



HAL
open science

Experimental cosmic statistics - II. Distribution

István Szapudi, Stéphane Colombi, Adrian Jenkins, Jörg Colberg

► **To cite this version:**

István Szapudi, Stéphane Colombi, Adrian Jenkins, Jörg Colberg. Experimental cosmic statistics - II. Distribution. Monthly Notices of the Royal Astronomical Society, 2000, 313, pp.725-733. 10.1046/j.1365-8711.2000.03256.x . hal-04110365

HAL Id: hal-04110365

<https://hal.science/hal-04110365>

Submitted on 3 Jun 2023

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.

Experimental cosmic statistics – II. Distribution

István Szapudi,^{1,2★†} Stéphane Colombi,³ Adrian Jenkins¹ and Jörg Colberg⁴

¹University of Durham, Department of Physics, South Road, Durham DH1 3LE

²Canadian Institute of Theoretical Astrophysics, 60 St George St, Toronto, Ontario M5S 3H8, Canada

³Institut d'Astrophysique de Paris, CNRS, 98bis bd Arago, F-75014 Paris, France

⁴Max-Planck-Institut für Astrophysik, D-85740 Garching, Germany

Accepted 1999 November 17. Received 1999 November 17; in original form 1999 May 24

ABSTRACT

Colombi et al. (Paper I) investigated the counts-in-cells statistics and their respective errors in the τ CDM Virgo Hubble Volume simulation. This extremely large N -body experiment also allows a numerical investigation of the *cosmic distribution function*, $Y(\tilde{A})$, itself for the first time. For a statistic A , $Y(\tilde{A})$ is the probability density of measuring the value \tilde{A} in a finite galaxy catalogue. Y was evaluated for the distribution of counts-in-cells, P_N , the factorial moments, F_k , and the cumulants, $\tilde{\xi}$ and S_N s, using the same subsamples as Paper I.

While Paper I concentrated on the first two moments of Y , i.e. the mean, the cosmic error and the cross-correlations, here the function Y is studied in its full generality, including a preliminary analysis of joint distributions $Y(\tilde{A}, \tilde{B})$. The most significant, and reassuring result for the analyses of future galaxy data is that the cosmic distribution function is nearly Gaussian provided its variance is small. A good practical criterion for the relative cosmic error is that $\Delta A/A \lesssim 0.2$. This means that for accurate measurements, the theory of the cosmic errors, presented by Szapudi & Colombi and Szapudi, Colombi & Bernardeau, and confirmed empirically by Paper I, is sufficient for a full statistical description and thus for a maximum likelihood rating of models. As the cosmic error increases, the cosmic distribution function Y becomes increasingly skewed and is well described by a generalization of the lognormal distribution. The cosmic skewness is introduced as an additional free parameter. The deviation from Gaussianity of $Y(\tilde{F}_k)$ and $Y(\tilde{S}_N)$ increases with order k , N and similarly for $Y(\tilde{P}_N)$ when N is far from the maximum of P_N , or when the scale approaches the size of the catalogue. For our particular experiment, $Y(\tilde{F}_k)$ and $Y(\tilde{\xi})$ are well approximated with the standard lognormal distribution, as evidenced by both the distribution itself and the comparison of the measured skewness with that of the lognormal distribution.

Key words: methods: numerical – methods: statistical – galaxies: clusters: general – large-scale structure of Universe.

1 INTRODUCTION

Precision higher order statistics will become a reality when the new wide field surveys, such as the Sloan Digital Sky Survey (SDSS) and the 2-degree Field Survey (2dF), become available in the near future. These prospective measurements contain information relating to the regime of structure formation, to the nature of initial conditions and to the physics of galaxy formation. The ability of such measurements to constrain models, in a broad sense, is inversely proportional to the overlap between the distribution of statistics predicted by different theories for a finite

galaxy survey. More precisely, maximum likelihood methods give the probability of the particular measurements for each theory, or after inversion, the likelihood of the theories themselves. This is an especially natural and fruitful procedure for a Gaussian distribution, where the first two moments are sufficient for a full statistical description. This simple case is assumed for most analyses in the literature and it motivates the special attention given to the investigation of the errors, or standard deviations. In general, however, the underlying distribution of measurements can be strongly non-Gaussian, in which case the correct shape for the distribution has to be employed for a maximum likelihood analysis. As a consequence, terms such as ‘ 1σ measurement’ lose their usual meaning: a few σ deviation from the average can be quite likely for a non-Gaussian distribution with a long tail. Therefore it is of utmost importance to ask two important questions.

★ E-mail: szapudi@cita.utoronto.ca

† Present address: Canadian Institute of Theoretical Astrophysics, 60 St George St, Toronto, Ontario M5S 3H8, Canada.

(i) In what regime is the Gaussian approximation valid for the distribution of the measured statistical quantities?

(ii) If the Gaussian limit is violated, is there any reasonably simple, practical assumption which would enable a maximum likelihood analysis?

This paper attempts to answer these questions by studying numerically the underlying distribution function of measurements for estimators of higher order statistics based on counts-in-cells. This complements the thorough numerical investigation of the errors undertaken by Colombi et al. (2000, hereafter Paper I, in this issue) and the theoretical investigation of the errors exposed in a suite of papers by Szapudi & Colombi (1996, hereafter SC), Colombi, Szapudi & Szalay (1998, hereafter CSS), and Szapudi, Colombi & Bernardeau (1999, hereafter SCB).

For a particular statistic A , $Y(\tilde{A})$ denotes the probability density of measuring a value \tilde{A} in a finite galaxy catalogue. We consider

the following counts-in-cells statistics: factorial moments F_k , cumulants ξ and S_N , void probability P_0 and its corresponding scaling function $\sigma \equiv -\ln(P_0)/F_1$, as well as the counts-in-cells distribution itself, P_N . A large τ CDM N -body experiment, \mathcal{E} , generated by the VIRGO consortium (e.g. Evrard et al., in preparation) was divided into $C_{\mathcal{E}} = 4096$ cubic subsamples, \mathcal{E}_i , $i = 1, \dots, C_{\mathcal{E}}$ for estimating numerically the cosmic distribution function, $Y(\tilde{A})$. This was rendered possible by the fact that this ‘Hubble Volume’ simulation involves 10^9 particles in a cubic box of size $2000 h^{-1}$ Mpc. A detailed description of the simulation and the method we used to extract counts-in-cells statistics in the full box \mathcal{E} and its each of subsamples \mathcal{E}_i can be found in Paper I.

Paper I concentrated entirely on the first two moments of $Y(\tilde{A})$, the average

$$\langle \tilde{A} \rangle = \int \tilde{A} Y(\tilde{A}) d\tilde{A}, \quad (1)$$

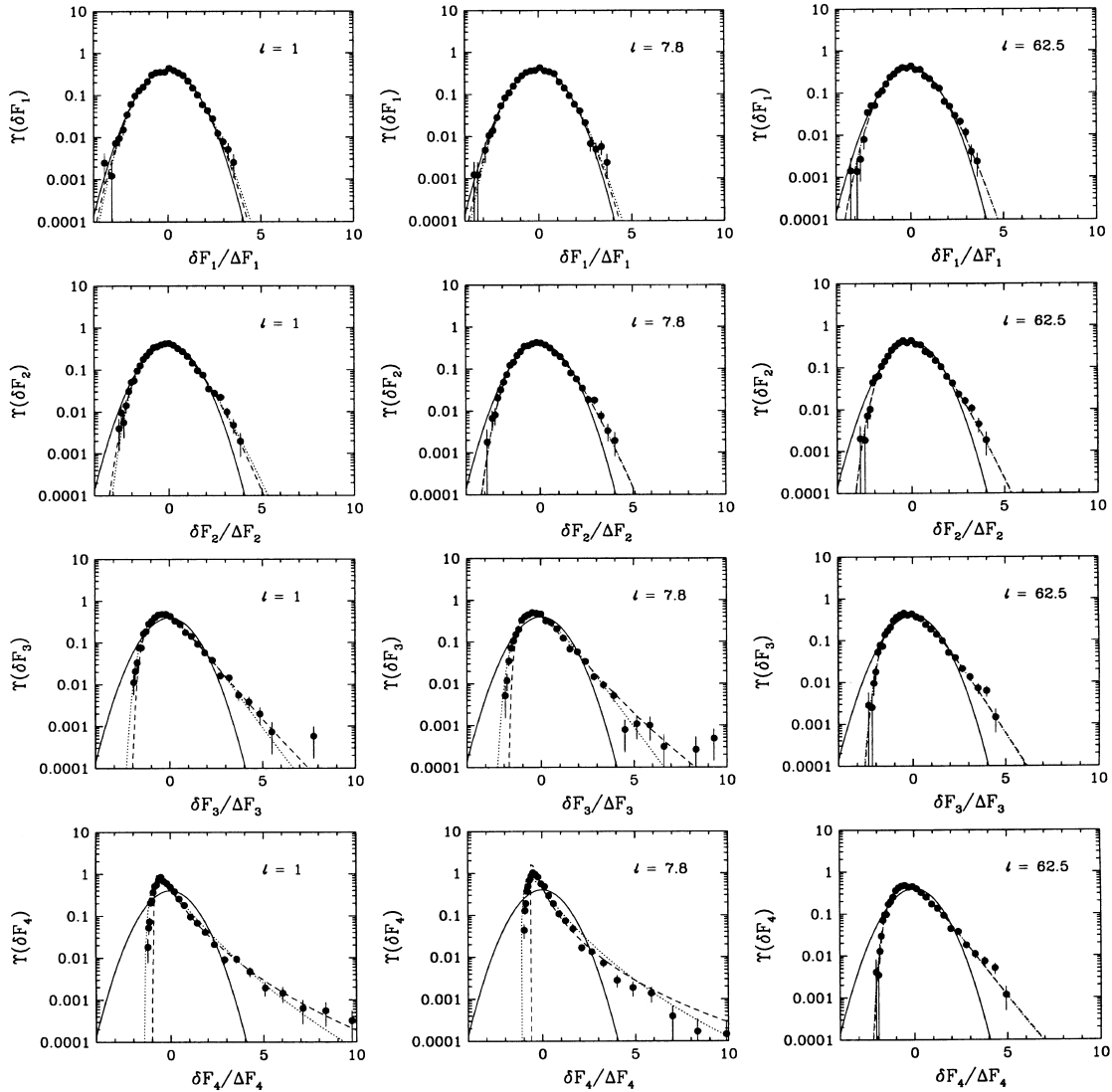


Figure 1. The cosmic distribution function of measurements $Y(\tilde{F}_k)$ shown as a function of $\delta\tilde{F}_k/\Delta F_k$ as explained in the text. The scale of the measurements $\ell = 1, 7.8$ and $62.5 h^{-1}$ Mpc is indicated on each panel. The order $k = 1, 2, 3, 4$ increases from top to bottom. The solid, dotted and dash curves correspond to the Gaussian, lognormal and generalized lognormal (equation 9) distributions, respectively. While the coordinate system of the figure does not display the value of the cosmic error directly, the amount of skewness of the lognormal distribution is an indicator of the magnitude $\Delta F_k/F_k$. The error bars show the measurement error as discussed at the beginning of Section 2.

and the cosmic error

$$(\Delta A)^2 \equiv \langle (\tilde{A} - \langle \tilde{A} \rangle)^2 \rangle = \int (\tilde{A} - \langle \tilde{A} \rangle)^2 Y(\tilde{A}) d\tilde{A}. \quad (2)$$

In the equations above, the mean $\langle \tilde{A} \rangle$ can differ from the true value. The cosmic bias is defined as

$$b_A \equiv \frac{\langle \tilde{A} \rangle}{A} - 1. \quad (3)$$

It is always present when indicators are constructed from unbiased

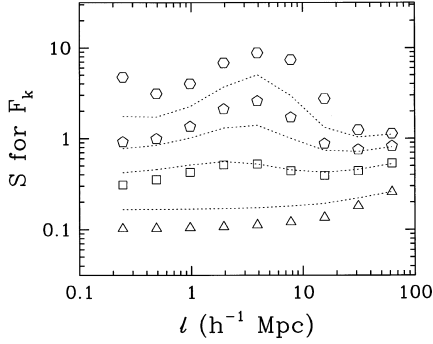


Figure 2. The skewness $S \equiv \langle (\tilde{F}_k - F_k)^3 \rangle / (\Delta F_k)^3$ as a function of scale. The triangles, squares, pentagons and hexagons respectively correspond to $k = 1, 2, 3$ and 4. There are also dotted lines corresponding to an underlying lognormal distribution (8); the orders increase from bottom to top. The errors on the measurement have not been estimated since it would require a complicated calculation depending on the estimate of up to the sixth moment of $Y(\tilde{F}_k)$.

estimators in a non-linear fashion, such as cumulants (e.g. SCB; Hui & Gaztañaga 1998, hereafter HG).

The most relevant results of Paper I are summarized next.

(i) *The measured average $\langle \tilde{A} \rangle$* is in excellent agreement with perturbation theory, one-loop perturbation theory and extended perturbation theory (EPT) in their respective range of applicability. These tests demonstrate the quality of our numerical experiment.

(ii) *The measured cosmic error $\Delta A/A$* is in accord with the theoretical predictions of SC and SCB in their respective domain of validity. A few per cent accuracy is achieved in the weakly non-linear regime for the factorial moments. On small scales the theory tends to overestimate the errors, perhaps by a factor of 2 in the worst case, owing to the approximate nature of the hierarchical models representing the joint moments (SCB).

(iii) *The cosmic bias* is negligible compared to the errors in the full dynamic range, as predicted by theory (SCB, see also HG for an opposing view).

(iv) *Cross-correlations between statistics of order k and l* are in general agreement with theory considering the preliminary nature of the measurements. The precision of the predictions, however, decreases with increasing difference of orders, $|k - l|$. This suggests that the local Poisson model (SC) loses accuracy, as expected.

The theory of the errors confirmed by Paper I provides an excellent basis for future maximum likelihood analyses of data whenever Y is Gaussian. While this was tacitly assumed by most previous works, this article examines for the first time the range of validity of this assumption. To this end the cosmic distribution function $Y(\tilde{A})$ is examined numerically. In particular, one of the

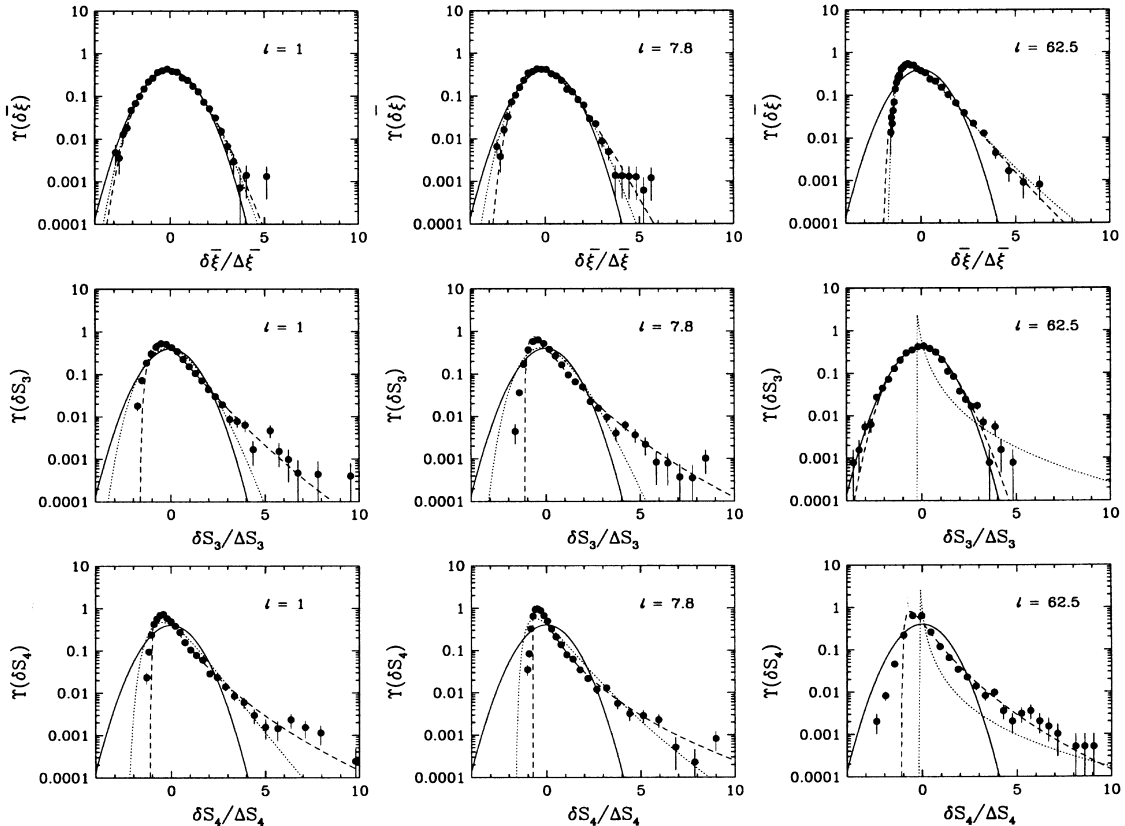


Figure 3. Analogous to Fig. 1 for $Y(\tilde{\xi})$, $Y(\tilde{S}_3)$ and $Y(\tilde{S}_4)$.

parameters determining its shape, the cosmic skewness

$$S \equiv \langle (\bar{A} - \langle \bar{A} \rangle)^3 \rangle / (\Delta A)^3, \quad (4)$$

is calculated as well. When Gaussianity is no longer a good approximation, new Ansätze are proposed for characterizing $Y(\bar{A})$. In addition we perform a preliminary analysis of the bivariate cosmic distributions $Y(\bar{A}, \bar{B})$.

The next section presents the estimates of Y for the factorial moments, the cumulants (including the variance of the counts), the void probability distribution and its scaling function and the counts-in-cells themselves. A universal shape is found for $Y(\bar{A})$ which is well described in all regimes by a generalized version of the lognormal distribution. In addition to the mean (1) and variance (2), this depends on a third parameter, the cosmic skewness (4). This is also investigated along with the resulting *effective cosmic bias*. Section 3 presents the measured bivariate distributions, with explicit comparison to theoretical predictions of SCB.

Finally, section 4 discusses the results in the context of maximum likelihood analysis of future surveys. Readers unfamiliar with counts-in-cells statistics can consult Appendix A in Paper I for a concise summary of definitions and notation.

2 THE COSMIC DISTRIBUTION FUNCTION

The main results of this section are displayed in Figs 1–6. For simplicity Figs 1, 3 and 5 will be referred to as type D, displaying distributions, while Figs 2, 4 and 6 as type S, showing skewness. A general description of each type is followed by the results obtained for the cosmic distribution of the factorial moments (Section 2.1), cumulants (Section 2.2), counts-in-cells (Section 2.3) and void probability with its scaling function σ (Section 2.4). The cosmic skewness and the resulting effective bias are discussed in Section 2.5.

In all figures of type D, the results are displayed in a convenient system of coordinates. For any statistic \bar{A} the normalized

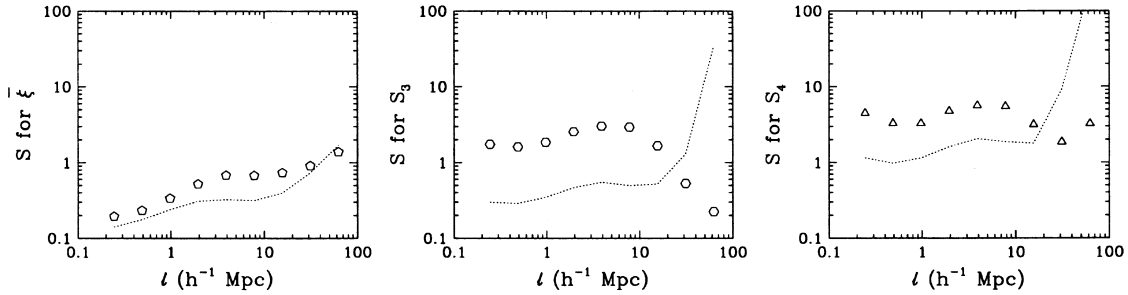


Figure 4. Same as in Fig. 1, but we consider here the skewness of $\xi^{\tilde{}}$ (left panel), \tilde{S}_3 (middle panel) and \tilde{S}_4 (right panel).

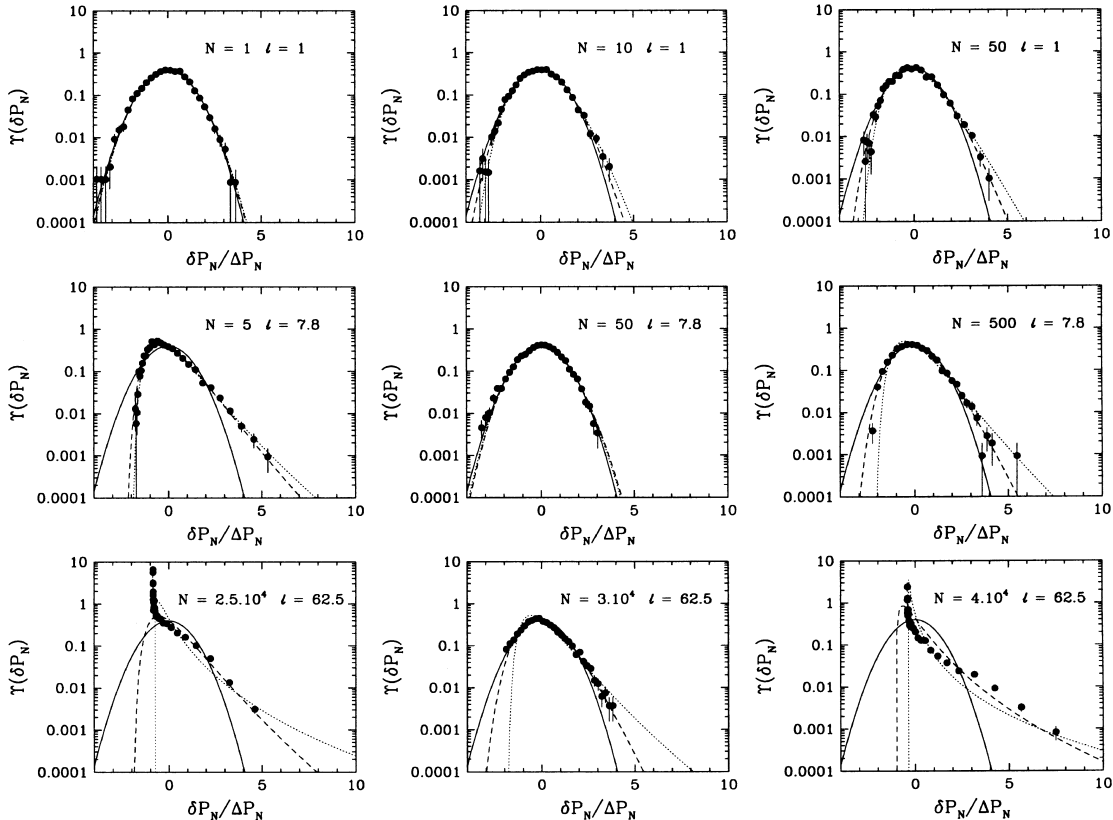


Figure 5. Same as in Fig. 1, but now the distribution function of measurements $Y(\bar{P}_N)$ is shown as a function of $\delta \bar{P}_N / \Delta P_N$ for various scales and values of N as indicated on each panel.

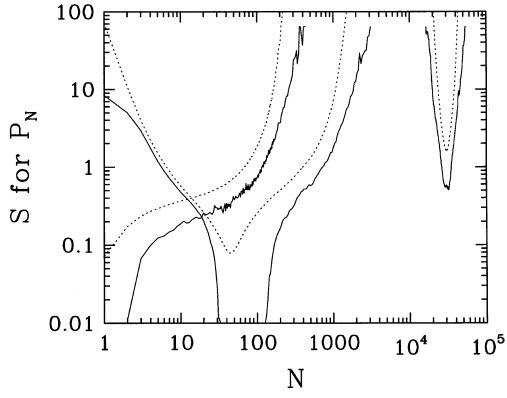


Figure 6. The skewness S of $Y(\tilde{P}_N)$ as a function of N for three scales, $\ell = 1 h^{-1} \text{Mpc}$ (left curve), $\ell = 7.8 h^{-1} \text{Mpc}$ (middle curve) and $\ell = 62.5 h^{-1} \text{Mpc}$ (right curve). The dotted curves give the lognormal prediction, which is always larger than the measurement.

quantity

$$\tilde{x}_A \equiv \frac{\delta \tilde{A}}{\Delta A} = \frac{\tilde{A} - A}{\Delta A} \quad (5)$$

is considered, where $A = \langle \tilde{A} \rangle$ to simplify notations. The average of \tilde{x}_A is zero and its variance is unity by definition, which facilitates the comparison of the plots. The disadvantage of this coordinate system is that the cosmic error $\Delta A/A$ is not directly shown.

For reference, each figure of type D displays a Gaussian (solid curve) and lognormal distribution with the same variance and average (dots, e.g. Coles & Jones 1991):

$$Y(\tilde{A}) = \frac{1}{\tilde{A} \sqrt{2\pi\kappa}} \exp\left\{-\frac{[\ln(\tilde{A}/A) + \kappa/2]^2}{2\kappa}\right\}, \quad (6)$$

with

$$\kappa = \ln[1 + (\Delta A/A)^2]. \quad (7)$$

The skewness of this distribution is given by

$$S = (\Delta A/A)^3 + 3\Delta A/A. \quad (8)$$

For comparison, the skewness of the lognormal assumption is plotted with dotted lines on figures of type S. The amount of skewness of the lognormal is a function of the cosmic error, i.e. more skewness on the figures indicates a larger cosmic error which is hidden by the choice of the coordinate system.

In addition, a ‘generalized lognormal distribution’ is introduced (dashes on figures of type D):

$$Y(\tilde{A}) = \frac{s}{\Delta A [s(\tilde{A} - A)/\Delta A + 1] \sqrt{2\pi\eta}} \times \exp\left(-\frac{\{\ln[s(\tilde{A} - A)/\Delta A + 1] + \eta/2\}^2}{2\eta}\right), \quad (9)$$

$$\eta = \ln(1 + s^2), \quad (10)$$

where s is an adjustable parameter. It is fixed by the requirement that the analytical function (9) has identical average, variance and skewness, $S = s^3 + 3s$, with the measured $Y(\tilde{A})$. It has more parameters, thus form (9) characterizes the shape of function $Y(\tilde{A})$ better than the other two functions, especially for the large $\delta \tilde{A}$ tail. As will be shown next, it is an excellent approximation for the

underlying probability distribution *in all regimes for all statistics*. This robust universality is the most striking result of this article.

The cosmic distribution function, as with any measurement from finite data, is subject to both measurement and cosmic errors (the ‘error on the error problem’, cf. SC). The measurement error on Y , owing to the finite number of subsamples extracted from the whole simulation, can be calculated via straightforward error propagation. It essentially corresponds to the usual $1/\sqrt{C_{\mathcal{E}}}$ factor, where $C_{\mathcal{E}}$ is the number of subsamples. This is plotted on all figures of type D as error bars. On figures of type S no error bars are shown, since this would require an accurate estimate up to the sixth moment of the cosmic distribution $Y(\tilde{A})$. The excellent agreement between cosmic error measurements and theory (Paper I) indicates that the number of subsamples is sufficient and thus the resulting error bars should be fairly small. Similar arguments suggest that the simulation volume was sufficiently large to render the cosmic error on the cosmic distribution negligible.

2.1 Factorial moments

Fig. 1 displays $Y(\tilde{F}_k)$ for $1 \leq k \leq 4$ and various scales $\ell = 1, 7.8, 62.5 h^{-1} \text{Mpc}$.

The agreement with the generalized lognormal distribution is excellent, but even the lognormal gives an adequate description. The deviation from a Gaussian is pronounced whenever the relative cosmic error $\Delta F_k/F_k$ is significantly larger than unity. While the figures do not show the cosmic error directly, the skewness of $Y(\tilde{F}_k)$ is a reliable indication. It increases with the order k since $\Delta F_k/F_k$ also increases with k . Fig. 2 shows directly the quantity S measured for $Y(\tilde{F}_k)$ along with the lognormal value (8). The agreement shows that the lognormal model yields an excellent approximation.

Fig. 1 in conjunction with the measurements of the cosmic error in Paper I suggests that

$$\Delta A/A \leq \Delta_{\text{crit}}, \quad \Delta_{\text{crit}} = 0.2, \quad (11)$$

is a practical criterion for the validity of the Gaussian approximation.

2.2 Cumulants

Fig. 3 is analogous to Fig. 1, showing functions $Y(\tilde{\xi})$, $Y(\tilde{S}_3)$ and $Y(\tilde{S}_4)$ for the biased estimators. As was shown in Paper I, the bias is negligible compared to the cosmic errors, thus correction is not necessary. The agreement with the lognormal is more approximate than for $Y(\tilde{F}_k)$, except for the variance $\tilde{\xi}^2$. Indeed, the skewness of $Y(\tilde{S}_N)$ is in general different from the lognormal prediction, as illustrated by Fig. 4. On small scales it is larger than predicted by equation (8) while on large scales, where edge effects dominate, it is much smaller. The generalized lognormal (9) can still account for the shape of $Y(\tilde{S}_N)$ quite well, especially for the large \tilde{S}_N tail.

The cosmic skewness of $Y(\tilde{S}_k)$ is fairly small on large scales. This is a natural consequence of the fact that cumulants are not subject to the positivity constraint $\tilde{S}_k \geq 0$, as is the case for factorial moments. On large scales, the measured \tilde{S}_k may well be positive or negative and similarly with $\tilde{\xi}$ on extremely large scales. As a result, the left-hand tail of the distribution is more pronounced in both lower right panels of Fig. 3 than the corresponding figure for factorial moments and $Y(S_3)$ is almost Gaussian in the middle-right panel.

Rule (11) for the Gaussian limit still applies, at least for $\tilde{\xi}$, and

perhaps a slightly more stringent condition should be chosen for cumulants of higher order. $Y(\tilde{S}_3)$ is fairly skewed even though the measured cosmic error is slightly below the threshold value for $\ell = 1 h^{-1}$ Mpc and $\ell = 7 h^{-1}$ Mpc (see Paper I).

2.3 Counts-in-cells

Fig. 5 shows the function $Y(\tilde{P}_N)$ in various cases. The upper panels focus on a small scale $\ell \approx 1 h^{-1}$ Mpc. In this regime, the CPDF and $-\Delta P_N/P_N$ are decreasing functions of N as demonstrated in Paper I. Once again, the validity of the Gaussian approximation depends on the size of the cosmic error. As a result, $Y(\tilde{P}_N)$ is nearly Gaussian for $N = 1$ and becomes more and more skewed as N increases. The lognormal approximation appears to be adequate within the errors, although it is slightly too skewed as illustrated by Fig. 6.

The middle panels show an intermediate scale $\ell \approx 7.8 h^{-1}$ Mpc. On these scales (cf. Paper I) both the CPDF and the cosmic error have a unimodal behaviour with an extremum (maximum for the CPDF and correspondingly minimum for the errors) for $N \sim N_{\max} = 26$. This explains why for the chosen values of $N = 5, 50$ and 500 , function $Y(\tilde{P}_N)$ is skewed, approximately Gaussian and skewed again respectively. For $N = 5$ lognormal is an excellent approximation, while the skewness for $N = 500$ is somewhat less than that of a lognormal.

Finally, the lower panels display the largest available scale $\ell = 62.5 h^{-1}$ Mpc. The behaviour of P_N and $\Delta P_N/P_N$ is similar as previously, with the extremum shifted to $N \sim N_{\max} \approx 30\,000$. In this case, the cosmic error is always large, at least of order 50 per cent (cf. Paper I). All the curves are thus significantly skewed for the chosen values of $N = 25\,000, 30\,000$ and $40\,000$. The agreement with the lognormal assumption is somewhat inaccurate, although the generalized lognormal improves the fit, especially for the left-hand panel. Note that the apparently abrupt limit for small values of $\delta\tilde{P}_N/\Delta P_N$ is the result of the positivity constraint $\tilde{P}_N \geq 0$. This constraint becomes quite severe when the average value is much smaller than the errors. While there is still plenty of dynamic range for upscattering, there is a hard restriction for down scattering. This is only partly taken into account in our generalized lognormal model and any modifications in this respect are left for future work. Finally, the practical criterion (11) is again valid for determining the Gaussian approximation.

Note that the finite number $C = 512^3$ of sampling cells (see Paper I), the CPDF is necessarily a multiple of $1/C$. This quantization could cause contamination of $Y(\tilde{P}_N)$ unless $P_N \gg 1/C \approx 10^{-8.13}$. The condition $P_N \geq 10^{-6}$ adopted corresponds to at least ~ 100 cells per subsample on average with N particles. Despite that, a small amount of contamination might still persist for $\delta\tilde{P}_N \gtrsim -P_N$, i.e. at the left side of the plots on Fig. 5. The same effect might also alter the tail of the counts-in-cells measurements presented in Paper I, although not significantly.

2.4 Void probability and scaling function

According to the investigations in Paper I, the cosmic error on P_0 and σ increases steadily with scale up to a sudden transition on scales $\ell \sim 5 h^{-1}$ Mpc, where it becomes large or infinite. This behaviour was studied extensively by Colombi, Bouchet & Schaeffer (1995, hereafter CBS) where more of the details can be found. The most relevant consequence here is that in the available dynamic range the cosmic error is small and $Y(\tilde{P}_0)$ and $Y(\tilde{\sigma})$ are nearly Gaussian. For this reason it would be superfluous to print the corresponding figures.

2.5 Cosmic skewness and cosmic bias

According to Figs 1–6, the degree of skewness of the cosmic distribution function increases with the order k and with $|N - N_{\max}|$, where N_{\max} is the value for which P_N reaches its maximum. The cosmic skewness is already significant for third-order statistics, F_3 and S_3 . An important consequence of the large cosmic skewness is that the maximum $Y(\tilde{A})$, i.e. the most likely measurement, is shifted to the left from the ensemble average on Figs 1, 3 and 5. Maximizing the Ansatz (9), which is always a good fit to the cosmic distribution function, yields

$$b_A = A_{\max}/A - 1 = \frac{\Delta A}{A s} \left(\frac{1}{(1+s^2)^{3/2}} - 1 \right), \quad (12)$$

where b_A is the *effective* cosmic bias. Since $s > 0$, it is negative and its absolute value is smaller than the cosmic error,

$$|b_A| \leq 0.66 \frac{\Delta A}{A}. \quad (13)$$

For a lognormal distribution, $s = \Delta A/A$,

$$|b_A| = 1 - [1 + (\Delta A/A)^2]^{-3/2} \leq 1. \quad (14)$$

The effective cosmic bias becomes increasingly significant when the cosmic error is large. Similarly to the cosmic bias (SCB), $b_A \sim -(3/2)(\Delta A/A)^2$ from expanding equation (14) in the small error regime.

The phenomenon of effective bias was already pointed out by SC (and preliminarily investigated by Colombi, Bouchet & Schaeffer 1994). Since A_{\max} is the most likely value of \tilde{A} , the only one available measurement in a catalogue of the neighbouring universe is likely to a yield lower than average value. This is true even for an unbiased indicator such as \tilde{F}_k or \tilde{P}_N . Unfortunately, this effect cannot be corrected for, but it can be taken into account in the framework of the maximum likelihood approach using the above results on the shape of $Y(\tilde{A})$.

3 BIVARIATE COSMIC DISTRIBUTION FUNCTION: A PRELIMINARY ANALYSIS

Figs 7 and 8 display contours of the joint cosmic distribution $Y(\tilde{A}, \tilde{B})$ (solid lines) for factorial moments and cumulants, respectively. For comparison the Gaussian limit is shown,

$$Y(\tilde{A}, \tilde{B}) = \frac{1}{2\pi\Delta A\Delta B\sqrt{1-\rho^2}} \exp\left[-\frac{1}{2}\mathcal{Q}(\tilde{A}, \tilde{B})\right], \quad (15)$$

$$\mathcal{Q}(\tilde{A}, \tilde{B}) = \frac{1}{1-\rho^2} [\tilde{x}_A^2 - 2\rho\tilde{x}_A\tilde{x}_B + \tilde{x}_B^2], \quad (16)$$

where $\rho \equiv \langle \delta\tilde{x}_A \delta\tilde{x}_B \rangle$ is the cross-correlation coefficient. Dot-dashes display the above function with the measured ρ , $\Delta A/A$ and $\Delta B/B$, while long dashes represent the same function but with the parameters inferred from the theory of SCB with the E^2PT model (see Paper I for details). The contours, corresponding in the Gaussian limit to the 1σ (thin curves) level, $\mathcal{Q}(\tilde{A}, \tilde{B}) = 1$, and the 2σ (thick curves) level, $\mathcal{Q}(\tilde{A}, \tilde{B}) = 4$, are displayed in the coordinate system of the measured \tilde{x}_A and \tilde{x}_B .

On $\ell = 7.1 h^{-1}$ Mpc scales the theoretical predictions are expected to match the second-order moments of Y for factorial moments and even the cross-correlations (see Paper I). This is illustrated by Fig. 7, where the long-dashed ellipses superpose well to the dot-dashed ones. For the cumulants the theory

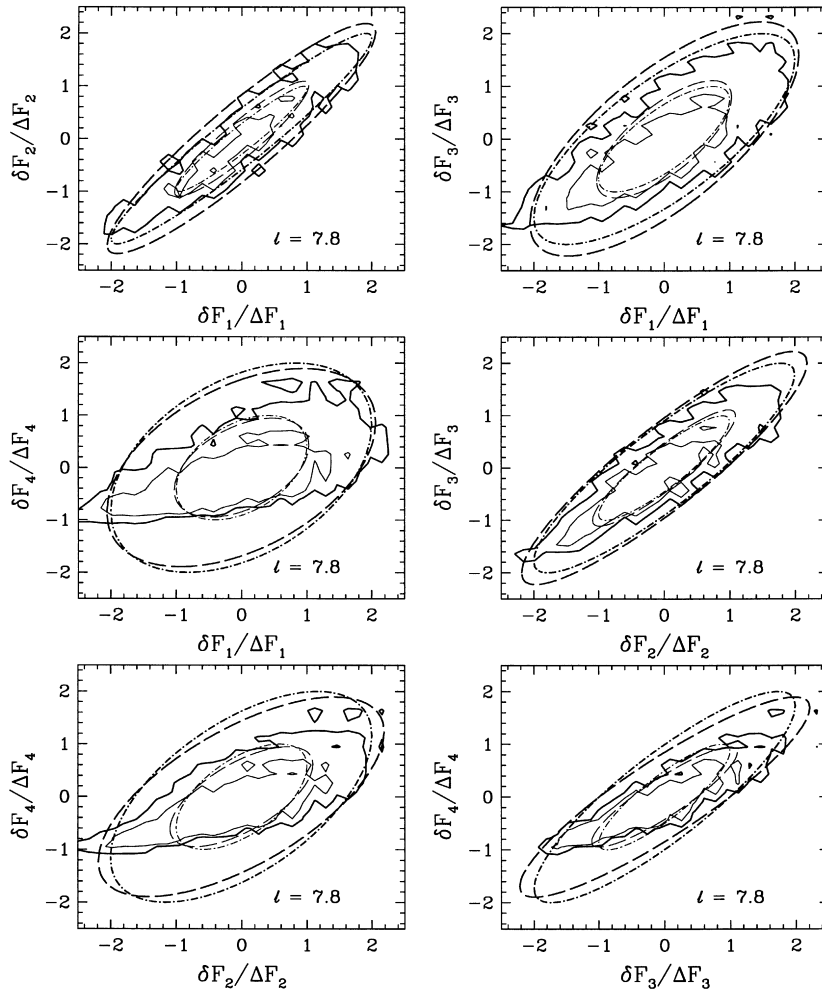


Figure 7. The joint cosmic distribution function for factorial moments, $Y(\tilde{F}_k, \tilde{F}_l)$. Thin and thick solid contours are displayed for two values of Y which would correspond respectively to 1σ and 2σ contours in the Gaussian limit. The latter is shown as thin and thick dot–dashes. For comparison, the analytic prediction of SCB for E^2PT is also plotted with thin and thick long dashes corresponding to the Gaussian limit with theoretical cosmic errors and cross-correlation coefficient. The scale of the measurement is $\ell = 7.8 h^{-1}$ Mpc as displayed on each panel. The image used to draw contour plots has 30^2 pixels. It was generated using bilinear interpolation from another array with logarithmic binning in each coordinate in order to reduce the errors on the estimate of function $Y(\tilde{A}, \tilde{B})$ in each bin.

overestimates the errors slightly, which is reflected in the contours of Fig. 8, although cross-correlations are still reasonable, as indicated by the orientation of the ellipses.

The departure from the Gaussian limit is significant, except for the upper left panel on Figs 7 and 8 and increases with order, in accord with the findings of the previous section. The contrast with Gaussianity increases with the cosmic error and thus with the order considered. With the exception of \tilde{N} , F_2 , $\tilde{\xi}$ and S_3 , the measured cosmic error violates (11) at $\ell = 7.1 h^{-1}$ Mpc (see Paper I). Moreover, as shown previously, criterion (11) should be strengthened for cumulants S_k , $k \geq 3$. In conclusion, condition (11) distinguishes the Gaussian limit for $Y(\tilde{A}, \tilde{B})$ adequately when applied to both statistics \tilde{A} and \tilde{B} .

Similarly to the monivariate distribution (Section 2), function $Y(\tilde{A}, \tilde{B})$ develops skewness and a significant tail for large values of $\tilde{x} = (\tilde{x}_A, \tilde{x}_B)$ when rule (11) is broken. There are three notable consequences.

(i) The effective cosmic bias (Section 2.5) is present again, i.e. the maximum of Y is shifted from the average towards the lower left corner of the panels.

(ii) The contours tend to cover a smaller area than for the Gaussian limit.

(iii) As a result of the positivity constraint, there is a well-defined lower vertical/horizontal bound in some panels, e.g. for $\tilde{x}_{F_4}, F_4 \geq 0$.

4 SUMMARY AND DISCUSSION

This paper has presented an experimental study of the cosmic distribution function of measurements $Y(\tilde{A})$, where \tilde{A} is an indicator of a statistic related to counts-in-cells. The cosmic distribution was considered for the factorial moments F_k , cumulants $\tilde{\xi}$ and S_N , the void probability P_0 with its scaling function, $\sigma \equiv -\ln(P_0)/F_1$, and finally the counts-in-cells P_N themselves. To analyse properties of the function $Y(\tilde{A})$, we used a state-of-the-art τ CDM simulation divided into 4096 subcubes, large enough themselves to represent a full galaxy catalogue. The statistics mentioned above were extracted from each subsample and the resulting distribution of measurements was used to estimate $Y(\tilde{A})$.

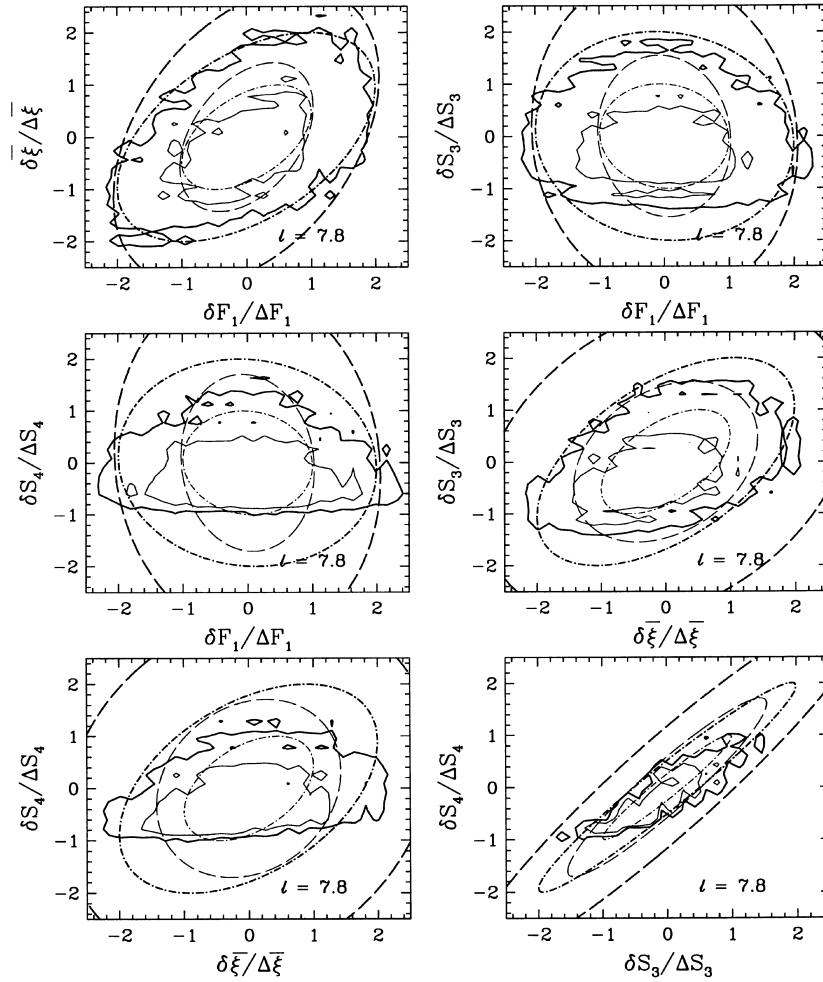


Figure 8. Same as in Fig. 7, but for the average count F_1 and the cumulants, $\bar{\xi}$, S_3 and S_4 .

While Paper I concentrated on the first two moments of the cosmic distribution, the average and the errors, here the focus was shifted towards the general shape of function Y itself, including its skewness, the cosmic skewness. The main results of this analysis are the followings.

(i) In contrast with popular belief, the cosmic distribution is *not Gaussian* in general. The most reassuring result is, however, that the Gaussian approximation appears to be valid whenever the cosmic errors are small, typically $\Delta A/A \lesssim 0.2$. This result is quite robust and it is insensitive to the particular statistic considered (except that a slightly more stringent condition might be chosen for cumulants S_k , $k \geq 3$). This means that for any quantity which can be reliably measured from a survey, a Gaussian error analysis should be valid.

When the relative cosmic error $\Delta A/A$ becomes significant, Y becomes increasingly skewed. Since $\Delta F_k/F_k$ and $\Delta S_k/S_k$ increase with k (SC; Paper I) and $\Delta P_N/P_N$ with $|N - N_{\max}|$, where N_{\max} is the maximum of the CPDF, so does the cosmic skewness, which eventually results in the breakdown of the Gaussian approximation. Functions $Y(\bar{F}_k)$ and $Y(\bar{\xi})$ are well approximated by a lognormal law. Otherwise, a third-order parametrization matching the average, the variance and the skewness of the observed distribution is necessary and in general sufficient. Such a generalization of lognormal distribution is proposed and found

to be in agreement with the measurements in all regimes investigated. Note that there are other alternatives such as the Edgeworth expansion (e.g. Juszkiewicz et al. 1995) or the skewed lognormal approximation of Colombi (1994). This latter consists of applying the Edgeworth expansion to $\log(\hat{A})$. This method, when applicable, improves significantly the domain of validity of the Edgeworth expansion, normally only useful in the weakly non-Gaussian limit $\Delta A/A \lesssim 0.5$.

(ii) While Paper I examined the cosmic bias resulting from the non-linear construction of certain estimators, here a new phenomenon was pointed out, which is similar in effect, but different in nature: the *effective cosmic bias*. It affects all estimators, including unbiased ones, and is a result of the cosmic skewness. Whenever the cosmic errors are large, the cosmic distribution function develops a skewness corresponding to a long tail. As a result, the most likely measurement will be smaller than the average. Such a phenomenon was pointed out earlier in SC, and here it has been found to be universal. As SCB and Paper I found that the cosmic bias is usually insignificant compared to the cosmic errors, it is likely that the effective cosmic bias is responsible for some of the conspicuously low measurements from small galaxy catalogues. This is in contrast with the conjecture of Hui & Gaztañaga (1998, hereafter HG), who assumed that the cosmic bias resulting from the use of biased estimators could explain this phenomenon. The effective cosmic bias renders

correction for the cosmic bias useless, in contrast with the proposition of HG. The effective cosmic bias (and the less significant cosmic bias if any) can be taken into account in the framework of a full maximum likelihood analysis, which relies on the shape of the cosmic distribution function approximated with sufficient accuracy.

(iii) A preliminary investigation of joint distribution $Y(\tilde{A}, \tilde{B})$ was performed for factorial moments and cumulants. It confirms the validity of the above points (i) and (ii) for cosmic bivariate distribution. In particular, a practical criterion for the validity of the Gaussian limit is that the cosmic error for both estimators be small enough, typically $\Delta A/A \lesssim 0.2$ and $\Delta B/B \lesssim 0.2$. This result can be safely generalized to N -variate distribution functions, thus providing the basis of full multivariate maximum likelihood analysis of data in the Gaussian limit.

We have not attempted to develop a more accurate multivariate approximation than (multivariate) Gaussian as this would go beyond the scope of this paper. However, we conjecture that an extension of our generalized lognormal distribution would be feasible (see the point of view of Sheth 1995). An alternate approach, proposed by Amendola (1996), would employ a multivariate Edgeworth expansion. However, similarly with point (i) above for monovariate distributions, this approximation is only valid when the errors are small; but this is precisely the criterion for the Gaussian limit as we have shown previously. A generalization of the lognormal distribution expanding the logarithm of the statistics via the multivariate Edgeworth technique provides a potential improvement of this method.

It is worth noting that the behaviour of the cosmic distribution function is expected to be extremely robust with respect to the particular model studied in this paper, τ CDM. For example, SC, in their preliminary investigations, found essentially the same universal behaviour in Rayleigh–Levy fractals. Moreover, as discussed more extensively in Paper I, the results are sufficiently stable that the usual worries of galaxy biasing (not to be confused with cosmic and effective cosmic bias) and redshift distortions are unlikely to change them qualitatively. Indeed the shape of the cosmic distribution function is almost entirely determined by the magnitude of the cosmic error and it is insensitive to which statistic is considered. The powerful universality found among entirely different statistics is likely to carry over when the two effects mentioned above, which are subtleties in comparison with the range of statistics investigated, are taken into account.

The results found in the present work and in Paper I are encouraging for investigations in future large galaxy catalogues and for problems related to data compression (e.g. Bond 1995; Vogeley & Szalay 1996; Tegmark, Taylor & Heavens 1996; Bond, Jaffe & Knox 1998; Seljak 1998). For example, the cosmic error on factorial moments is expected to be small on a large dynamic range in the SDSS (see, e.g. CSS), implying according to the above findings that the cosmic distribution function should be nearly Gaussian in this regime. In that case, theory of the cosmic

errors and cross-correlations, outlined in SC, CSS and SCB and thoroughly tested in Paper I, will be sufficient for full multivariate maximum likelihood analyses. Preliminary investigations on current surveys are being undertaken by Szapudi, Colombi & Bernardeau (in preparation) and Bouchet, Colombi & Szapudi (in preparation). Similarly the theoretical background is currently being developed for future weak lensing surveys (Bernardeau, Colombi & Szapudi, in preparation), where statistical analyses will be conducted with indicators very close to counts-in-cells (see, e.g. Bernardeau, Van Waerbeke & Mellier 1997; Mellier 1999; Jain, Seljak & White 1999).

ACKNOWLEDGMENTS

We thank F. Bernardeau, P. Fosalba, C. Frenk, R. Scoccimarro, A. Szalay and S. White for useful discussions. It is a pleasure to acknowledge support for visits by IS and SC to the MPA, Garching and by SC to the Department of Physics, Durham, during which part of this work was completed. IS and AJ were supported by the PPARC rolling grant for Extragalactic Astronomy and Cosmology at Durham. The Hubble Volume simulation data was made available by the Virgo Supercomputing Consortium (<http://star-www.dur.ac.uk/~frazierp/virgo/virgo.html>). The simulation was performed on the T3E at the Computing Centre of the Max-Planck Society in Garching. We would like to give our thanks to the many staff at the Rechenzentrum who have helped us to bring this project to fruition.

REFERENCES

- Amendola L., 1996, MNRAS, 283, 983
 Bernardeau F., Van Waerbeke L., Mellier Y., 1997, A&A, 322
 Bond J. R., 1995, Phys. Rev. Lett., 74, 4369
 Bond J. R., Jaffe A. H., Knox L., 1998, Phys. Rev. D, 57, 2117
 Coles P., Jones B., 1991, MNRAS, 248, 1
 Colombi S., 1994, ApJ, 435, 536
 Colombi S., Bouchet F. R., Schaeffer R., 1994, A&A, 281, 301
 Colombi S., Bouchet F. R., Schaeffer R., 1995, ApJS, 96, 401 (CBS)
 Colombi S., Szapudi I., Szalay A. S., 1998, MNRAS, 296, 253 (CSS)
 Colombi S., Szapudi I., Jenkins A., Colberg J., 2000, MNRAS, 313, 711 (Paper I, this issue)
 Hui L., Gaztañaga E., 1999, ApJ, 519, 622 (HG)
 Jain B., Seljak U., White S. D. M., 1999, preprint (astro-ph/9901191)
 Juszkiewicz R., Weinberg D. H., Amsterdamski P., Chodorowski M., Bouchet F. R., 1995, ApJ, 442, 39
 Mellier Y., 1999, ARA&A, 37, 127
 Seljak U., 1998, ApJ, 503, 492
 Sheth R. K., 1995, MNRAS, 277, 933
 Szapudi I., Colombi S., 1996, ApJ, 470, 131 (SC)
 Szapudi I., Colombi S., Bernardeau F., 1999, MNRAS, 310, 428 (SCB)
 Tegmark M., Taylor A. N., Heavens A. F., 1996, ApJ, 480, 22
 Vogeley M. S., Szalay A. S., 1996, ApJ, 465, 34

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