

Constrained correlation functions from the Millennium Simulation

P. Wilking, R. Röseler, and P. Schneider

Argelander-Institut für Astronomie, Universität Bonn, Auf dem Hügel 71, 53121 Bonn, Germany
e-mail: [pwilking;peter]@astro.uni-bonn.de

Received 16 February 2015 / Accepted 13 July 2015

ABSTRACT

Context. In previous work, we developed a quasi-Gaussian approximation for the likelihood of correlation functions that incorporates fundamental mathematical constraints on correlation functions, in contrast to the usual Gaussian approach. The analytical computation of these constraints is only feasible in the case of correlation functions of one-dimensional random fields.

Aims. In this work, we aim to obtain corresponding constraints in the case of higher dimensional random fields and test them in a more realistic context.

Methods. We develop numerical methods of computing the constraints on correlation functions that are also applicable for two- and three-dimensional fields. To test the accuracy of the numerically obtained constraints, we compare them to the analytical results for the one-dimensional case. Finally, we compute correlation functions from the halo catalog of the Millennium Simulation, check whether they obey the constraints, and examine the performance of the transformation used in the construction of the quasi-Gaussian likelihood.

Results. We find that our numerical methods of computing the constraints are robust and that the correlation functions measured from the Millennium Simulation obey them. Even though the measured correlation functions lie well inside the allowed region of parameter space, i.e., far away from the boundaries of the allowed volume defined by the constraints, we find strong indications that the quasi-Gaussian likelihood yields a substantially more accurate description than the Gaussian one.

Key words. methods: statistical – cosmological parameters – large-scale structure of Universe – galaxies: statistics – cosmology: miscellaneous

1. Introduction

The two-point correlation function ξ is been a very common tool in cosmology, although an increasing amount of astronomical literature deals with higher order statistics. Whenever correlation function measurements are used in a Bayesian framework to determine cosmological parameters, the probability distribution function (PDF) of the correlation function is needed. Usually, this likelihood, $\mathcal{L}(\xi)$, is assumed to be a multivariate Gaussian distribution (see, for example, an analysis of the correlation function of the cosmic microwave background by [Seljak & Bertschinger 1993](#); or common methods of baryon acoustic oscillations detection e.g., by [Labatie et al. 2012](#)).

However, the Gaussian approximation of $\mathcal{L}(\xi)$ is not necessarily well justified in all cases and may not always provide the level of precision required from statistical tools that are used to analyze state-of-the-art astronomical data, for example, non-Gaussianities in the cosmic shear likelihood were detected by [Hartlap et al. \(2009\)](#). In the case of third-order cosmic shear statistics, however, [Simon et al. \(2015\)](#) recently have found, at least in current state-of-the-art surveys, that a Gaussian likelihood is a reasonably good approximation. This agrees with results for the bispectrum covariance put forward by [Martin et al. \(2012\)](#). As an additional remark, objections against the use of Gaussian likelihoods as a “safe default” have been raised in cases where one lacks knowledge of the exact form of the likelihood, as pointed out, for example, in power spectrum analyses by [Carron \(2013\)](#) and [Sun et al. \(2013\)](#).

A very strong argument against the Gaussianity of $\mathcal{L}(\xi)$ is the existence of fundamental constraints that stem from the

non-negativity of the power spectrum and was put forward by [Schneider & Hartlap \(2009, hereafter SH2009\)](#). That correlation functions cannot take arbitrary values immediately implies that the Gaussian approximation cannot be fully correct, since a Gaussian distribution has infinite support. To remedy this, one might be tempted to use a Gaussian likelihood for ξ and include the constraints by simply incorporating priors that are zero outside the allowed region. However, as shown in previous work (see Figs. 4 and 5 in SH2009), the shape of the distributions of ξ are strongly affected by the constraints, even well inside the admissible range and thus a more comprehensive solution is needed.

Of course, it would be preferable to obtain the true PDF of ξ analytically, which is feasible only for the uni- and bivariate cases, even assuming one-dimensional Gaussian random fields, as shown by [Keitel & Schneider \(2011\)](#). Their results are a crucial ingredient of the quasi-Gaussian approach introduced in [Wilking & Schneider \(2013, hereafter WS2013\)](#). There, we use the aforementioned constraints to transform the correlation function into an unconstrained variable, where the Gaussian approximation is expected to hold to higher accuracy. Using numerical simulations, we show that for the correlation functions of one-dimensional Gaussian fields, this “quasi-Gaussian transformation” performs very well, meaning that it transforms ξ into a variable that is highly Gaussian. When we make use of the analytical univariate $p(\xi)$ from [Keitel & Schneider \(2011\)](#), this transformation can then be exploited to construct the quasi-Gaussian likelihood for ξ . As presented in WS2013, the new description of $\mathcal{L}(\xi)$ agrees well with the distributions obtained from simulations and

has an impact on the results of Bayesian parameter estimation, as shown in their toy-model analysis.

To date, a major caveat of the quasi-Gaussian approach stems from the fact that the analytical computation of the constraints presented in SH2009 is only optimal for one-dimensional random fields. This severely limits the set of possible applications of the results presented in WS2013. In Sect. 2, we develop numerical methods to compute the constraints on correlation functions that are also applicable to higher dimensional random fields, check their robustness, and compare the numerically obtained constraints to the analytical results for the one-dimensional case. In Sect. 3, we then apply the derived methods in an astrophysical context, i.e., to correlation functions measured from the halo catalog of the Millennium Simulation. We discuss some practical aspects of measuring ξ and show that the correlation functions obtained from the simulation clearly obey the constraints. Furthermore, we examine the performance of the quasi-Gaussian transformation: by comparing the skewness and kurtosis of the transformed and the untransformed correlation functions, we argue that the quasi-Gaussian PDF is a better description of the likelihood of correlation functions than the Gaussian one. We conclude with a brief summary and outlook in Sect. 4.

2. Numerical computation of the constraints on correlation functions

We consider the two-point correlation function of a random field $g(\mathbf{x})$, which is defined as $\xi(\mathbf{x}, \mathbf{y}) = \langle g(\mathbf{x}) g^*(\mathbf{y}) \rangle$. It is related to the power spectrum via Fourier transform. If assuming isotropy, this can be written as

$$\xi(s) = \int \frac{d^n k}{(2\pi)^n} P(|\mathbf{k}|) \exp(i\mathbf{k} \cdot \mathbf{x}) = \int \frac{dk}{2\pi} k^{n-1} P(k) Z_n(ks), \quad (1)$$

where $s \equiv |\mathbf{x}|$, and the dimensionality n of the underlying random field determines the function $Z_n(\eta)$. For a one-dimensional field, Eq. (1) becomes a cosine transform; in the 2D case, $Z_2(\eta) = J_0(\eta)$ is the Bessel function of the first kind of zero order; and for a 3D random field, $Z_3(\eta) = j_0(\eta)$ is the spherical Bessel function of zero order.

As SH2009 show, correlation functions obey fundamental constraints, which arise from the non-negativity of the power spectrum and are best expressed in terms of the correlation coefficients $r_n \equiv \xi(s_n)/\xi(0)$. As it turns out, the constraints can then be written in the form

$$r_{nl}(r_1, r_2, \dots, r_{n-1}) \leq r_n \leq r_{nu}(r_1, r_2, \dots, r_{n-1}), \quad (2)$$

meaning that the upper and lower boundaries on r_n are functions of the r_i with $i < n$.

SH2009 use the fact that the covariance matrix $C_{ij} = \langle g_i g_j^* \rangle = \xi_{|i-j|}$ (where $g_i = g(i \Delta x)$ for a one-dimensional random field evaluated at discrete grid points) has to be positive semi-definite, to explicitly calculate the constraints in the case of homogeneous, isotropic random fields, and show that the constraints they obtain are optimal for a one-dimensional random field, meaning that no stricter bounds exist for a general power spectrum. For higher dimensional fields, the bounds found for the one-dimensional case are still obeyed; however, owing to the isotropy of the field and the multidimensional integration in Eq. (1), tighter constraints hold that have to be computed numerically.

The procedure to obtain the optimal constraints numerically is outlined in SH2009. Rewriting Eq. (1) and applying a quadrature formula for the integral yields

$$r(s) \equiv \xi(s)/\xi(0) = \sum_{j=1}^K V_j Z_n(k_j s), \quad (3)$$

where the coefficients fulfill $0 \leq V_j \leq 1$ and $\sum V_j = 1$. We note that this approximation becomes arbitrarily accurate as $K \rightarrow \infty$.

When measuring correlation coefficients for N different separations s_i , each point $\mathbf{r} = (r_1, r_2, \dots, r_N)$, with $r_i = \xi(s_i)/\xi(0)$, in this N -dimensional space can be written as a weighted sum along the curve $\mathbf{c}(\lambda) = (Z_n(\lambda s_1), \dots, Z_n(\lambda s_N))$, where we used a continuous variable λ with $0 \leq \lambda < \infty$ instead of discrete wave numbers k_j :

$$\mathbf{r} = \sum_{j=1}^K V_j \mathbf{c}(\lambda_j). \quad (4)$$

Since $0 \leq V_j \leq 1$ and $\sum V_j = 1$, each point \mathbf{r} has to lie within the convex envelope of the curve $\mathbf{c}(\lambda)$, which corresponds to the constraints on the correlation coefficients; for example, by constructing the convex envelope of the curve $\mathbf{c}(\lambda)$ for two lags (r_1, r_2) in the one-dimensional case, this reproduces the analytically known bounds $r_{2u,1}(r_1)$.

As a result, to find the constraints only requires describing the convex envelope of the curve $\mathbf{c}(\lambda)$. Unfortunately, there does not seem to be a general analytical solution for this problem, which means resorting to numerical methods; for example, the qhull algorithm (Barber et al. 1996)¹ provides an efficient implementation for computing, among other things, convex hulls. It is, however, limited to inputs of dimensionality lower than nine, meaning that it is only applicable for a maximum number of separations of $N = 8$. Although this is not a requirement for computing the constraints, we use equidistant lags throughout this work, denoting $s_n = n \Delta s$.

As an example of determining the constraints, Fig. 1 shows the curve $\mathbf{c}(\lambda)$ in the $r_1 - r_2$ -plane, plotted in black up to as high as $\lambda = 50$ for illustrative purposes, as well as its convex hull. For a given r_1 , the upper and lower bounds on r_2 are given as intersections with the red hull. This method can, of course, be generalized to higher dimensions, e.g. the determination of $r_{5u,1}$ from r_1, \dots, r_4 , where the convex hull is a hypersurface in a five-dimensional space. Following this procedure, we developed a code to compute the constraints for one-, two-, and three-dimensional fields that we use in our analysis of correlation functions measured from the Millennium Simulation in Sect. 3.

2.1. Comparison of the constraints for one-dimensional fields

Below, we test the numerical method to obtain the constraints. To do so, we compare the numerical results to the analytically computed bounds. Since an analytical calculation of the constraints is only possible for one-dimensional random fields and equidistant lags, we limit our comparison to this case. Throughout this work, we use a gridded approach and denote $\xi_n \equiv \xi(s_n) = \xi(n \Delta s)$, where $\Delta s = L/N$ is the separation between adjacent grid points, and L denotes the length of the field.

There are several ways of testing our methods of computing the constraints: Most straightforward is to compare the constraints from the two methods directly, i.e., to compute the upper and lower bounds r_{nu} and r_{nl} both analytically and numerically, and check how much they differ. An alternative approach

¹ <http://www.qhull.org>

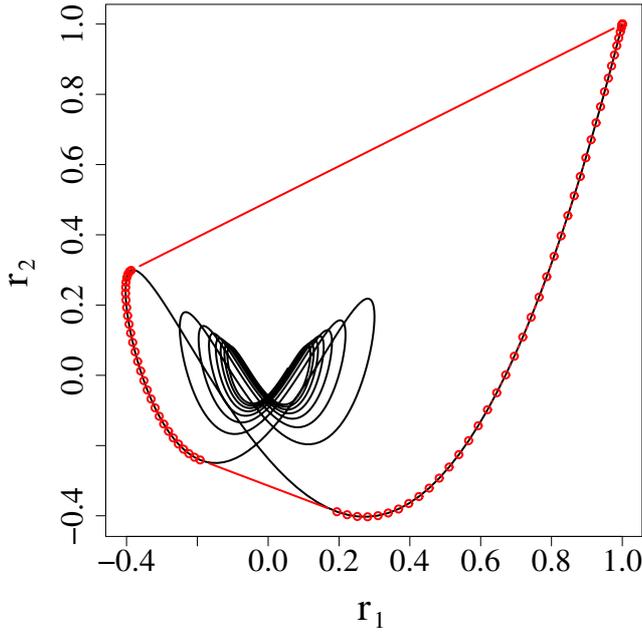


Fig. 1. Example of the curve $c(\lambda)$ for a two-dimensional random field in the $r_1 - r_2$ -plane, where $c(\lambda) = (J_0(\lambda), J_0(2\lambda))$. The red circles and the line connecting them show the convex hulls determined by qhull.

involves the quasi-Gaussian transformation $r_n \rightarrow y_n$, which, as briefly explained in Sect. 1, is a central ingredient of the quasi-Gaussian approximation for the likelihood of the correlation functions introduced in WS2013:

$$y_n = \operatorname{atanh} \frac{2r_n - r_{n_u} - r_{n_l}}{r_{n_u} - r_{n_l}}. \quad (5)$$

Since this transformation is the main application for the constraints, it can – and should be – applied as a means to compare the analytically and numerically obtained bounds, namely by using the different sets of constraints in the transformation and comparing the resulting y_n .

As previously described, the constraints on r_n are functions of the correlation coefficients with lower lags, and as such, we need input values for r_1, \dots, r_{n-1} to be able to compute and compare the different r_{n_l} and r_{n_u} . Again, two possibilities arise: to provide input values that are close to “real-life” applications, we can use realizations of correlation coefficients obtained from numerical simulations (see WS2013 for an efficient way to generate realizations of the correlation function of a one-dimensional Gaussian random field). However, this obviously requires assumptions about the underlying random field and, in particular, its power spectrum. Consequently, a more general approach is to draw the input correlation coefficients for computing the constraints randomly, i.e., from a uniform distribution over the allowed range, $r_n \in]r_{n_l}, r_{n_u}[$. Due to the nature of the constraints, this is an iterative procedure, meaning that one has to draw $r_1 \in]r_{1_l}, r_{1_u}[$, compute $r_{2_l, u}$ from this r_1 , then draw $r_2 \in]r_{2_l}, r_{2_u}[$ to determine $r_{3_l, u}$, and so on.

A comparison of the analytically and numerically obtained bounds r_{n_u} and r_{n_l} is shown in Fig. 2, where we plot the differences $r_{u,l}^{\text{ana}} - r_{u,l}^{\text{num}}$ as functions of n . For each bound $r_{n_u, l}$, the required input values of the correlation coefficients r_i with $i < n$ are drawn uniformly, as previously described. To perform a statistically significant check, this procedure is repeated 500 times, meaning that we generate 500 realizations of the input correlation coefficients and compute the upper and lower bounds both

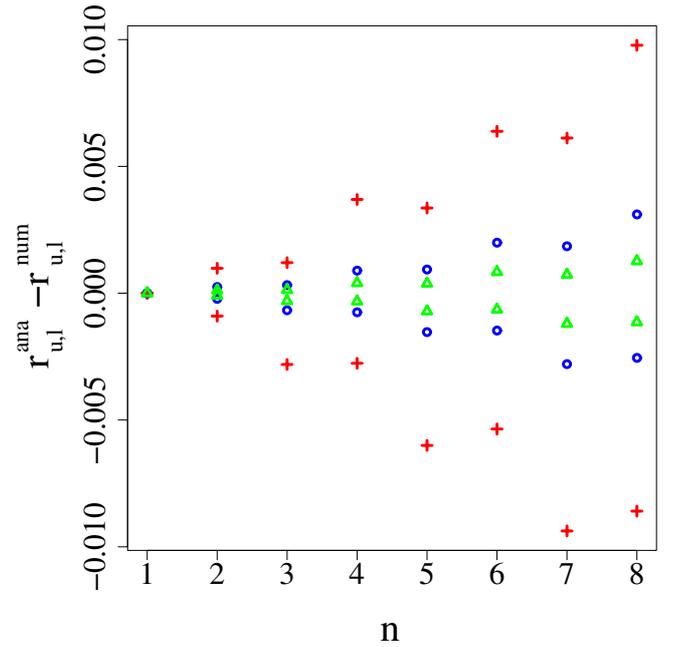


Fig. 2. Difference between the analytically and numerically obtained bounds, averaged over 500 realizations. The upper three sets of points correspond to the difference $r_{u,l}^{\text{ana}} - r_{u,l}^{\text{num}}$, whereas the ones with negative values show $r_{l,l}^{\text{ana}} - r_{l,l}^{\text{num}}$. Furthermore, the different symbols denote the number of steps used to sample the convex hull of the curve $c(\lambda)$ for values of $0 \leq \lambda \leq 2\pi$, namely 100 (red crosses), 200 (blue circles), and 300 (green triangles) steps.

numerically and analytically for each realization. The values plotted in the figure are obtained by averaging the difference between the analytical and the numerical values over the 500 realizations. In addition, we also investigate how much impact the sampling of the convex hull of the curve $c(\lambda)$ has on the accuracy of the numerical bounds. It turns out that it is sufficient in all cases to sample the curve $c(\lambda)$ for values of $0 \leq \lambda \leq 2\pi$, since going to higher values of λ has no impact on the volume within the convex hull because of the periodicity of the function $Z_n(ks)$ in Eq. (1). By way of comparison, we vary the sampling rate of the convex hull; i.e., we sample it using 100 (red crosses), 200 (blue circles), and 300 (green triangles) steps. The upper three sets of points show the difference $r_{u,l}^{\text{ana}} - r_{u,l}^{\text{num}}$ between the analytical and numerical upper bounds, whereas the lower ones depict the deviation between the lower bounds, i.e., $r_{l,l}^{\text{ana}} - r_{l,l}^{\text{num}}$.

Three conclusions can be drawn from Fig. 2: first, the deviation of the numerically obtained bounds from the analytical ones shows a tendency to grow with n , which is to be expected, since the sampling of the convex hull becomes more challenging with increasing dimensionality. (The numerical and analytical bounds on r_1 do not differ at all, since $r_{1_u, l} = \pm 1$.) Second, the impact of this sampling has a strong impact on the accuracy of the numerical calculation of the bounds; namely, the difference between the numerical and the analytical results decreases by about a factor of three when doubling the number of steps used for the convex hull sampling. Actually, this sampling is the limiting factor for the accuracy of the numerical bounds, as can be seen from the third observation: in the case of the upper bounds, the numerical results are systematically smaller than the exact analytical values, whereas for the lower bounds, the numerical values are too high. This effect is an expected consequence of the non-continuous approximation for the smooth hull; as a result of

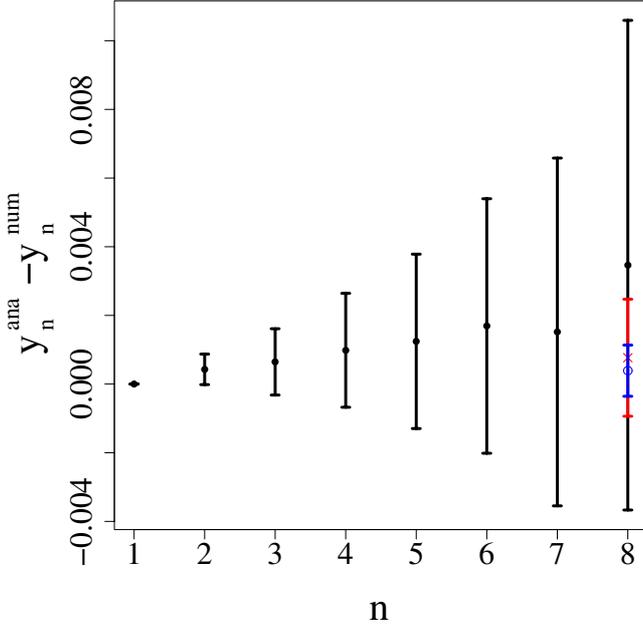


Fig. 3. Difference between the y_n computed using the analytically and the numerically obtained bounds, averaged over 500 realizations, with error bars showing the standard deviations. The input values for the computation of the bounds used in the transformation $r_n \rightarrow y_n$ stem from 500 simulated realizations of the correlation function on a one-dimensional Gaussian field of length L with $N = 32$ grid points and a Gaussian power spectrum of width k_0 , with $Lk_0 = 80$. In the case of the black, solid points, the convex hull of the curve $c(\lambda)$ is sampled using 100 points for the range $0 \leq \lambda \leq 2\pi$. By way of comparison, we show the corresponding results for sampling rates of 200 (red cross) and 300 (blue circle) in the case $n = 8$.

convexity, the hyperplanes that link the points used to sample the hull always have to be located inside the hull.

In summary, the accuracy of the numerical constraints can be increased by improving the sampling of the hull. While a larger number of steps would, presumably, improve the results even further, using more than 300 steps for the sampling becomes impractical because of the computational costs. However, as the following tests demonstrate, using 300 steps seems sufficiently accurate.

As mentioned before, another important check for the accuracy of the numerical methods that are used to compute the bounds is to apply the quasi-Gaussian transformation $r_n \rightarrow y_n$ and to compare the resulting y_n . Below, we adopt correlation coefficients from simulations instead of uniformly drawn ones as input for the computation of the bounds. In particular, we use 500 realizations of the correlation function on a one-dimensional Gaussian field of length L with $N = 32$ grid points and a Gaussian power spectrum, where the field length and the width of the power spectrum are related by $Lk_0 = 80$. (For details on the simulations used in this work, we refer to WS2013.)

For each simulated realization of the correlation coefficients, we compute the bounds r_{nu} and r_{nl} for each n , both numerically and analytically, and use them to transform r_n to y_n , as defined in Eq. (5). To compare the resulting values for y_n , we plot the differences $y_n^{\text{ana}} - y_n^{\text{num}}$ as a function of n in Fig. 3. Here, the plotted values are the average over the 500 realizations, while the error bars denote the standard deviations. For the sake of clarity, we only show the impact of the number of steps used in the convex hull sampling for the case $n = 8$, where the standard deviations are largest.

The accuracy of the numerical approximation again shows a strong dependence on the number of steps used to sample the convex hull. Nevertheless, the difference between the values of y_n , computed using the analytical and the numerical bounds, becomes very small when using 300 steps. As a result, we conclude that the problem of the non-continuous approximation of the curve $c(\lambda)$ and its convex hull can be tackled and that using 300 steps in the sampling yields sufficiently accurate bounds.

3. Application to the Millennium Simulation

So far, our studies about the constraints on correlation functions and the quasi-Gaussian likelihood have been performed in a general, mathematical framework. In this section, we investigate our results in a more astrophysical context by applying them to cosmological correlation functions measured in the Millennium Simulation (Springel et al. 2005). Thus, we aim to check the relevance of the constraints that originally stem from purely mathematical properties, since they are based on the fact that ξ is the Fourier transform of a positive quantity (the power spectrum), in a more practical context, where ξ is measured using an estimator. The size of the Millennium Simulation enables us to easily measure multiple realizations of ξ , thereby providing an approximate determination of the underlying probability distribution and, consequently, a statistical analysis.

3.1. Computing correlation functions

Below, we compute the correlation function of dark matter halos in the Millennium Simulation. Because we are not interested in redshift evolution, we only use the halo catalog from the $z = 0$ simulation snapshot, from which we then select typical galaxy-mass halos by choosing a mass cut $M_{200}^{\text{crit}} > 10^{12} h^{-1} M_{\odot}$, which yields a total number of $\sim 440\,000$ halos. However, to perform a statistical analysis, we require different realizations of the correlation function. For this reason, as a first attempt, we divide the full simulation cube into 1000 subcubes of volume $50^3 (h^{-1} \text{Mpc})^3$ and measure ξ in each of the subcubes.

Along with the halo catalog from the simulation we also need a random catalog, so, for each subcube, we draw halo positions uniformly. We then determine the number of halo pairs for given pair separations in both the data and the random catalog, as well as the cross-correlation. From the count rates $DD(s)$, $RR(s)$, $DR(s)$ (normalized to account for different numbers of halos in both the random catalog and the halo catalog from each subcube) at different pair separations s , we compute ξ using an estimator. While Landy-Szalay (LS) is the most widely used estimator (Landy & Szalay 1993), we also aim to test the impact of the choice of estimator on the constraints, and adapt the following common set of estimators from Vargas-Magaña et al. (2013):

$$\xi_{\text{PH}} = \frac{DD}{RR} - 1, \quad \text{Peebles \& Hauser (1974);}$$

$$\xi_{\text{Hew}} = \frac{DD-DR}{RR}, \quad \text{Hewett (1982);}$$

$$\xi_{\text{DP}} = \frac{DD}{DR} - 1, \quad \text{Davis \& Peebles (1983);}$$

$$\xi_{\text{H}} = \frac{DD \times RR}{DR^2} - 1, \quad \text{Hamilton (1993);}$$

$$\xi_{\text{LS}} = \frac{DD-2DR+RR}{RR}, \quad \text{Landy \& Szalay (1993).}$$

As mentioned in Sect. 2, to calculate and test the constraints, we measure the correlation function at equidistant lags, i.e., $\xi_n \equiv \xi(n \cdot \Delta s)$, where the maximum number of lags is $n = 8$, an effect of the limitations of the numerical computation of the constraints.

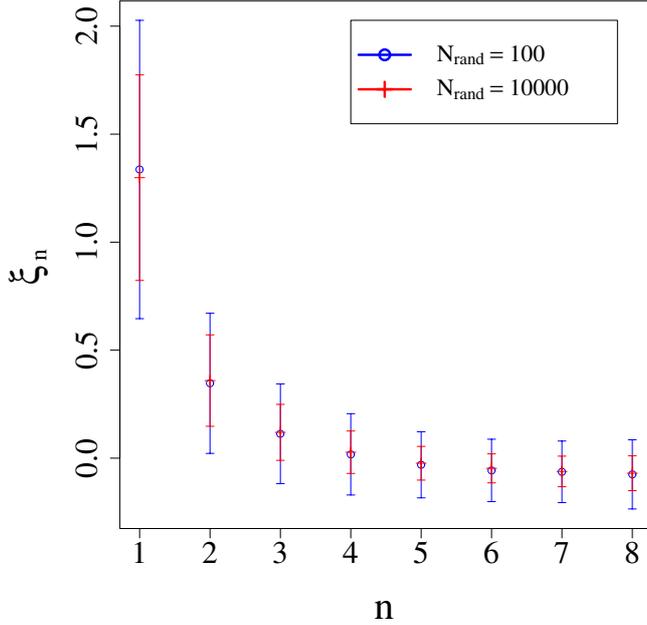


Fig. 4. Correlation function from 1000 subcubes of the Millennium Simulation, computed using the LS estimator. The points and error bars show the correlation function $\xi_n \equiv \xi(n \cdot \Delta s)$ for $\Delta s = 5 h^{-1}$ Mpc (see text for details) averaged over the 1000 subcubes of side length $50 h^{-1}$ Mpc, as well as the standard deviation. For the blue circles, the random catalog for each subcube contains 100 halos, as opposed to 10 000 halos for the red crosses.

The size of the random catalog also merits some discussion: ideally, it should be infinitely large, i.e., $N_{\text{rand}} \rightarrow \infty$. However, the computation of the pair separations is the most time-consuming step in the calculation of ξ and, consequently, the number of halos in the random catalogs for each of the 1000 subcubes is subject to practical limitations. We study the impact of the random catalog size in Fig. 4: here, we show the correlation function for an exemplary choice of lags with $\Delta s = 5 h^{-1}$ Mpc, i.e., we measure $\xi_1, \xi_2, \dots, \xi_8$ at lags of $5, 10, \dots, 40 h^{-1}$ Mpc. In practice, we need to allow for a range of pair separations to obtain sufficiently large numbers of pairs. To do this, we adapt a bin size of width $1 h^{-1}$ Mpc, so, for example, in the computation of ξ_1 we use all pairs with separations ranging from 4.5 to $5.5 h^{-1}$ Mpc. For the auto-correlation ξ_0 , i.e., the correlation function at zero lag (which we do not plot in the figure, but which is required for the calculation of the constraints), we count all pairs with very small separations, e.g., $s \leq 1 h^{-1}$ Mpc. (We refer to the next section for a discussion on the measurement of ξ_0 and on the choice of lags.) The points and error bars show the mean and standard deviation over the 1000 subcubes of side length $50 h^{-1}$ Mpc; we use the LS estimator, which has been shown to be less sensitive to the size of the random catalog than others (see Kerscher et al. 2000). For the blue circles, a small random catalog ($N_{\text{rand}} = 100$ halos for each subcube) was used, whereas the choice of $N_{\text{rand}} = 10\,000$ for the red crosses results in noticeably smaller standard deviations over the 1000 realizations, at the cost of a longer computation time.

Hence, we aim to find a trade-off between those two values. First, although the mean of the correlation functions for the two random catalog sizes plotted in Fig. 4 do not seem to differ very much at first sight, choosing the catalog size as small as $N_{\text{rand}} = 100$ is a quite extreme case. This is because a large fraction of the realizations yield a diverging auto-correlation ξ_0 as a result of the count rates in RR or DR being zero, at least when measuring ξ_0

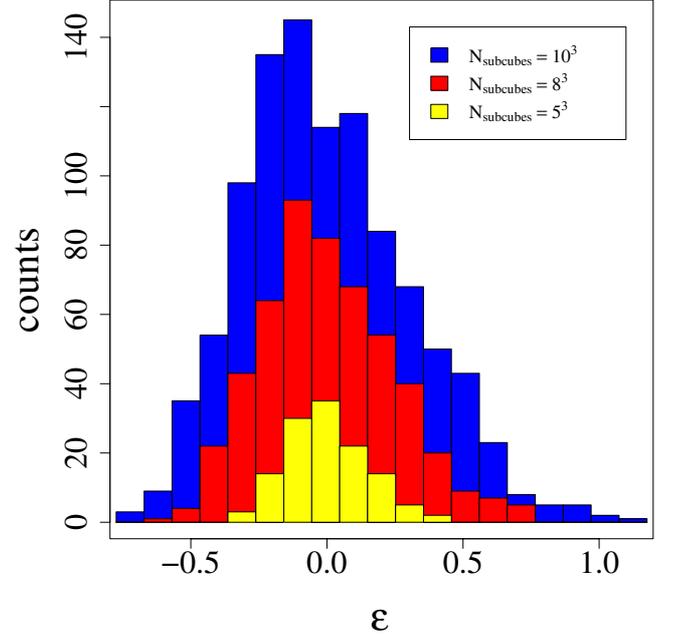


Fig. 5. Histograms of the distributions of overdensities in the subcubes of the simulation volume, where ϵ denotes the number overdensity of halos in the subcube relative to the mean halo number density in the whole simulation box, see Eq. (6). The colors indicate the number of subcubes the simulation box was divided into.

as previously described. However, when increasing the random catalog to 1000 halos per subcube, the mean correlation function for non-zero lag shows a deviation of only $\sim 1\%$ compared to the mean ξ for $N_{\text{rand}} = 10\,000$ (and even here, about a tenth of the realizations show a diverging ξ_0). At $N_{\text{rand}} = 5000$, no such problems occur, while even the error bars, as plotted in the figure, become indistinguishable from those at $N_{\text{rand}} = 10\,000$. This means that a random catalog size of 5000 is a reasonable trade-off between accuracy and computational expenses.

An additional observation from Fig. 4 is that ξ becomes negative for larger lags, i.e., around 20 – $25 h^{-1}$ Mpc. The reason for this is an integral constraint (see, for example, Landy & Szalay 1993) that arises when measuring correlation functions in finite volumes, where the global mean density is unknown and is usually approximated by the mean observed density.

In our case, one way to assess this issue is to decrease the number of subcubes in our analysis, and make them larger and more representative for the whole box, while at the same time measuring ξ at the same lags as before. Beside lessening the impact of the integral constraint on small lags, this has two additional effects: first, with fewer subcubes, we obtain fewer realizations of ξ , which can pose a challenge for a statistical analysis, and second, as the number of halo pairs per cube becomes larger, the scatter over the measured realizations of ξ decreases. To decide on the number of subcubes necessary to make the subcubes as representative as possible for the whole simulation volume, we estimate the overdensity ϵ in each subcube by comparing the mean number density in the subcube to the one from the whole simulation volume, using

$$\bar{n}_{\text{sub}} = \bar{n}_{\text{box}} (1 + \epsilon), \quad (6)$$

and examine the distributions of ϵ by plotting them as histograms in Fig. 5. Here, we slice the simulation volume into different numbers of subcubes and compute the overdensity in each subcube. It is clear that the distribution $p(\epsilon)$ is very broad for the

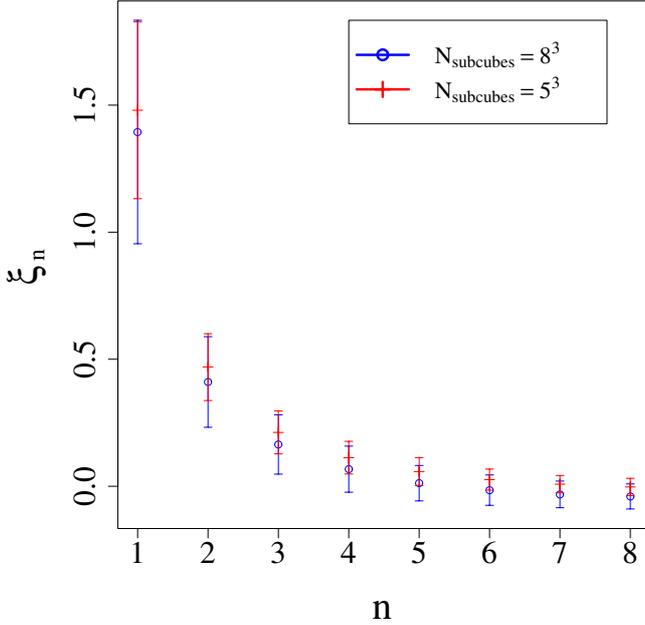


Fig. 6. Correlation function from the Millennium Simulation, computed using the LS estimator, a random catalog size of 5000, and a lag separation of $\Delta s = 5 h^{-1}$ Mpc. The points and error bars show the mean and standard deviation computed over the subcubes of the simulation, where the simulation box was sliced into 8^3 subcubes for the blue data points, as opposed to 5^3 for the red ones.

value $N_{\text{subcubes}} = 10^3$ used so far, and, as expected, it becomes quite narrow for the case of 5^3 subcubes, which indicates that the integral constraint does not pose a large problem in this case. The resulting correlation functions for the cases of 8^3 and 5^3 subcubes (and otherwise the same parameters as before, i.e., same lags and random catalog size) are shown in Fig. 6. As can be seen, slicing the simulation volume into 5^3 subcubes yields reasonable results, i.e., a non-negative correlation functions with a sufficiently small variance.

Finally, we briefly evaluate the choice of estimator. To do so, instead of measuring ξ_n at a few different lags n , it is advisable to compute $\xi(s)$ for all lags in each subcube. Since this is obviously not practicable, we compute it at a high number of lags, meaning that we divide the range of pair separations into adjacent bins of width $0.2 h^{-1}$ Mpc.

The correlation functions (average over the 5^3 subcubes and using a fixed size of random catalog for each subcube, namely $N_{\text{rand}} = 5000$) are shown in Fig. 7. Here, the different colored lines denote the five estimators, and the gray-shaded region depicts the standard deviation over the 125 realizations in the case of the most commonly used LS estimator. For clarity, in the left panel, we plot scales from 8 to $40 h^{-1}$ Mpc, whereas the right panel shows the correlation function for very small lags, i.e., 4 – $8 h^{-1}$ Mpc. Clearly, the numerous estimators yield very similar results, in particular compared to the standard deviation of the 125 realizations.

3.2. Testing the constraints

In this section, we investigate whether the correlation functions computed from the halo catalog of the Millennium Simulation obey the numerically obtained constraints. While we make the case in Sect. 2.1, that using 300 points to sample the hull yields sufficiently good agreement between the numerical and the exact

analytical values of the constraints, we have to restrict ourselves to 270 points in the case of a 3D random field owing to the computational costs. Although the convex hull only has to be computed once and can then be used to determine the constraints for all sets of correlation coefficients, sampling the hull for a 3D random field with the given accuracy poses memory problems for the qhull algorithm, and is beyond the scope of this research. However, this does not pose a problem: when we compare the accuracy of the numerical constraints, as plotted in Fig. 2, it is apparent that the improvement in accuracy, when going from 200 to 300 steps, is far smaller than the one from 100 to 200, and thus we expect 270 points to be accurate enough.

To test the constraints, for each realization we compute the correlation coefficients $r_n \equiv \xi_n/\xi_0$ as well as the upper and lower bounds, r_{nu} and r_{nl} . It turns out that the width of the ξ_0 -bin has a strong influence, in particular on the width of the distributions of the correlation functions r_n . For example, we first choose a relatively broad bin, i.e., we measure ξ_0 by averaging over all pair separations from 0 to $2 h^{-1}$ Mpc. This choice is primarily motivated by the fact that increasing the spread of the correlation coefficients over the 125 realizations allows us to test how close to the edges of the allowed region the r_n move. Toward the end of this section, we study the impact of the width of the ξ_0 -bin in more detail.

One question that arises is how to visualize the constraints. The simplest approach to this is to use scatter-plots, with dots for the individual realizations. An example in the r_1 – r_2 -plane is shown in Fig. 8. Here, the red dots show the different realizations of r_1 and r_2 , computed for the subcubes using the LS estimator. Additionally, we plot iso-density contours that contain 68, 95, and 99.7% of these realizations. For the lefthand panel, we sliced the simulation volume into 1000 subcubes, as opposed to 125 for the righthand panel. As explained in the previous section, the higher number of subcubes greatly increases the spread of the correlation functions, which can also be clearly observed in r -space. (Even for the high number of subcubes, the integral constraint is expected to be negligible for the correlation functions at small lags.) In both panels, the upper and lower blue lines represent the constraints, i.e., $r_{2u,1}(r_1)$, which we compute numerically for each realization of r_1 shown in the figure and plot as connected lines. All realizations clearly lie well inside the constraints, particularly when compared to results for purely Gaussian (one-dimensional) random fields with fiducial power spectra (see, for example, similar figures in related works, such as Fig. 3 from WS2013).

As an additional way of depicting the constraints, we apply a part of the quasi-Gaussian transformation from Eq. (5) to map the allowed range of the correlation coefficients to $(-1, +1)$, namely by transforming the correlation coefficients r_n to

$$x_n = \frac{2r_n - r_{nu} - r_{nl}}{r_{nu} - r_{nl}}. \quad (7)$$

To visualize the constraints more clearly, we use a modified version of box-and-whisker plots, meaning that we display our samples $\{r_n\}$ and $\{x_n\}$ as boxes, where the upper and lower borders show the first and third quartiles of the sample, i.e., the values that split off the upper and lower 25% of the data. As is common practice, we also show the sample median (i.e., the second quartile) as a line inside the box, as well as two whiskers. Usually the ticks at the end of the whiskers denote the minimum and maximum of the data (in the most widely used type of box-and-whisker plot). Here, we use them to denote the upper and lower constraints: Since $r_{nu,i}$ are functions of all r_i with $i < n$, we show the mean r_{nu} and r_{nl} over all realizations for plots in r -space. For

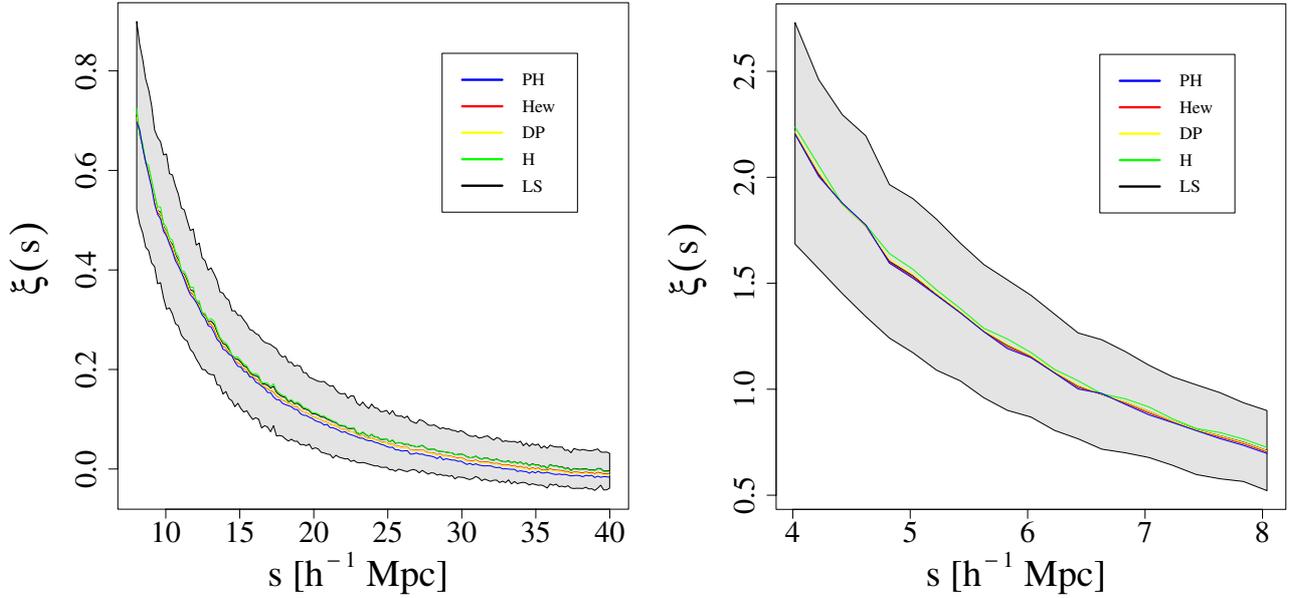


Fig. 7. Correlation function measured for “all” lags (see text for details) as a function of the pair separation s , with a random catalog size of $N_{\text{rand}} = 5000$. The lines denote the mean $\xi(s)$ measured in the 125 subcubes using the different estimators listed in Sect. 3.1, and the gray shaded region shows the standard deviation. In the *left panel*, the s -range from 8 to 40 h^{-1} Mpc is plotted, and *the right panel* shows the results for very small lags, from 4 to 8 h^{-1} Mpc.

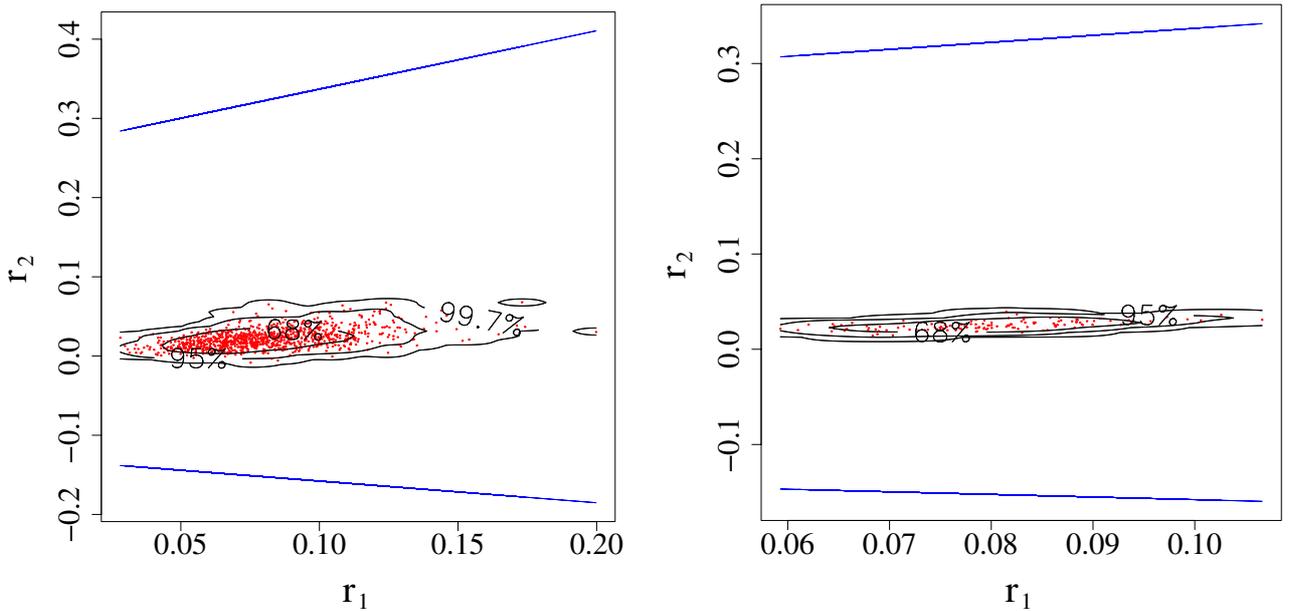


Fig. 8. Correlation coefficients r_1 and r_2 measured from the halo catalogs in the subcubes of the Millennium Simulation, using the LS estimator, where we slice the simulation volume into 1000 subcubes for the *left panel* and 125 for the *right one*, and the random catalog for each subcube contains 5000 halos. In both cases, we measure ξ at lags of separation $\Delta s = 5 h^{-1}$ Mpc, and use all halo pairs with pair separations of 0 to 2 h^{-1} Mpc to compute ξ_0 . The red dots show the 1000 (125) realizations, while the black lines are iso-density contours that contain the given percentages of the realizations. The upper and lower constraints $r_{2u,l}(r_1)$, computed individually for each realization of r_1 , are shown as blue lines.

the transformed values x_n , the bounds are simply ± 1 , so there is no need to average over the realizations.

Figure 9 shows box-and-whisker plots of r_n and x_n at all eight lags n , where we use the same lags and random catalog size as before, as well as the LS estimator. It can be seen that the constraints are clearly obeyed, and although the distributions becoming broader for increasing lags, the boxes showing the upper and lower quartiles only occupy a small portion of the allowed region. The distributions are not necessarily centered within the allowed region, which is not surprising, since their exact shape

and position also depend on the underlying power spectrum. The choice of estimator has barely any impact on the widths and positions of the distributions within the allowed region, which is to be expected, since Fig. 7 already illustrates that the different estimators yield quite similar results.

As mentioned at the beginning of this section, the main influence on the variance of the distributions in ξ , and correspondingly in r - and x -space, seems to be the width of the ξ_0 -bin. In Fig. 10 we investigate this observation and also study the impact of the choice of the separation between the lags at which

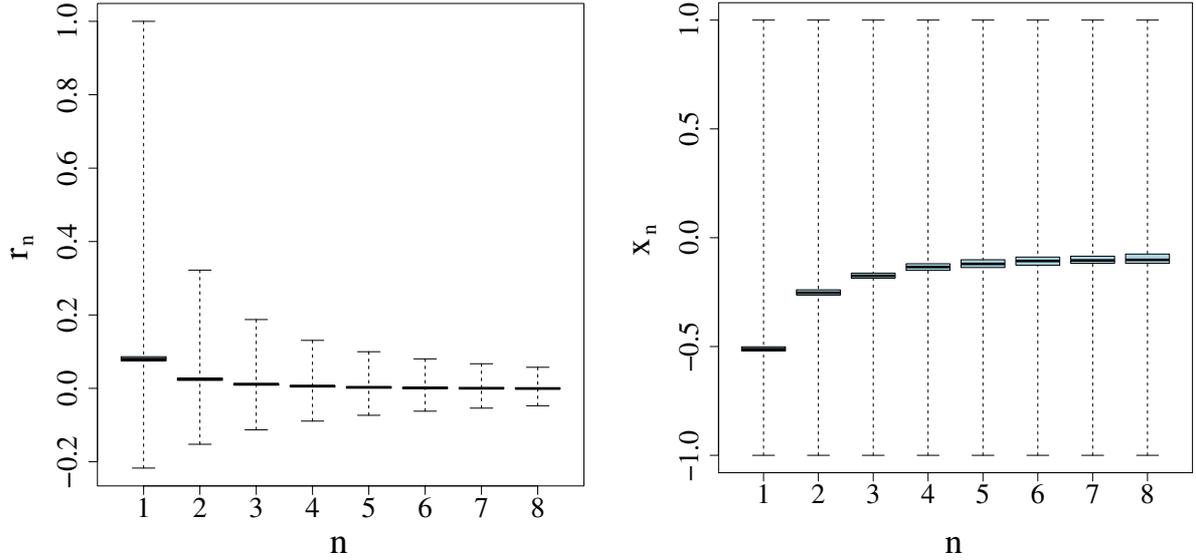


Fig. 9. Box-and-whisker plots for the correlation coefficients r_n and the transformed values x_n , as defined in Eq. (7), where each data point shows the upper and lower quartile (edges of the box), median (line inside the box), and mean upper and lower boundaries (whiskers) for the 125 realizations measured from the simulation subcubes. These use the same lags, random catalog size, and estimator as before (see text and previous figure captions for details).

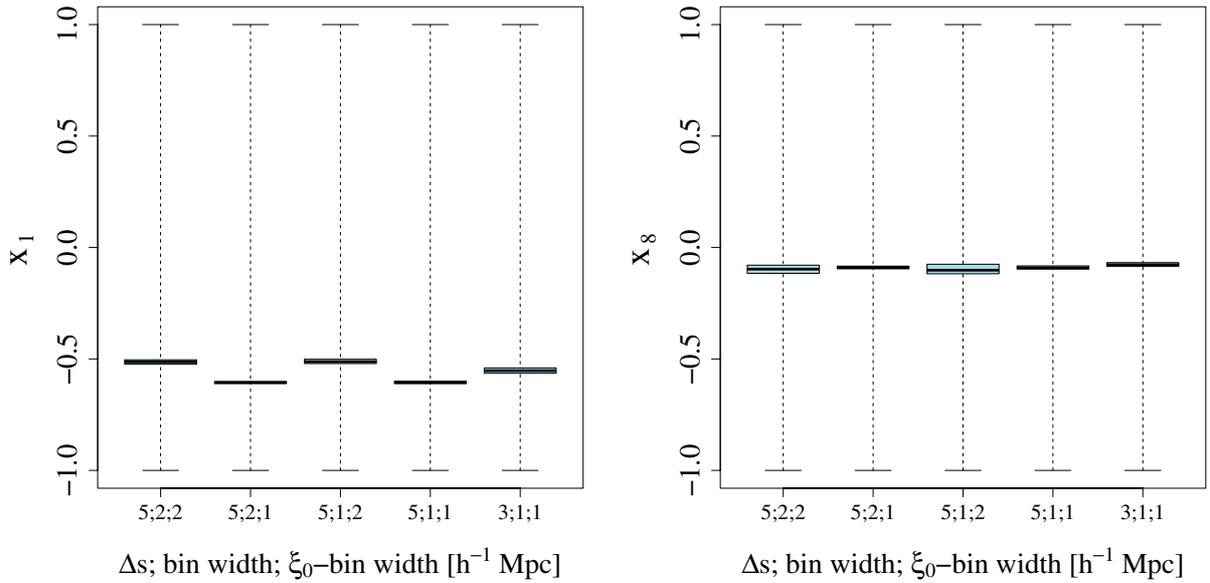


Fig. 10. Box-and-whisker plots of the transformed correlation coefficients at smallest and largest lag for varying lag separation and bin width. The triple labeling of each distribution gives the lag separation Δs , the bin width of pair separations at non-zero lag (i.e., for $\xi_1 \dots \xi_8$), and the width of the ξ_0 -bin. For example, for the second case shown in each panel, we measure ξ_0 from halo pairs with separations from 0 to $1 h^{-1}$ Mpc, ξ_1 from those with separations from 4 to $6 h^{-1}$ Mpc, ξ_2 from 9 to $11 h^{-1}$ Mpc, and so on.

we measure ξ . In the two panels of the figure, we show box-and-whisker plots of the transformed correlation coefficients at smallest and largest lag, i.e., x_1 and x_8 , and we vary the separation Δs of the lags as well as the bin widths of the pair separations used to measure ξ_0 and the correlation functions at non-zero lag, $\xi_1 \dots \xi_8$. In the case of the four left-most distributions in each panel, we use a lag separation of $\Delta s = 5 h^{-1}$ Mpc, where we adapt a bin width of $1 h^{-1}$ Mpc for $\xi_1 \dots \xi_8$ for the first and second distribution, and a bin width of $2 h^{-1}$ Mpc for the third and fourth ones. In both cases, we separately use a narrow and a broad bin width for the measurement of ξ_0 (also 1 and $2 h^{-1}$ Mpc). The figure illustrates that the width of the distributions of x_n is mainly determined by the ξ_0 -bin size, whereas the

width of the bins for ξ_n at lags $n > 0$ barely has any influence. This is due to the structure of the quasi-Gaussian transformation, where ξ_0 appears in every correlation coefficient r_n , and, as a result, in the computation of every lower and upper bound. The impact of the width of the ξ_0 -bin is particularly strong for small-lag distributions, which it also shifts, as can be seen from the distribution of x_1 . In particular, this shift is larger than a case where we measure ξ_n at different lags altogether, as illustrated for a lag separation of $\Delta s = 3 h^{-1}$ Mpc in the fifth distribution shown in the figure. In this context, it is important to stress that the problem of how to measure ξ_0 in practice is well-known since, in most applications, it is difficult to measure ξ at very small lags. As we show, however, this poses a particularly difficult

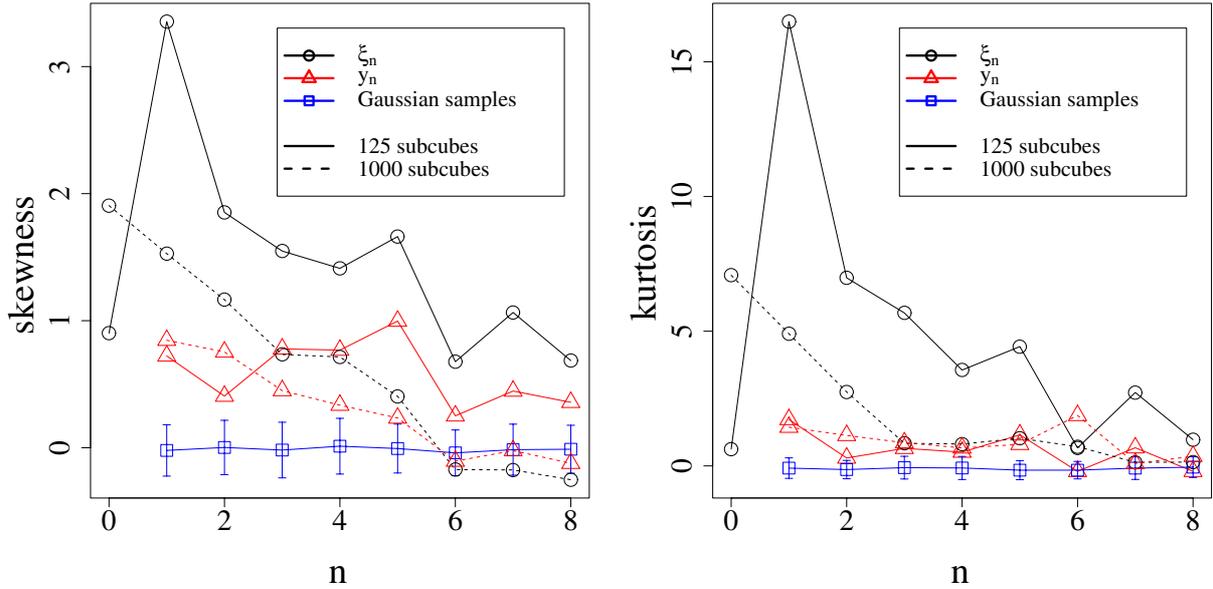


Fig. 11. Test for the univariate Gaussianity of the $\{\xi\}$ - and $\{y\}$ -samples obtained from the Millennium Simulation, using a lag separation of $\Delta s = 5 h^{-1}$ Mpc and bin widths of $1 h^{-1}$ Mpc for all ξ_n , including ξ_0 . The black circles and red triangles show the univariate skewness and kurtosis of the distributions $p(\xi_n)$ and $p(y_n)$, computed over 125 (solid curves) and 1000 (dashed curves) subcubes of the simulation volume. For the blue curves, we draw 100 Gaussian samples with the same mean, covariance matrix and sample size as the distributions $p(y_n)$ in the case of 125 subcubes, and plot the mean and standard deviation of their skewness and kurtosis.

challenge when analyzing measured correlation functions in a quasi-Gaussian framework, since here, the exact determination of ξ_0 is vital; the auto-correlation function enters everywhere, because one would always transform ξ to y (or at least to r) for an analysis involving the constraints.

In summary, all correlation functions we measured from the Millennium Simulation are quite far away from the edge of the allowed region. This finding seems to hold irrespective of the choice of estimator, lags, etc., providing that ξ is measured in a “sensible” way. As an example, using very small random catalog sizes does indeed yield single realizations outside the allowed region.

3.3. Quality of the Gaussian approximation in ξ and y -space

In this section, we use the correlation function samples measured from the Millennium Simulation to assess the quality of a quasi-Gaussian approach. Similar to the tests shown in WS2013 for simulated correlation functions, we transform ξ to y as defined in Eq. (5) and test the Gaussianity of the distributions in y and ξ , because the Gaussianity in y -space is a central prerequisite for the accuracy of the quasi-Gaussian likelihood.

While it would be preferable to assess the quality of the quasi-Gaussian approximation directly, i.e., to check how well the quasi-Gaussian PDF agrees with $p(\xi)$, as obtained from the Millennium Simulation, computing the quasi-Gaussian PDF still requires measuring the underlying power spectrum, which is beyond the scope of this work. In real life, however, one would usually transform the measured correlation function to y -space to perform a Bayesian analysis and, thus, the Gaussianity of $p(y)$ is pivotal. Even so, knowledge about the underlying power spectrum would still be required to make use of the analytically known $p(\xi_0)$.

In the literature, various tests for Gaussianity exist. In this study, we focus on the calculation of moments; in particular, we compute the skewness and kurtosis, which are defined in such a way that they are zero for a Gaussian distribution. In the

univariate case, the skewness γ of a distribution $p(x)$ reads

$$\gamma = \left\langle \frac{(x - \mu)^3}{\sigma^3} \right\rangle \equiv \frac{m_3}{m_2^{3/2}}, \quad (8)$$

where $m_i = \langle (x - \mu)^i \rangle$ denotes the central i th-order moment. Thus, γ is essentially the (renormalized) third-order moment, and the kurtosis

$$\kappa = \left\langle \frac{(x - \mu)^4}{\sigma^4} \right\rangle - 3 \equiv \frac{m_4}{m_2^2} - 3 \quad (9)$$

is closely related to the fourth-order moment. In the multivariate case, we use the definitions established by [Mardia \(1970, 1974\)](#), who define the skewness of a d -variate distribution as

$$\gamma_d = \frac{1}{n^2} \sum_{j=1}^n \sum_{i=1}^n \left\{ (x_i - \mu)^T C^{-1} (x_j - \mu) \right\}^3, \quad (10)$$

where n is the sample size, and μ and C are the sample mean and covariance matrix. The kurtosis measure reads

$$\kappa_d = \frac{1}{n} \sum_{j=1}^n \left\{ (x_j - \mu)^T C^{-1} (x_j - \mu) \right\}^2 - d(d + 2), \quad (11)$$

where we subtract the last term to ensure that a perfectly Gaussian sample yields $\kappa_d = 0$.

To test the impact of the quasi-Gaussian transformation on Gaussianity, we transform each realization of the correlation function (measured for eight lags of separation $\Delta s = 5 h^{-1}$ Mpc with bins of width $1 h^{-1}$ Mpc for $\xi_0 \dots \xi_8$) to y and compute skewness and kurtosis of the distributions in ξ - and y -space. Analogously to our tests in WS2013, we also draw Gaussian samples with the same mean and covariance matrix as our samples $\{y\}$, both for comparison and to account for small sample sizes. The results for the univariate distributions are plotted in [Fig. 11](#). Here, we show the skewness and kurtosis of the distributions $p(\xi_0), \dots, p(\xi_8)$ and $p(y_1), \dots, p(y_8)$; for the solid lines,

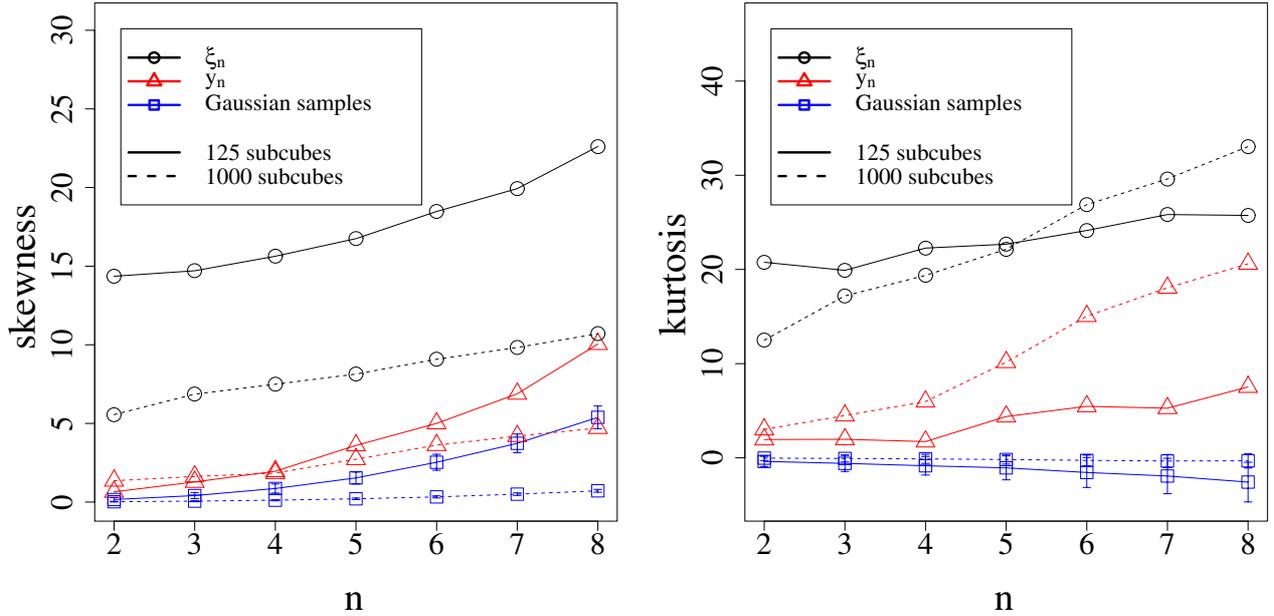


Fig. 12. Multivariate Mardia’s skewness and kurtosis of the n -variate distributions of the $\{\xi\}$ - and $\{y\}$ -samples, obtained from the Millennium Simulation, and of corresponding Gaussian samples, using the same parameters as before (see previous figure caption). In contrast to the previous figure, we plot two curves for Gaussian samples: for the solid (dashed) curve, we draw Gaussian samples with the same mean, covariance matrix, and sample size as the corresponding distributions in y -space for the case of 125 (1000) subcubes; the blue squares and error bars show the mean and standard deviation of the skewness and kurtosis of the 100 samples.

we sliced the simulation volume into 125 subcubes, whereas the dashed lines correspond to the case of 1000 subcubes. By way of comparison, the blue curves show the skewness and kurtosis of corresponding Gaussian samples, where for the sake of clarity, we only plot the curves for 125 subcubes. Because skewness and kurtosis fluctuate quite significantly for this small sample size, we draw 100 Gaussian samples of size 125 and compute the skewness and kurtosis of each sample; the blue squares and error bars show the mean and standard deviation of the skewness and kurtosis of the 100 samples. We only include the purely Gaussian samples for comparison and to check how close to zero the measured skewness and kurtosis are for these small sample sizes. Specifically, they are not meant to provide any insight into the absolute scale of the non-Gaussianity observed in ξ . Evidently, the distributions in y are far more Gaussian than those in ξ (with the exception of $p(\xi_0)$ in the case of 125 subcubes) and, of particular note, show a kurtosis comparable to the Gaussian samples.

Since the Gaussianity of the univariate distributions does not imply Gaussianity of the multivariate PDFs, we also compute the moments of the n -variate distributions $p(\xi_0, \dots, \xi_{n-1})$, $p(y_1, \dots, y_n)$ and of corresponding multivariate Gaussian samples, and plot them as functions of n , as in Fig. 12. Here, we show the results for corresponding Gaussian samples for both 125 and 1000 subcubes, where the plotted values and error bars are the mean and standard deviation of the skewness and kurtosis computed over 100 Gaussian samples. While the multivariate moments of the Gaussian samples of size 125 are not consistent with zero, this is indeed the case for the larger sample size of 1000. For the dashed curves, i.e., the case of 1000 subcubes, in the case of higher n , the integral constraint has a non-negligible impact on the measured correlation functions, as explained in the previous section. As a consequence, the corresponding skewness and kurtosis results should only be considered quantitatively. Nevertheless, it is clear that the difference between the level of Gaussianity in ξ - and y -space becomes even larger for

the multivariate case, reaching about one order of magnitude in γ and κ in the case of 125 subcubes.

As we demonstrated in the previous section, the width of the ξ_0 -bin, i.e., the range of pair separations used to measure the auto-correlation function, has an impact on the distributions of the correlation coefficients, and thus on those of the y_n . Accordingly, we vary the ξ_0 -bin width and again study the multivariate moments of the corresponding distributions. Figure 13 shows a similar plot to Fig. 12, but we use a ξ_0 -bin width of 2 instead of $1 h^{-1}$ Mpc. As it turns out, this yields distributions in y -space which are almost perfectly Gaussian, at least in the case of 125 subcubes, where their moments are hardly distinguishable from those of the corresponding Gaussian samples with same sample size. As before, it seems that the width of the ξ_0 -bin has a far higher impact on the results than the bin widths for $\xi_1 \dots \xi_8$. Actually, using bins of $2 h^{-1}$ Mpc for the higher-lag correlation functions barely influences the outcome.

In summary, all tests shown in this section indicate that the distributions in y -space are far more Gaussian than those in ξ , and in some cases even have skewness and kurtosis comparable to those of Gaussian samples of the same size. This demonstrates the validity of the quasi-Gaussian approach independent of the specific parameters used to measure the correlation function.

4. Conclusions and outlook

Building on SH2009, we have developed numerical methods to compute the fundamental constraints on correlation functions. We have shown these methods, which are applicable also in the case of two- and three-dimensional random fields, to be robust and precise, since the numerical computation of the constraints for the one-dimensional case reproduces the analytically known bounds. We then applied our results to samples of correlation functions measured from the halo catalog of the Millennium Simulation. After discussing some challenges in the

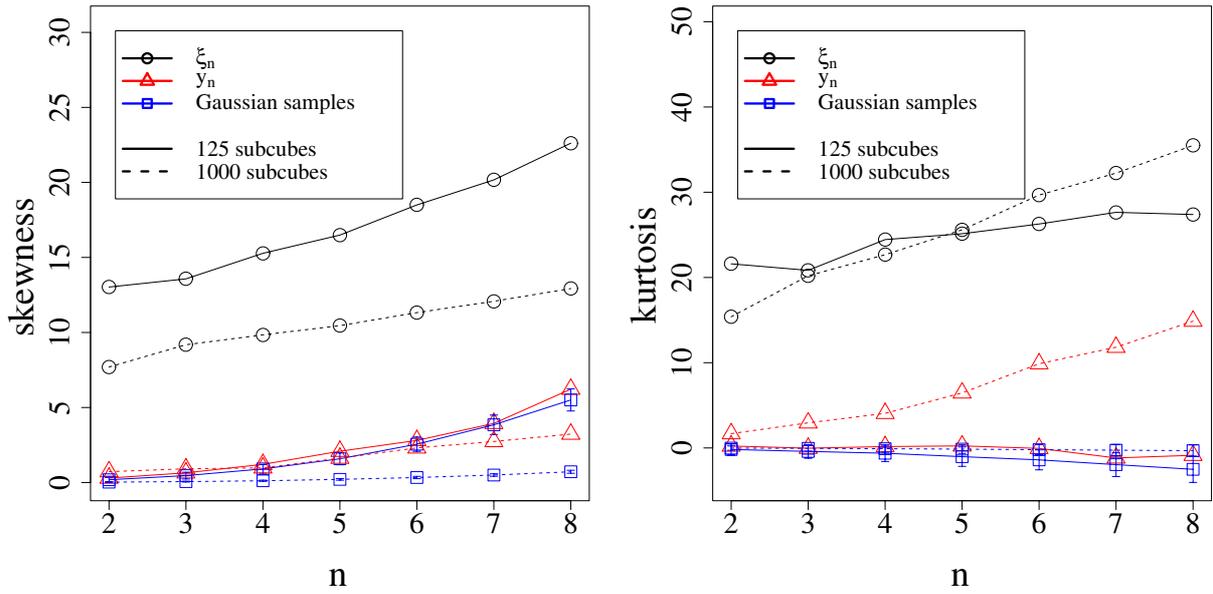


Fig. 13. Multivariate skewness and kurtosis of the $\{\xi\}$ - and $\{y\}$ -samples obtained from the Millennium Simulation and of corresponding Gaussian samples. When compared to the previous figure, we adapt a broader ξ_0 -bin, i.e., we measure the auto-correlation function from all halo pairs with pair separations from 0 to $2 h^{-1}$ Mpc.

measurement of ξ , such as the choice of random catalog size and lag separation, as well as the question of how to overcome the integral constraint, we have shown that the correlation functions measured from the simulation very clearly obey the constraints. Even though all measured correlation functions lie far away from the edges of the allowed region, we have demonstrated that the quasi-Gaussian quantity y yields significantly smaller non-Gaussian signatures (i.e., skewness and kurtosis) than the original correlation function ξ , giving further support to the claim that the quasi-Gaussian approximation for the correlation function likelihood, introduced in WS2013, is a far better description than the Gaussian one.

As a brief outlook on possible future work, one vital improvement would be to bypass the current limitation to eight lags in the numerical computation of the constraints, since modern astronomical observations usually measure ξ at far more lags. Furthermore, the performance of the quasi-Gaussian likelihood in the three-dimensional case should be assessed and compared to the classical Gaussian approach. While this would be testable on the samples of correlation functions measured from the Millennium Simulation, the most significant advance would be the application of our methods to real data, and an investigation of their impact on cosmological parameter estimation. Aside from the current limitation to only eight lags, this would pose additional challenges, depending on the area of application: in the case of a redshift survey, for example, different constraints on the correlation function, measured along and perpendicular to the line-of-sight, would hold as a result of redshift space distortions. Nonetheless, the constraints on correlation functions of three-dimensional random fields are, in principle, treatable, despite open challenges and room for improvements. Taking these

things into consideration, this work opens up a vast field of applications where Gaussian likelihoods for ξ have previously been used.

Acknowledgements. We would like to thank Cristiano Porciani for useful discussions and help. We also thank our anonymous referee for constructive comments and suggestions. The Millennium Simulation databases used in this paper, and the web application providing online access to them, were constructed as part of the activities of the German Astrophysical Virtual Observatory (GAVO). This work was in part supported by the Deutsche Forschungsgemeinschaft, under the project SCHN 342/11.

References

- Barber, C. B., Dobkin, D. P., & Huhdanpaa, H. 1996, *ACM Transactions on Mathematical Software*, 22, 469
- Carron, J. 2013, *A&A*, 551, A88
- Davis, M., & Peebles, P. J. E. 1983, *ApJ*, 267, 465
- Hamilton, A. J. S. 1993, *ApJ*, 417, 19
- Hartlap, J., Schrabback, T., Simon, P., & Schneider, P. 2009, *A&A*, 504, 689
- Hewett, P. C. 1982, *MNRAS*, 201, 867
- Keitel, D., & Schneider, P. 2011, *A&A*, 534, A76
- Kerscher, M., Szapudi, I., & Szalay, A. S. 2000, *ApJ*, 535, L13
- Labatie, A., Starck, J.-L., & Lachièze-Rey, M. 2012, *ApJ*, 746, 172
- Landy, S. D., & Szalay, A. S. 1993, *ApJ*, 412, 64
- Mardia, K. 1970, *Biometrika*, 57, 519
- Mardia, K. 1974, *Sankhyā. Series B. Methodological*, 36, 115
- Martin, S., Schneider, P., & Simon, P. 2012, *A&A*, 540, A9
- Peebles, P. J. E., & Hauser, M. G. 1974, *ApJS*, 28, 19
- Schneider, P., & Hartlap, J. 2009, *A&A*, 504, 705
- Seljak, U., & Bertschinger, E. 1993, *ApJ*, 417, L9
- Simon, P., Semboloni, E., van Waerbeke, L., et al. 2015, *MNRAS*, 449, 1505
- Springel, V., White, S. D. M., Jenkins, A., et al. 2005, *Nature*, 435, 629
- Sun, L., Wang, Q., & Zhan, H. 2013, *ApJ*, 777, 75
- Vargas-Magaña, M., Bautista, J. E., Hamilton, J.-C., et al. 2013, *A&A*, 554, A131
- Wilking, P., & Schneider, P. 2013, *A&A*, 556, A70