$

Such a probability, given the assumptions of Section 2.2.1, is also the probability that none of the local maxima contained in the patch is above the threshold. In other words, if we consider the population of local maxima satisfying $\nu_i > \nu$, none of them belongs to the patch. This happens with a probability $P_0(\nu)$, where P_0 is the probability of finding no maxima with normalized density larger than ν inside a spherical cell or a disc of radius L :

$$P_G(\nu) = P_0(\nu); \quad (16)$$

hence

$$p_G(\nu) = \frac{dP_0}{d\nu}. \quad (17)$$

The calculation of such a void probability can be performed using standard count-in-cell formalism if the number density $n(\nu_i > \nu)$ and the connected N -point correlation functions,² $\xi_N^p(x_1, \dots, x_N)$, of local maxima above the threshold, are known (White 1979; Fry 1985; Balian & Schaeffer 1989; Szapudi & Szalay 1993).

In particular, one can define the averaged correlations over a cell of size L and volume $V = (4\pi/3)L^3$ or area $V = \pi L^2$,

$$\bar{\xi}_N^p(L) \equiv \frac{1}{V} \int_V d^D x_1 \cdots d^D x_N \xi_N^p(x_1, \dots, x_N), \quad (18)$$

¹ See Adler & Taylor (2007) for a rigorous formulation corresponding to a more general patch shape than just a sphere.

² The connected N -point correlation functions are equal to zero for a Gaussian field if $N \geq 3$.

the normalized cumulants,

$$S_N^p(L) \equiv \frac{\bar{\xi}_N^p(L)}{\bar{\xi}_2^p(L)^{N-1}}, \quad S_1^p \equiv S_2^p \equiv 1, \quad (19)$$

and

$$N_c \equiv nV \bar{\xi}_2^p(L). \quad (20)$$

Each of these expressions contains an implicit ν -dependence. The number N_c represents the typical number of peaks above the threshold per overdense patch in excess to the average. It measures the deviation from a pure Poisson distribution due to clustering.

With these definitions, the void probability can be written as

$$P_0(\nu_{\max}) = \exp[-nV\sigma(N_c)], \quad (21)$$

with

$$\sigma(y) = \sum_{N \geq 1} (-1)^{N-1} \frac{S_N^p}{N!} y^{N-1}. \quad (22)$$

The challenge is now to relate the statistical properties of the local maxima to that of the underlying density field. This is made difficult by the fact that the void probability depends on the full hierarchy of correlations up to any order: in particular one has to relate the N -point correlation functions of the peaks to the N -point correlation functions of the underlying density field. We denote the latter by $\xi_N(x_1, \dots, x_N)$, and similarly the averaged N -point correlation functions of the density field by $\bar{\xi}_N(L)$.

This exercise has been performed in detail for random Gaussian fields by Bardeen et al. (1986) (hereafter BBKS) and Bond & Efstathiou (1987) (hereafter BE) in the 3D and 2D cases, respectively, extending earlier calculations of Kaiser (1984) and Politzer & Wise (1984). Note that these latter calculations did not consider statistics of peaks, but more generally of regions of δ_ℓ above the density threshold. However, in the rare event regime considered here, the two approaches should become equivalent (this is discussed in detail in BBKS).

The non-Gaussian case has been examined as well for a quite general class of hierarchical models by (Bernardeau & Schaeffer 1999, hereafter BS). The statistics under consideration in that work was that of overdense cells of size ℓ and not of peaks of the density field smoothed with a Gaussian window of size ℓ . Again, the approach of BS should give the same results as those obtained for peaks in the rare event regime.

2.2.3 Asymptotic expression of $\sigma(y)$

A fundamental result of the calculations of BS is that in the high-peak limit, i.e. $\nu \gg 1$, and at large enough separations, i.e.

$$\frac{\xi_2(x_i, x_j)}{\sigma_0^2} \ll 1, \quad (23)$$

$$\xi_N^p(x_1, \dots, x_N) \simeq \sum_{\text{trees}} \sum_{\text{labels}} \prod \xi_2^p(x_i, x_j), \quad (24)$$

in the notation of these authors. This expression is valid at least in the framework of the minimal hierarchical-tree model. The trees refer to ensembles of distinct pair associations of elements in the ensemble $\{1, \dots, N\}$ such that a fully connected structure containing exactly the N elements is constructed without any loop. The labels take into account all the possible combinations of elements in $\{1, \dots, N\}$ that lead to the same tree topology. In each tree topology, there are always $N - 1$ links, by definition. The total number of combinations of all the trees and the labels yields N^{N-2} possibilities. Fig. 1 of

BS can be examined to understand the process. For instance, for the four-point correlation function there are two tree topologies: (i) the ‘star’ where one point is connected to all the others and (ii) the ‘line’ where one point is connected to one or two others depending on its position (at the end or in the middle). There are four possible labellings of the star and 12 possible labellings for the line.

Equation (24) also applies to Gaussian fields, independently of the shape of the smoothing window, at least if $\nu \gg 1$ and the following condition, more restrictive than (23) holds:

$$\nu^2 \frac{\xi_2(x_i, x_j)}{\sigma_0^2} \ll 1. \quad (25)$$

Indeed, the unconnected part of the N -point correlation function (the moment), μ_n , is given by

$$\mu_n = \prod_{i>j} [\xi_2^p(x_i, x_j) + 1] \quad (26)$$

in the high-threshold regime (Politzer & Wise 1984; Cline et al. 1987). Extracting the connected part from this expression consists exactly in extracting the ensemble of distinct pair associations in $\{1, \dots, N\}$ such that the corresponding topology is fully connected. In the large-separation limit, i.e. at leading order in $\xi_2^p(x_i, x_j)$, or equivalently if condition (25) is verified, only the tree topologies remain (because they correspond to the minimum power in ξ_2^p while being fully connected), and each label for each tree is given the same weight in equation (26), hence leading to equation (24) in that regime.

Equation (24) reads, after volume averaging in a sphere of radius L (BS),

$$S_N^p(L) \simeq N^{N-2}. \quad (27)$$

This result applies as well to the general tree-hierarchical model (BS). This means that the function σ defined in equation (22) reads (BS)

$$\sigma(y) = \left(1 + \frac{1}{2}\theta\right) e^{-\theta}, \quad \theta e^\theta = y. \quad (28)$$

Note that when $N_c \ll 1$, which occurs at some point for a large enough value of ν for which there are very few peaks above the threshold in average per patch,

$$\sigma(N_c) \simeq 1 - N_c/2 \quad (29)$$

by definition (equation 22). Therefore, even though we expect the low-end tail of the Gumbel statistics to be affected by the potential crudeness of our approximation of the function $\sigma(y)$, the high-end tail should still be quite well described.

3 THE GAUSSIAN FIELD CASE

Once the function $\sigma(y)$ is specified, one needs to carry out the calculation of the number density of peaks above the threshold as well as their averaged two-point correlation function. The detailed expressions can be found for a Gaussian field in BBKS and BE.

3.1 Shape parameters of the power spectrum: γ and R_*

The important parameters that control the number density of peaks above threshold ν and their two-point correlation function are the moments

$$\sigma_j^2 \equiv \int \frac{k^{D-1} dk}{2\pi^{D-1}} P(k) W_\ell^2(k) k^{2j}, \quad (30)$$

where $D = 2$ or 3 is the dimension of the space considered. Then BBKS and BE define the coherence parameter γ as

$$\gamma \equiv \frac{\sigma_1^2}{\sigma_0 \sigma_2} \quad (31)$$

and the scalelength R_* as

$$R_* = \sqrt{D} \frac{\sigma_1}{\sigma_2}. \quad (32)$$

For a Gaussian smoothing window and a scale-free power spectrum $P(k)$ given by

$$P(k) = Ak^n, \quad (33)$$

the integrals in equation (30) can be performed analytically, yielding the simple expressions

$$\sigma_0^2 = \left(\frac{\ell}{\ell_0}\right)^{-(n+D)}, \quad \ell_0 = \left[\frac{A}{4\pi^{D-1}} \Gamma\left(\frac{n+D}{2}\right)\right]^{1/(n+D)}, \quad (34)$$

$$\gamma = \sqrt{\frac{n+D}{n+D+2}}, \quad R_* = \sqrt{\frac{2D}{n+D+2}} \ell, \quad (35)$$

valid for $n > -D$.³

For a top-hat smoothing the scaling law (34) remains valid with a different correlation length, which for the 3D case can be written (see e.g. Lokas et al. 1996) as

$$\ell_0 = \left\{ \frac{9A}{8\pi^{3/2}} \frac{\Gamma[(n+3)/2] \Gamma[(1-n)/2]}{\Gamma[1-n/2] \Gamma[(5-n)/2]} \right\}^{1/(n+3)}. \quad (36)$$

On the other hand, we did not find any simple analytic expression of the correlation length for a top-hat smoothing in two dimensions.

3.2 Number density of peaks

The number density of peaks above the threshold is

$$n(\nu) = \int_\nu^\infty d\nu' \mathcal{N}(\nu'). \quad (37)$$

This integral is easily performed numerically, using the fact that the function $\mathcal{N}(\nu)$ is given by

$$\mathcal{N}(\nu) = \frac{1}{(2\pi)^{(D+1)/2} R_*^D} e^{-\nu^2/2} G_D(\gamma, \gamma\nu) \quad (38)$$

in D dimensions, with G_3 approximated by equation (4.4) of BBKS and G_2 given by equation (A1.9) of BE. For completeness, we rewrite these equations here:

$$G_3(\gamma, w) \simeq \frac{w^3 - 3\gamma^2 w + (Bw^2 + C_1) \exp(-Aw^2)}{1 + C_2 \exp(-C_3 w)}, \quad (39)$$

with

$$A = \frac{5}{2(9 - 5\gamma^2)}, \quad (40)$$

$$B = \frac{432}{\sqrt{10\pi}(9 - 5\gamma^2)^{5/2}}, \quad (41)$$

$$C_1 = 1.84 + 1.13(1 - \gamma^2)^{5.72}, \quad (42)$$

$$C_2 = 8.91 + 1.27 \exp(6.51\gamma^2), \quad (43)$$

³ The case $D = 3$ was derived in BBKS, with the expression for the correlation length ℓ_0 given in e.g. Lokas et al. (1996).

$$C_3 = 2.58 \exp(1.05\gamma^2) \quad (44)$$

and

$$G_2(\gamma, w) = (w^2 - \gamma^2) \left\{ 1 - \frac{1}{2} \operatorname{erfc} \left[\frac{w}{\sqrt{2(1-\gamma^2)}} \right] \right\} + \frac{w(1-\gamma^2)}{\sqrt{2\pi(1-\gamma^2)}} \exp \left[-\frac{w^2}{2(1-\gamma^2)} \right] + \frac{1}{\sqrt{3-2\gamma^2}} \exp \left(-\frac{w^2}{3-2\gamma^2} \right) \times \left\{ 1 - \frac{1}{2} \operatorname{erfc} \left[\frac{w}{\sqrt{2(1-\gamma^2)(3-2\gamma^2)}} \right] \right\}. \quad (45)$$

3.3 Correlation function of peaks

In the large-separation regime the two-point correlation function of the peaks (25) reads, if one neglects contributions from higher order derivatives of $\xi_2(r)$, which is a fair approximation according to BBKS and BE if ξ_2 is a power law of negative index,

$$\xi_2^p = \frac{\langle \tilde{\nu} \rangle^2}{\sigma_0^2} \xi_2, \quad (46)$$

where

$$\langle \tilde{\nu} \rangle = \frac{\int_v^\infty \tilde{\nu}(v') \mathcal{N}(v') dv'}{\int_v^\infty \mathcal{N}(v') dv'}, \quad (47)$$

and the effective threshold $\tilde{\nu}(v)$ writes

$$\tilde{\nu} = v - \frac{\gamma\theta}{1-\gamma^2}, \quad (48)$$

with θ approximated by equation (6.14) of BBKS in three dimensions:

$$\theta \simeq \frac{3(1-\gamma^2) + (1.216 - 0.9\gamma^4) \exp[-\gamma/2(\gamma\nu/2)^2]}{\sqrt{3(1-\gamma^2) + 0.45 + (\gamma\nu/2)^2} + \gamma\nu/2}, \quad (49)$$

and given in two dimensions by

$$\theta = (1-\gamma^2) \frac{H(\gamma, \gamma\nu)}{G_2(\gamma, \gamma\nu)}, \quad (50)$$

where $H(\gamma, w) \equiv \partial G_2/\partial w$ is given by equation (A4.7a) of BE. For completeness,

$$H(\gamma, w) = 2w \left\{ 1 - \frac{1}{2} \operatorname{erfc} \left[\frac{w}{\sqrt{2(1-\gamma^2)}} \right] \right\} + \frac{4(1-\gamma^2)^2}{(3-2\gamma^2)\sqrt{2\pi(1-\gamma^2)}} \exp \left[-\frac{w^2}{2(1-\gamma^2)} \right] - \frac{2w}{(3-2\gamma^2)^{3/2}} \exp \left(-\frac{w^2}{3-2\gamma^2} \right) \times \left\{ 1 - \frac{1}{2} \operatorname{erfc} \left[\frac{w}{\sqrt{2(1-\gamma^2)(3-2\gamma^2)}} \right] \right\}. \quad (51)$$

In the large-threshold limit, we simply have

$$\langle \tilde{\nu} \rangle \rightarrow v, \quad v \rightarrow \infty, \quad (52)$$

as derived by Kaiser (1984). After averaging over volume V , the expression of $\bar{\xi}_2^p(L)$ is thus simply given by

$$\bar{\xi}_2^p(L) = \frac{\langle \tilde{\nu} \rangle^2}{\sigma_0^2} \bar{\xi}_2(L), \quad (53)$$

where $\bar{\xi}_2(L)$ is the averaged two-point correlation function of the underlying density field. It can be derived easily from the power spectrum of the underlying (smoothed) density field, δ_ℓ , using equation (30) with the top-hat window and replacing ℓ with L . The largeness of the patch size, L , compared to the smoothing scale, ℓ , should guarantee that the large-separation approximation (25) is verified in practice. This can be checked a posteriori by examining the range of values of ν where $p_G(\nu)$ is significant.

Note that Heavens & Sheth (1999) and Desjacques (2008) (the latter in the large-separation limit) performed the exact calculation of the two-point correlation function of peaks in two and three dimensions, respectively, by taking into account corrections depending on second and fourth derivative of the two-point correlation function of the underlying field. These corrections can be significant on large scales, where the power spectrum of the underlying field significantly deviates from a power law. However, they get progressively smaller with increasing threshold ν , and in the rare event limit in which we are working here, they are probably irrelevant, except perhaps for the low- ν tail of the Gumbel statistics. Still, this assumption should be checked explicitly for the cold dark matter case by assessing the differences in the distribution introduced by computing the two-point correlation function of peaks with and without them. Such an investigation is beyond the scope of this paper, but it should be kept in mind when comparing analytic estimates of the Gumbel statistics to real data. Also note that in this more rigorous context, the simple proportionality relation (53) does not apply anymore.

3.4 Asymptotic regime

A particularly interesting case corresponds to the regime $\nu \gg 1$ and the Poisson limit $N_c \ll 1$, where $\sigma(N_c) \simeq 1$. Such a regime is expected to be reached if L/ℓ is sufficiently large and has been studied previously (Aldous 1989). However, since they will prove to be very useful to understand the intermediate (as opposed to asymptotic) regime that we discuss later, we recast the main results using our formalism. Here, the calculations will be facilitated by examining the cumulative Gumbel distribution, $P_G(\nu)$.

In the large- ν regime, the number density of peaks is proportional to the Euler characteristic \mathcal{E} of the underlying density field⁴ (e.g. BBKS, BE). Just how large a value of ν is required for this to be true depends on the level of accuracy one aims to reach in the description of the function $P_G(\nu)$. For instance, BBKS suggest $\gamma\nu > 2.5$ in three dimensions for a 10 per cent accuracy on approximating $n(\nu)$ by the Euler characteristic.

With the additional assumption that $N_c \ll 1$, equations (4.14) of BBKS and equation (3.3) of BE read, respectively, in three and two dimensions,

$$P_{G,3}(\nu) \simeq \exp(-\mathcal{E}_3 V) = \exp \left[-U_3(\nu^2 - 1) \exp \left(-\frac{\nu^2}{2} \right) \right], \quad (54)$$

$$P_{G,2}(\nu) \simeq \exp(-\mathcal{E}_2 V) = \exp \left[-U_2 \nu \exp \left(-\frac{\nu^2}{2} \right) \right], \quad (55)$$

with

$$U_D = \frac{\gamma^D V}{(2\pi)^{(D+1)/2} R_*^D}. \quad (56)$$

⁴ The Euler characteristic is seen here as an alternate count of critical point number densities of various kinds included in the excursion in regions with normalized density larger than ν , see e.g. Colombi, Pogosyan & Souradeep (2000) and Adler & Taylor (2007).

For scale-free power spectra we have

$$U_D = \left(\frac{4}{3}\right)^{D-2} \frac{\pi}{(2\pi)^{(D+1)/2}} \left(\frac{n+D}{2D}\right)^{D/2} \left(\frac{L}{\ell}\right)^D. \quad (57)$$

Note that in the original derivation (Aldous 1989), the right-hand part of equation (54) contains a term scaling like v^2 and not $v^2 - 1$. Recall however that these expressions still assume v to be sufficiently large compared to unity. Note also that in the very large threshold limit, one recovers the classical result (see Adler 1981; Adler & Taylor 2007)

$$1 - P_{G,D}(v \gg 1) \simeq \mathcal{E}_D V. \quad (58)$$

In the calculations presented in Adler & Taylor (2007), though, the edge effects discussed in Section 2.2.1 are not neglected, i.e. $\mathcal{E}_D V$ must be in fact viewed as the ensemble average of the Euler characteristic of the intersection between the excursion and the patch.⁵

An interesting value of v corresponds to

$$n(v_*)V = 1, \quad (59)$$

or $P_G = 1/e$. Obviously we must have v_* sufficiently large compared to unity for equations (54) and (55) to hold, as well as

$$N_c(v_*) = \frac{v_*^2}{\sigma_0^2} \xi_2(L) \ll 1, \quad (60)$$

to remain in the Poisson limit. The last condition imposes a constraint on the size L of the patch, which must be large enough compared to the smoothing scale ℓ . This obviously depends on spectral index: the ratio L/ℓ should be larger if n is small since there is more power on large scale.

Asymptotically, v_* reads

$$v_* \simeq \sqrt{2 \ln U_D} \left[1 + \frac{(D-1) \ln(2 \ln U_D)}{4 \ln U_D} \right]. \quad (61)$$

This equation shows that v_* grows rather slowly with L/ℓ .

Now we compare the expressions (54) and (55) to the standard cumulative form (1). To determine the parameters a , b and γ_G of equations (1) and (2), we perform a second-order Taylor expansion near $y = 0$, where $G_{\gamma_G}(y = 0) = 1/e$.

At second order in γ_{GY} we have

$$\ln(-\ln G_{\gamma_G}) \simeq -y + \gamma_G \frac{y^2}{2}. \quad (62)$$

Similarly we have

$$\begin{aligned} \ln(-\ln P_{G,3}) &\simeq -\frac{a^2}{2} + \ln[U(a^2 - 1)] \\ &\quad - \frac{ab(a^2 - 3)}{a^2 - 1} y - \frac{b^2(a^4 + 3)}{(a^2 - 1)^2} \frac{y^2}{2}, \end{aligned} \quad (63)$$

$$\begin{aligned} \ln(-\ln P_{G,2}) &\simeq -\frac{a^2}{2} + \ln(Ua) \\ &\quad - \frac{b(a^2 - 1)}{a} y - \frac{b^2(a^2 + 1)}{a^2} \frac{y^2}{2}. \end{aligned} \quad (64)$$

Our particular choice of the expansion is convenient since it implies

$$a = v_*. \quad (65)$$

⁵ Note that having a very accurate determination of the high- v tail of the Gumbel distribution has been paid a lot of attention by mathematicians and numerous methods have been employed to do so, estimating the Euler characteristic being one amongst them (see the introduction of Azaïs & Delmas 2002, for a panorama on various methods).

Then

$$b_3 = \frac{1}{v_*} \frac{v_*^2 - 1}{v_*^2 - 3}, \quad (66)$$

$$b_2 = \frac{v_*}{v_*^2 - 1}, \quad (67)$$

$$\gamma_{G,3} = -\frac{v_*^4 + 3}{v_*^2(v_*^2 - 3)^2} < 0, \quad (68)$$

$$\gamma_{G,2} = -\frac{v_*^2 + 1}{(v_*^2 - 1)^2} < 0, \quad (69)$$

where the labels 2 and 3 refer to the number of dimensions considered, D . It is interesting to study the asymptotic values of these parameters when $\ln U_D \gg 1$; hence when v_* is very large,

$$b \sim 1/v_*, \quad \gamma_G \sim -1/v_*^2, \quad (70)$$

thus

$$g_{\gamma_G} \sim \frac{v}{v_*} - 1. \quad (71)$$

We can thus understand that the range of validity of the Taylor expansion translated in terms of the v variable is $v \in [v_*(1 - \varepsilon), v_*(1 + \varepsilon)]$, where ε is a fraction of unity, for both $P_G(v)$ and $G_{\gamma_G}(v)$. While v_* does not in general correspond exactly to the position of the maximum of $p_G(v)$ and

$$g_{\gamma_G}(v) \equiv \frac{dG_{\gamma_G}}{dv}, \quad (72)$$

it is rather close to it, and increasingly so as L/ℓ becomes larger. This means in practice that the functions $p_G(v)$ and $g_{\gamma_G}(v)$ should be well described close to their maximum by our second-order Taylor expansion in an interval corresponding to confidence levels up to ~ 90 – 95 per cent. Hence, the functions $p_G(v)$ and $g_{\gamma_G}(v)$ should match each other quite well in that interval, but not in the tails, especially in the large- v one.

Note as well that γ_G is negative, as measured experimentally by e.g. Mikelson et al. (2009) and that it converges to zero, so that the function $P_G(v)$ converges to equation (3) as demonstrated more rigorously by e.g. Bickel & Rosenblatt (1973) (see also Rosenblatt 1976).

This asymptotic result was first exploited in cosmology by Coles (1988) to analyse the hottest hotspots in temperature fluctuations of the CMB. However, it is important to note that convergence to form (3) is rather slow. On the other hand, form (1) with values of a , b and γ_G given by equations (65), (66), (67), (68) and (69) along with implicit equation (59) to determine v_* remains always a good fit of equations (54) and (55) in the 90–95 per cent confidence region, but again, not in the tails. These analytical results will be illustrated explicitly in Section 4.

4 MEASUREMENTS

To check the validity of the theoretical calculations, we performed numerical experiments in the 2D case. We generated scale-free random Gaussian fields on sets of 100 realizations on a grid of size 4096^2 for each value of the spectral index we considered, $n = 0, -0.5, -1$ and -1.5 . Smoothing was performed with a Gaussian window of size $\ell = 5$ pixels and 400 non-overlapping circular patches of radius $L = 100$ pixels were extracted from each realization, amounting to a grand total of 40 000 patches to measure $p_G(v)$ for each value of n .

The results are displayed in Fig. 1 for $n = 0, -0.5$ and in Fig. 2 for $n = -1.0$ and -1.5 . Agreement between the measurements and

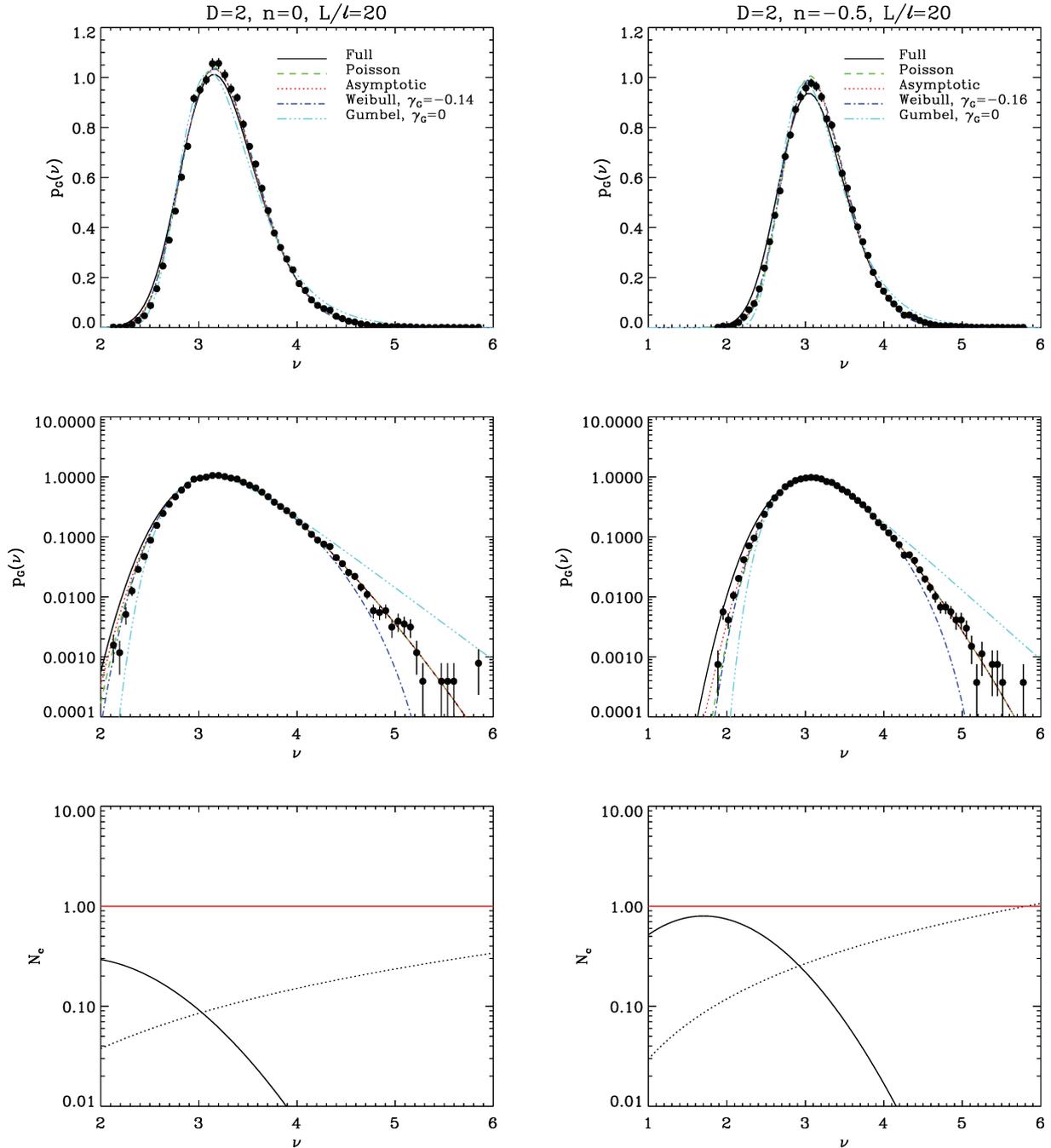


Figure 1. The Gumbel statistics measured in the case $L/l = 20$ for 2D random Gaussian fields with power spectra $P(k) \propto k^n$. The spectral index, $n = 0$ and -0.5 , is indicated on each panel. The values $n = -1$ and -1.5 are examined in Fig. 2. First and second row of panels: $p_G(v)$ and its logarithm as a function of v . The symbols correspond to the measurements in simulated data as described in the text. Vertical error bars show 1σ errors calculated from 100 independent realizations of the field. As indicated in the top panels, the solid curves correspond to our theoretical prediction (equations 16, 21, 28, 38, 45, 47, 48, 50, 51); the short-dashed ones are the same but assume that peaks are unclustered (Poisson limit, or $N_c = 0$, equivalently $\sigma = 1$ in equation 21, but still use equations 38 and 45 to determine the peak abundance); the dotted ones further assume that the number density of peaks in the excursion is approximated by the Euler characteristic (equation 55); the dot-dashed curves give form (1) fitted on the dotted curves, with the value of γ_G obtained from matching the Taylor expansion discussed in Section 3.4 (equations 59, 65, 67, 69); finally, the last curves (three dots–one dash repeated) correspond to the asymptotic behaviour (3) expected when the ratio L/l approaches infinity: they are the same as the dot-dashed curves but with $\gamma_G = 0$. Third row of panels: N_c (solid curves) and $\xi_2^p \equiv v^2 \bar{\xi}_2(L)/\sigma_0^2$ (dotted curves) as functions of v . When $N_c \gtrsim 1$, one expects the effect of the clustering of peaks to become significant. On the other hand, when $\xi_2^p \gtrsim 1$, our description of the N -point correlation functions of peaks becomes inaccurate, but this only has a significant impact on the analytical calculation of $p_G(v)$ if $N_c \gtrsim 1$. Note that the intersection of the solid and the dotted curves is expected to be in the vicinity of the maximum of $p_G(v)$.

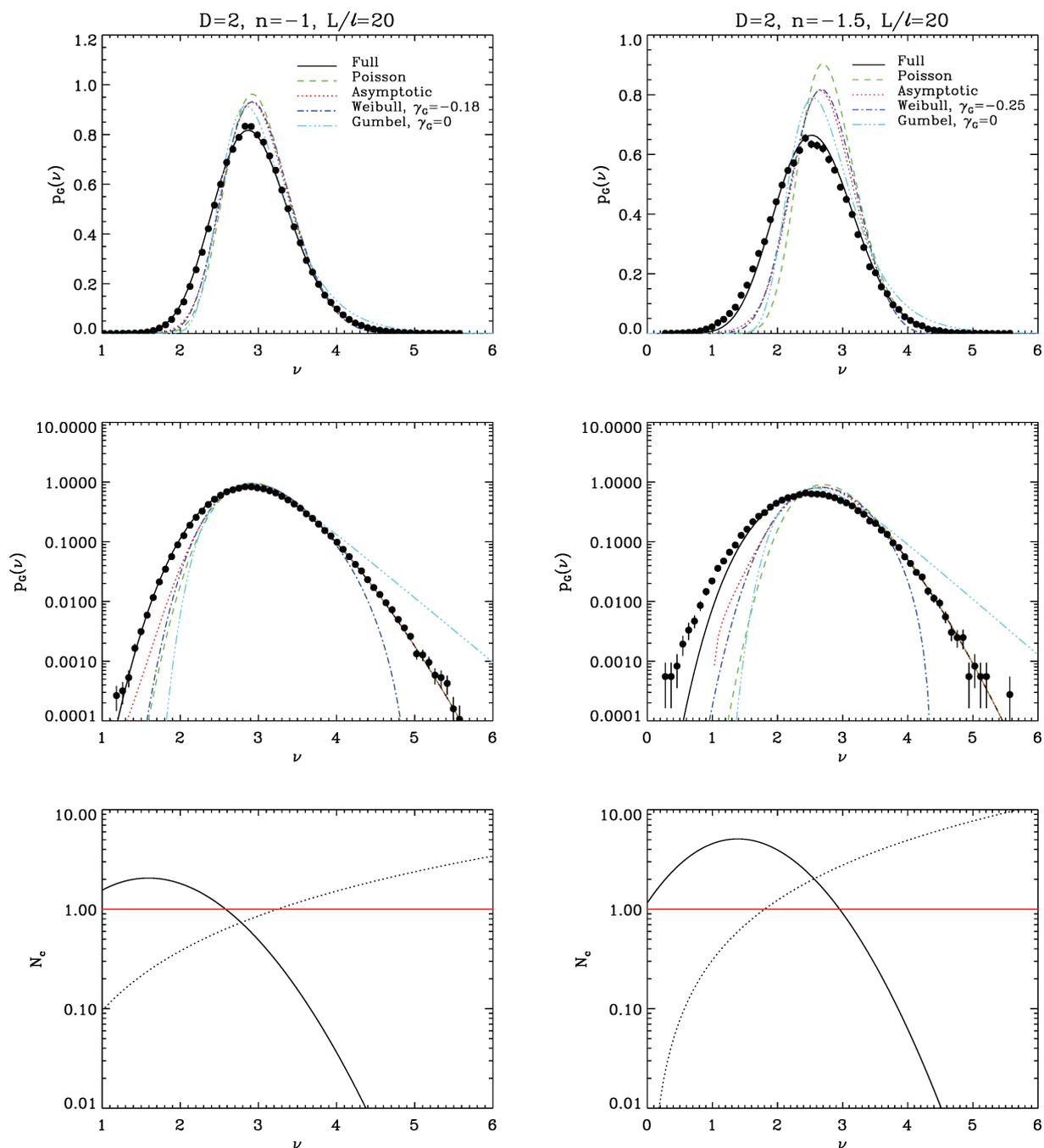


Figure 2. Same as Fig. 1, but for $n = -1$ and -1.5 .

theory is spectacular – with the best results in the high- ν region, as expected. Even the case $n = -1.5$ is well described by theoretical predictions despite the fact that condition (25) is broken while $N_c \gtrsim 1$. Except for $n = -1$, the low- ν tail is slightly off the theory, which overestimates a bit the measurements for $n > -1$ and significantly underestimates them for $n = -1.5$. Still, our analytic calculations are sufficiently accurate to define confidence regions with an accuracy level of a few per cent. As expected, the effect of clustering between peaks decreases with increasing n and becomes rather small for $n > -1$, given the choice of $L/\ell = 20$. In that regime, the number density of peaks is in fact well approximated by the Euler characteristics, and the function $p_G(\nu)$ is well fitted by a negative Weibull-type form

(equation 1) with the parameters derived in Section 3.4, except in the high- ν tail, as expected. In agreement with the predictions of Section 3.4, $|\gamma_G|$ decreases with n and the asymptotic regime (3) is approached slowly although not reached yet, especially in the tails. It was indeed argued in Section 3.4 that convergence to it is rather slow and requires an increasingly large value of L/ℓ as $-n$ becomes larger. Note that the data points or the solid curves can also be fitted easily by form (1) with appropriate choice of a , b and γ_G (Mikelson et al. 2009), except of course for the high- ν tail. For simplicity we did not show that fit because it is purely phenomenological and there is no simple analytic expression for a , b and γ_G except in the asymptotic regime studied in Section 3.4.

5 SUMMARY AND DISCUSSION

We have computed analytically the Gumbel statistics for random Gaussian fields smoothed with a Gaussian window of size ℓ . The Gumbel statistics, $p_G(v_{\max})dv_{\max}$, represents the probability distribution of the maximum v_{\max} of the field estimated in a patch of size L thrown at random. Our important results can be summarized as follows.

(i) For L sufficiently large in front of ℓ , v_{\max} can be approximated by the maximum value of the density estimated at the positions of the peaks included in the patch. As a result, the cumulative Gumbel distribution, $P_G(v) = \int_{-\infty}^v p_G(v_{\max})dv_{\max}$, can be seen as the void probability P_0 of finding no peak with density larger than v in the patch.

(ii) We have made use of the standard count-in-cell formalism (White 1979; Fry 1985; Balian & Schaeffer 1989; Szapudi & Szalay 1993) to compute this void probability as a function of the average number of peaks above the threshold in the patch, $n(v)V$, and their correlation functions averaged over the patch, $\bar{\xi}_N^p(L)$. These quantities were themselves calculated using results of the literature: BBKS and BE to estimate $n(v)$ and $\bar{\xi}_2^p$; BS and Politzer & Wise (1984) to evaluate higher order correlations through a hierarchy of normalized cumulants given by $S_N^p \equiv \bar{\xi}_N^p / (\bar{\xi}_2^p)^{N-1} = N^{N-2}$. Rigorously speaking, these calculations are only valid in the large-separation limit, $\bar{\xi}_2^p \ll 1$ and $\nu \gg 1$. They also neglect contributions from higher order derivatives of the correlation function of the density field, which can in principle be taken into account following Heavens & Sheth (1999) and Desjacques (2008).

(iii) In the regime $\nu \gg 1$ and in the Poisson limit, $N_c \equiv nV\bar{\xi}_2^p \ll 1$, the quantity $-\ln(P_G)$ is simply proportional to the Euler characteristic of the excursion (Aldous 1989). This allows one to derive tractable analytical expressions for the Gumbel statistics (equations 54, 55, 56). We have shown that in this case $P_G(v)$ is well fitted by a negative Weibull-type distribution (1) (with $\gamma_G < 0$), except in the high- ν tail. As expected, $\gamma_G \rightarrow 0$ when $L/\ell \rightarrow \infty$ and one converges slowly to the Gumbel-type distribution, equation (3), as shown long ago (Bickel & Rosenblatt 1973).

(iv) Our analytical calculations were successfully tested against numerical experiments of 2D scale-free Gaussian random fields, in particular in a regime where both $N_c \gtrsim 1$ and $\bar{\xi}_2^p \gtrsim 1$, i.e. where the validity of the ‘exact’ calculations mentioned in point (ii) remains questionable.

Note that our calculations can be easily extended to non-Gaussian fields, using e.g. the formalism of Pogosyan, Gay & Pichon (2009) to estimate the number density of peaks, $n(v)$, and modifications of the Press & Schechter formalism (Press & Schechter 1974) to compute $\bar{\xi}_2^p$ in the high- ν regime (see e.g. Desjacques & Seljak 2010; Valageas 2010, and references therein). The hierarchical relation $S_N^p \simeq N^{N-2}$ should still hold if $\nu \gg 1$ (BS), as extensively discussed at the end of Section 2.2.2 and in the beginning of Section 2.2.3. In the Poisson limit and for $\nu \gg 1$, the result obtained in point (iii) above should still hold: $-\ln(P_G)$ should be simply given by the Euler characteristic, which itself can be easily estimated in the non-Gaussian case (Pogosyan et al. 2009; Matsubara 2010). In fact, one expects that P_G should still be well fitted by the family of distributions (1) but with a different value of γ_G (Mikelson et al. 2009). However, convergence to the asymptotic form (3) in the limit $L/\ell \rightarrow \infty$ remains to be proven. In the intermediary regime probed by the Euler characteristic, the Gumbel statistics provides an interesting test of non-Gaussianity, as shown experimentally by Mikelson et al. (2009) on the simulated temperature maps of the

CMB. In a companion paper (Davis et al. 2011), we have studied applications of the Gumbel statistics to clusters of galaxies, where the quantity of interest is the probability distribution function of the mass of the most massive cluster in the patch (see also Cayón et al. 2010; Holz & Perlmutter 2010). We plan to apply 2D Gumbel statistics to the analysis of CMB data, including non-Gaussian corrections, in the near future.

ACKNOWLEDGMENTS

The authors thank F. Bernardeau for very useful discussions. OD acknowledges the support of an STFC studentship and JD’s research is supported by Adrian Beecroft, the Oxford Martin School and STFC.

REFERENCES

- Adler R. J., 1981, *The Geometry of Random Fields*, Wiley Series in Probability and Mathematical Statistics. Wiley, Chichester
- Adler R. J., Taylor J. E., 2007, *Random Fields and Geometry*, Springer Monographs in Mathematics. Springer, New York
- Aldous D., 1989, *Probability Approximations via the Poisson Clumping Heuristic*, Applied Mathematical Sciences, Vol. 77. Springer-Verlag, New York
- Ayaita Y., Weber M., Wetterich C., 2010, *Phys. Rev. D*, 81, 3507
- Azañs J.-M., Delmas C., 2002, *Extremes*, 5, 181
- Balian R., Schaeffer R., 1989, *A&A*, 220, 1
- Bardeen J. M., Bond J. R., Kaiser N., Szalay A. S., 1986, *ApJ*, 304, 15 (BBKS)
- Bernardeau F., Schaeffer R., 1999, *A&A*, 349, 697 (BS)
- Bhavsar S. P., Barrow J. D., 1985, *MNRAS*, 213, 857
- Bickel P., Rosenblatt M., 1973, in *Krishnaiah P. R., ed., Multivariate Analysis III*. Academic Press, New York, p. 3
- Bond J. R., Efstathiou G., 1987, *MNRAS*, 226, 655 (BE)
- Cayón L., Gordon C., Silk J., 2010, preprint (arXiv:1006.1950)
- Coles P., 1988, *MNRAS*, 231, 125
- Coles S., 2001, *An Introduction to Statistical Modeling of Extreme Values*. Springer, New York
- Colombi S., Pogosyan D., Souradeep T., 2000, *Phys. Rev. Lett.*, 85, 5515
- Cornell C. A., 1968, *Bull. Seismological Soc. America*, 58, 1583
- Cruz M., Martínez-González E., Vielva P., Cayón L., 2005, *MNRAS*, 356, 29
- Davis O., Devriendt J. E. G., Colombi S., Silk J., 2011, *MNRAS*, doi:10.1111/j.1365.2966.2011.18286.x
- Desjacques V., 2008, *Phys. Rev. D*, 78, 103503
- Desjacques V., Seljak U., 2010, *Class. Quantum Grav.*, 27, 124011
- Embrechts P., Schmidli H., 1994, *Math. Methods Operations Res.*, 39, 1
- Fry J. N., 1985, *ApJ*, 289, 10
- Gott J. R., III, Juric M., Schlegel D., Hoyle F., Vogeley M., Tegmark M., Bahcall N., Brinkmann J., 2005, *ApJ*, 624, 463
- Gumbel E. J., 1958, *Statistics of Extremes*. Columbia Univ. Press, New York (reprinted in 2004 by Dover Press, New York)
- Heavens A. F., Sheth R., 1999, *MNRAS*, 310, 1062
- Holz D. E., Perlmutter S., 2010, preprint (arXiv:1004.5349)
- Kaiser N., 1984, *ApJ*, 284, L9
- Katz R. W., Brown B. G., 1992, *Climatic Change*, 21, 289
- Katz R. W., Parlange M. B., Naveau P., 2002, *Advances Water Resources*, 25, 1287
- Larson D. L., Wandelt B. D., 2004, *ApJ*, 613, L85
- Leadbetter M. R., Lindgren G., Rootzén H., 1983, *Extremes and Related Properties of Random Sequences and Processes*. Springer-Verlag, New York
- Lokas E. L., Juszkiewicz R., Bouchet F. R., Hivon E., 1996, *ApJ*, 467, 1

- Matsubara T., 2010, *Phys. Rev. D*, 81, 3505
Mikelsons G., Silk J., Zuntz J., 2009, *MNRAS*, 400, 898
Piterbarg V. I., 1996, *Asymptotic Methods in the Theory of Gaussian Processes and Fields*, *Translations of Mathematical Monographs*, Vol. 148. American Mathematical Society, Providence, RI
Pogosyan D., Gay C., Pichon C., 2009, *Phys. Rev. D*, 80, 1301
Poltzer H. D., Wise M. B., 1984, *ApJ*, 285, L1
Press W. H., Schechter P., 1974, *ApJ*, 187, 425
Rosenblatt M., 1976, *Ann. Probability*, 4, 1009
Szapudi I., Szalay A. S., 1993, *ApJ*, 408, 43
Valageas P., 2010, *A&A*, 514, 46
Vielva P., Martínez-González E., Barreiro R. B., Sanz J. L., Cayón L., 2004, *ApJ*, 609, 22
White S. D. M., 1979, *MNRAS*, 186, 145
Yamila Yaryura C., Baugh C. M., Angulo R. E., 2011, *MNRAS*, doi:10.1111/j.1365-2966.2011.18233.x

This paper has been typeset from a $\text{\TeX}/\text{\LaTeX}$ file prepared by the author.