

A statistical formalism for alignment analysis

F. Dávila-Kurbán,^{1,2,3,4} M. Lares,^{1,2,3}★ D. Garcia Lambas^{1,2,3}

¹*Instituto de Astronomía Teórica y Experimental (IATE, CONICET/UNC), Córdoba, Argentina*

²*Observatorio Astronómico Córdoba, Argentina*

³*Consejo de Investigaciones Científicas y Técnicas (CONICET), Argentina*

⁴*Facultad de Matemática, Astronomía y Física, Universidad Nacional de Córdoba, Argentina*

Accepted XXX. Received YYY; in original form ZZZ

ABSTRACT

The detection of anisotropies with respect to a given direction in a vector field is a common problem in astronomy. Several methods have been proposed that rely on the distribution of the acute angles between the data and a reference direction. Different approaches use Monte Carlo methods to quantify the statistical significance of a signal, although often lacking an analytical framework. Here we present two methods to detect and quantify alignment signals and test their statistical robustness. The first method considers the deviance of the relative fraction of vector components in the plane perpendicular to a reference direction with respect to an isotropic distribution. We also derive the statistical properties and stability of the resulting estimator, and therefore does not rely on Monte Carlo simulations to assess its statistical significance. The second method is based on a fit over the residuals of the empirical cumulative distribution function with respect to that expected for a uniform distribution, using a small set of harmonic orthogonal functions, which does not rely on any binning scheme. We compare these methods with others commonly used in the literature, using Monte Carlo simulations, finding that the proposed statistics allow the detection of alignment signals with greater significance.

Key words: Methods: statistical – Methods: numerical

1 INTRODUCTION

Shapes and orientations of galaxies and the large scale structures in which they are embedded may have significant coherence given the effects of accretion and mergers as well as tidal, stripping and other combined actions. As a consequence of these processes, the statistical properties of the galaxy orientations with respect to the cosmic web structures may differ from those expected for randomly oriented galaxies (e.g. see Mo et al. 2010). Thus, studies of intrinsic alignment signals allow to explore the links between the joint evolution of galaxies and their surrounding structures (e.g. Panko et al. 2013, and references therein). Taking into account these facts, the analysis of the orientations of galaxies in the context of both the local environment and the large-scale structures may be crucial to test scenarios of galaxy formation and evolution, in particular for theoretical predictions of the angular momentum of galaxies (e.g., Peebles 1969). The alignment signals, however, are somewhat elusive, given the variety of preferred directions that arise from the actual distribution of surrounding structures and the fact that galaxy orientations with respect to any direction are mainly random. For these reasons a robust statistical method to detect and assess alignment signals and their significance is a key tool in studies of alignments between galaxies and the large-scale distribution of structures.

In the case of spiral galaxies, the spatial distribution of stars in a disc defines a preferred plane whose normal is oriented roughly onto the rotation axis. The tidal field exerted by regions characterized by

structures such as clusters, filaments or voids, are present during a considerable extent of galaxy evolution, and could produce observable features in their original spin vector. The fact that galaxies rotate are indicative of the physical conditions under which they formed, and the rotation itself is certainly an important test of any theory for the origin of the galaxies (Peebles 1969). Galaxy angular momentum is widely believed to arise from gravitational torque due to misalignment of the gravitational shear tensor and the inertia tensor in early formation stages (Doroshkevich 1970; White 1984). Thus, the galaxy spin field holds information about the gravitational shear field and can be used, for example, for a statistical reconstruction thereof (Lee & Pen 2000). Furthermore, it is commonly assumed that, during early stages of formation, baryonic and dark matter shared a similar evolution and likely gained the same specific angular momentum prior to the formation of the disc (e.g., Fall & Efstathiou 1980). The study of alignment dark matter haloes (hereafter DM haloes) was possible, and became a popular subject, after N-body simulations had enough resolution to perform studies of this nature (e.g. Cuesta et al. 2008; Libeskind et al. 2013; Forero-Romero et al. 2014; Joachimi et al. 2015; Kiessling et al. 2015, and references therein).

The methods of alignment detection and the results obtained are diverse. For example, Forero-Romero et al. (2014) studied alignment of shape, angular momentum, and peculiar velocity of DM haloes with respect to the cosmic web, as described by using the tidal field or velocity shear, employing the Bolshoi simulation (Riebe et al. 2011). They quantify the alignments by measuring the fraction of haloes that is preferentially aligned with one of the eigenvectors in the local definition of the cosmic web, and with the average value of the angle

★ E-mail: marcelo.lares@unc.edu.ar

between an eigenvector and the vector of interest. They found the strongest alignment for halo shapes with filaments and walls defined by the tidal field, but when defined by velocity shear they found anti-alignment with massive haloes. For the angular momentum, they only found a weaker signal for the most massive haloes to be anti-aligned with filaments, and being aligned along the sheets of the velocity shear. There is a discrepancy with previous works (Aragón-Calvo et al. 2007; Hahn et al. 2007; Aragon-Calvo & Yang 2014) which indeed detect alignments for less massive haloes. Forero-Romero et al. (2014) argues that this might be due to high sensitivity of the alignment signal to the small-scale cosmic web description. Additionally, they find peculiar velocities to be preferentially parallel to walls and filaments. These results indicate that the alignment properties of DM haloes can depend on the physical definition of the cosmic web, with tidal field versus velocity shear approach yielding complementary information.

With the greater computing power of recent years, the study of alignments of galaxies in simulation has been possible. Codis et al. (2018) and Kraljic et al. (2019), for example, study the distribution of angles measured between the spin of galaxies and haloes and the different elements of the surrounding cosmic web in the Horizon-AGN and SIMBA simulations, respectively. Their results agree on the spin of low-mass galaxies being more likely to lie within the plane of sheets while massive galaxies preferentially having a spin perpendicular to the sheets.

The search for galactic alignment has been analyzed also in observations in the context of structures that, to a reasonable extent, can be described with spherical symmetry, like clusters of galaxies or voids. The observational aspect of this topic of study has its own difficulties to face, mainly the small sample sizes and line-of-sight projection effects. Earlier works focused on the orientation of galaxies with respect to the Local Supercluster and other clusters such as Virgo and Coma (e.g. Kashikawa & Okamura 1992; Godłowski 1993; Godłowski 1994; Hu et al. 1995; Wu et al. 1997; Yuan et al. 1997; Hu et al. 1998; Godłowski & Ostrowski 1999) relied on the "position angle (PA) – inclination method" (Jaanieste & Saar 1978; Flin & Godłowski 1986). In this method, the measured PAs of galaxies (usually on photographic plates) are converted into 3-dimensional vectors using inclination angles obtained from the measured projected minor-to-major axial ratios, b/a . The distribution of these vectors could then be compared with a null-hypothesis, e.g. isotropic spatial distribution, and thus assess whether the data is isotropic or anisotropic by comparison. However, the shape of these isotropic distributions can be significantly affected by selection criteria (Aryal & Saurer 2000). These effects can be large when the sample is selected from incomplete datasets (e.g. a limited portion of the sky) and lead to artificial structures in the data. Therefore, a statistically robust method that reliably describes not only the data, but the comparison sample as well, is of crucial importance in these analyses in order to conclude in favor of either isotropy or anisotropy in the data.

Other observational studies employ similar methods that also rely on binning statistics such as the normalized pair count, $P(\cos\theta)$ in bins of the measured angle θ between the subject of interest (e.g. satellite or central galaxies or otherwise) and a preferred direction determined by some other structure such as cluster centers or elements of the cosmic web (e.g. Brainerd 2005; Yang et al. 2006; Zhang et al. 2015). The significance of these statistics is usually assessed by comparison with a large number of Monte Carlo simulations.

Varela et al. (2012, hereafter V12) performed a rigorous assessment of an analytical model for the distribution of θ , and its behaviour in the isotropic case for the estimation of the statistical significance, based on previous works (e.g. Betancort-Rijo & Trujillo 2009; Brunino et al. 2007; Cuesta et al. 2008; Lee 2004). This work tackled some

discrepancies that emerged in previous observational studies of the alignment of galaxies around voids, namely Trujillo et al. (2006) and Slosar & White (2009) (hereafter T06 and S09, respectively). T06 analyzed 201 face-on and edge-on galaxies using data from the SDSS-DR3 and the 2dFRGS (Colless et al. 2001) and found a significant tendency of the spin of the galaxies to be in the direction perpendicular to the void radial direction. On the other hand, S09 using two samples of 578 and 258 galaxies from the SDSS-DR6 with similar selection criteria found no statistical evidence for departure from random orientations. V12, used the SDSS-DR7 and a statistical procedure robust enough to overcome the problem of the indeterminacy of the real inclination of galaxies computed from their apparent axial ratio, and assess the validity of the procedure with extensive Monte Carlo simulations. They detect a statistically significant tendency of galaxies around large voids (with radii of over $15 h^{-1} \text{Mpc}$) to have their angular momenta aligned with the radial direction of the voids. This highlights the importance of, not only a bigger sample size, but the use of robust and reliable statistical methods to correctly assess the validity of the alignment signal detected.

In this paper we present two formal methods to analyze the alignments of a sample of particles with respect to a center. The first method consists on the definition of simple metrics from the radial and tangential components of the vectors, while the second one relies on the parametrization of a residual function between the data obtained from the sample and from an isotropic distribution. On either case, we do not assume any binning scheme. Instead, we use all the information in the data and apply robust estimations of the uncertainties in the alignments metrics. By deriving the theoretical distribution of the parameters that measure alignment signal we can not only determine its statistical significance with accuracy, but we can do so without investing computational resources and time into the Monte Carlo simulations usually needed to estimate this.

The outline of this paper is as follows. In Section 2 we develop the statistical formalism for the aforementioned two new methods for the study of vector alignments, and introduce the parameters with which to measure alignment signal. In Section 3 we apply the methods to synthetic data corresponding to different scenarios of alignment to test how well the new parameters recover the alignment signal, and compare it to more traditional methods. Finally, Section 4 presents our main conclusions.

2 METHODS

The symmetry of structures such as filaments, haloes, The symmetry of structures such as filaments, haloes, clusters or voids, both in their geometry as in their dynamics,

clusters or voids, both in their geometry as in their dynamics, allows the consideration of a preferential direction to analyze the orientation of galaxies. In the case of spherical symmetry, this is the radial direction. The objective is to develop a statistical formalism to measure in a robust manner the distribution of the orientation of galaxies and detect possible excesses with respect to a completely random distribution. Given the problem of vector orientations with respect to a central point we want to define a statistical parameter and obtain its distribution in order to know the significance of a hypothesis test.

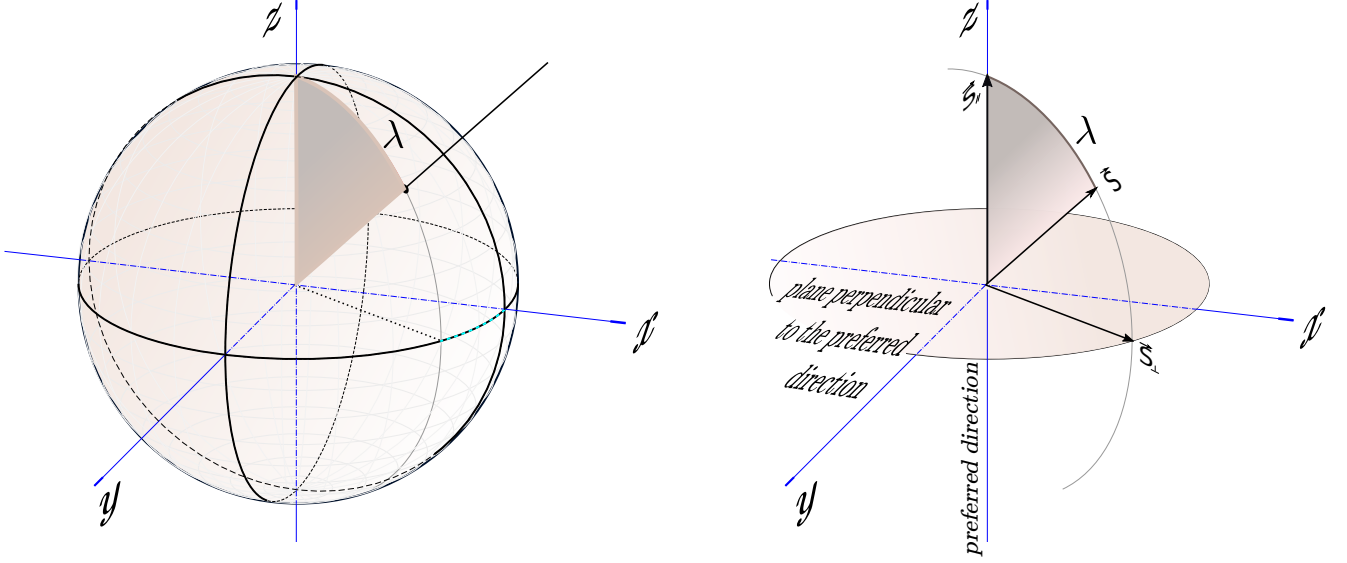


Figure 1. Coordinate system used. The z axis is the radial outward direction of the void. The angle θ is formed between the z and the vector \vec{S} , and takes values in the range $[0, \pi]$.

2.1 Ratio of vector components

Given a radial direction \hat{z} of unit norm (see Fig. 1), perpendicular and parallel components of vector \vec{S} can be calculated as:

$$\vec{S}_{\parallel} = \vec{S} \cdot \hat{z}, \quad \text{and} \quad \vec{S}_{\perp} = \vec{S} - \vec{S}_{\parallel}, \quad (1)$$

where \vec{S}_{\perp} is the perpendicular component to the radial direction \hat{z} , \vec{S}_{\parallel} is the parallel component to the radial direction and $\vec{S} = \vec{S}_{\perp} + \vec{S}_{\parallel}$.

The angle θ formed by the radial direction and the direction of the vector \vec{S} is related to the components:

$$S_{\perp} = |\vec{S}| \sin(\theta); \quad S_{\parallel} = |\vec{S}| \cos(\theta). \quad (2)$$

The distribution of this angle can be used to analyze alignments, and given its relation to the components, the latter can also be used to determine the orientations. To that effect we define:

$$\mathcal{B} = \frac{S_{\perp}}{S_{\parallel}} = \frac{S \sin(\theta)}{S \cos(\theta)} = \tan(\theta). \quad (3)$$

The parameter \mathcal{B} is also a measure of the orientation of the vector \vec{S} . Note that the ranges of these two parameters are as follows:

$$0 \leq \theta \leq \pi; \quad -\infty \leq \mathcal{B} \leq \infty$$

Using the symmetry of the problem, we can define the parameters considering the acute angle between the directions \hat{z} and \hat{S} , as well as the norm of the component \vec{S}_{\parallel} ,

$$\beta = |\mathcal{B}| = \frac{S_{\perp}}{|S_{\parallel}|}, \quad \text{with} \quad \lambda = \min(\theta, \pi - \theta) \quad (4)$$

for which:

$$0 \leq \lambda \leq \frac{\pi}{2}; \quad 0 \leq \beta \leq \infty$$

Vectors with $\beta > 1$ have a preference in the perpendicular direction and $\pi/4 < \lambda < \pi/2$, while vectors with $\beta < 1$ have a preference in the radial direction and $0 < \lambda < \pi/4$.

To quantify the direction of \vec{S} the following statistics could be considered:

- the angle θ
- the acute angle λ
- the ratio of perpendicular to parallel components, \mathcal{B} (with θ)
- the ratio of perpendicular to parallel components, β (with λ)

We will explore in the following subsections the use of the angles or the ratios. It is of utmost importance to establish the distributions of these parameters for the case in which there is no alignment signal whatsoever. This way we determine the amplitude of the statistical fluctuations and establish a measurement of the signal in a data sample calculating its statistical significance. To this end, we define the null hypothesis

H_0 : the distribution of vectors are random with spherical symmetry

i.e., there is no alignment signal whatsoever. This hypothesis can also be used to generate control samples with Monte Carlo procedures if need be.

Note that the regions of β greater or lesser than 1 are different, so it is expected that for a random distribution there would be more "perpendicular" than "parallel" vectors (see below Fig. 2, upper panel).

2.2 Distribution of the ratio of vector components

In this Section we derive the distribution of the test statistic β , defined as the ratio between the perpendicular and parallel vector components (Eq. 4), under the null hypothesis.

The distribution of β can be deduced from the change of random variables theorem (Gillespie 1983), which, in its general form, can be enunciated as follows:

Let $\{X_i\}_{i=1}^n$ be a R.V. with known $f_{\vec{X}}(\vec{x})$, and let m R.V. $\vec{Y} = \varphi(\vec{x})$, where $\varphi = (\varphi_1, \varphi_2, \dots, \varphi_m)^t$ and $\varphi_k : \mathbb{R}^n \rightarrow \mathbb{R}$ real functions. The joint probability function $f_{\vec{Y}}(\vec{y})$ is given by:

$$f_{\vec{Y}}(\vec{y}) = \int_{-\infty}^{\infty} d\vec{x} f_{\vec{X}}(\vec{x}) \prod_{i=1}^m \delta(y_i - \varphi_i(\vec{x})), \quad (5)$$

where δ is the Dirac Delta function. For the particular case of a unidimensional variable, $X : \Omega \rightarrow \mathbb{R}$ and $Y : \Omega \rightarrow \mathbb{R}$, with $Y = \varphi(X)$,

$$f_Y(y) = \int_{-\infty}^{\infty} dx f_X(x) \delta(y - \varphi(x)) \quad (6)$$

Then, we can use this theorem to find the distribution of β from $F_X(x) = U(0, 1)$ with the transformation:

$$\beta = \tan(\arccos(x)) \quad (7)$$

as well as from $f_{\Lambda}(\lambda) = \sin(\lambda)$ with the transformation

$$\beta = \tan(\lambda), \quad 0 < \lambda < \pi/2. \quad (8)$$

Using the latter, we have:

$$\begin{aligned} f_B(\beta) &= \int_{-\infty}^{\infty} d\lambda f_{\Lambda}(\lambda) \delta(\beta - \tan(\lambda)) \\ &= \int_0^{\pi/2} d\lambda \sin(\lambda) \delta(\beta - \tan(\lambda)) \end{aligned} \quad (9)$$

To solve this integral, we perform the change of variables:

$$z = \tan(\lambda) \implies \lambda = \arctan(z), \quad d\lambda = \frac{dz}{1+z^2}$$

Therefore,

$$f_B(\beta) = \int_0^{\infty} dz \frac{\sin(\arctan(z))}{1+z^2} \delta(\beta - z) = \frac{\sin(\arctan(\beta))}{1+\beta^2} \quad (10)$$

This expression can be simplified using the properties of trigonometric functions. Indeed, if $\beta = \tan(z)$ for a number z , then:

$$\begin{aligned} \beta^{-2} + 1 &= \frac{1}{\tan(z)^2} + 1 = \frac{\cos(z)^2}{\sin(z)^2} + 1 = \frac{\sin(z)^2 + \cos(z)^2}{\sin(z)^2} \\ &= \frac{1}{\sin(z)^2} \end{aligned}$$

$$\implies \frac{1}{\sin(z)} = \sqrt{\beta^{-2} + 1} = \sqrt{\frac{1+\beta^2}{\beta^2}} = \frac{\sqrt{1+\beta^2}}{\beta}$$

$$\implies \sin(z) = \frac{\beta}{\sqrt{1+\beta^2}}$$

$$\implies \sin(\arctan(\beta)) = \frac{\beta}{\sqrt{1+\beta^2}}$$

Replacing in Eq. 10 we have,

$$\begin{aligned} f_B(\beta) &= \frac{\sin(\arctan(\beta))}{1+\beta^2} \\ &= \frac{\beta}{\sqrt{1+\beta^2}} \frac{1}{1+\beta^2} \\ &= \beta(1+\beta^2)^{-3/2} \end{aligned} \quad (11)$$

Therefore, the probability function is:

$$\begin{aligned} F_B(\beta) &= \int_0^{\beta} f_B(b) db \\ &= \int_0^{\beta} b(1+b^2)^{-3/2} db \\ &= -\frac{1}{\sqrt{1+b^2}} \Big|_0^{\beta} \\ &= 1 - \frac{1}{\sqrt{1+\beta^2}} \end{aligned} \quad (12)$$

Figure 2, top panel, shows the theoretical distribution of β , with the Monte Carlo sampling of the random variable shown with the histogram.

Knowing the distribution f_B one can perform analyses of the orientations of vectors with respect to a particular direction. Generally, it is not useful to measure a single value of the R.V. β , given that it is subject to random fluctuations. Therefore, we calculate the values of the estimator β in a sample of observations. I. e., we analyze a random sample (R.S.) of values in order to determine if it differs from the expected results for a random distribution (a R.S. under the null hypothesis) of vectors. To formalize these analyses we need to establish some basic properties of the distribution f_B .

The first moment of the distribution, if it exists, is:

$$\begin{aligned} E[B] &= \int_0^{\infty} t f_B(t) dt \\ &= \int_0^1 t f_B(t) dt + \int_1^{\infty} t f_B(t) dt \end{aligned} \quad (13)$$

where

$$\int_1^{\infty} t f_B(t) dt = \int_1^{\infty} t \frac{t}{(1+t^2)^{3/2}} dt$$

Keeping in mind that for a real number $x > 1$ we have $x^n > x$, and $x > \sqrt{x}$, therefore $x^{3/2} = x\sqrt{x} < x$. Then, for $\beta > 1$, $1 + \beta^2 > 2 > 1$ and

$$\frac{1}{(1+t^2)^{3/2}} > \frac{1}{(1+t^2)}$$

Then we can limit the integral:

$$\begin{aligned} \int_1^{\infty} t \frac{t}{(1+t^2)^{3/2}} dt &> \int_1^{\infty} \frac{t^2}{(1+t^2)} dt \\ &> \int_1^{\infty} \frac{t}{(1+t^2)} dt \\ &> \lim_{M \rightarrow \infty} \frac{1}{2} \ln(1+t^2) \Big|_1^M = \infty. \end{aligned}$$

We see, then, that the expectation value $E[B]$ is undefined. Indeed, no moment of this distribution is defined. In fact, keeping in mind that:

$$\begin{aligned}
 E[B^n] &= \int_0^\infty t^n f_B(t) dt \\
 &= \int_0^1 t^n f_B(t) dt + \int_1^\infty t^n f_B(t) dt
 \end{aligned} \tag{14}$$

and that:

$$\beta > 1 \implies 1 + \beta^1 > 1 \implies \frac{\beta^n}{1 + \beta^2} > \frac{\beta}{1 + \beta^2}$$

for $n \geq 1$. Then,

$$\int_1^\infty t^n f_B(t) dt > \int_1^\infty t f_B(t) dt > \infty.$$

The distribution $f_B(\beta)$ is a pathological distribution where the moments are undefined. The properties of this distribution are similar to the properties of the Cauchy distribution. This limitation prevents from using Monte Carlo procedures to estimate the distribution of $\bar{\beta}$ because it is not possible to ensure that the average values of β follow a stable distribution. In order to work with this distribution, we could devise a truncated distribution, between arbitrary values l_1 and l_2 , with the condition that $l_1 \sim 0$ and l_2 be much larger than the region of interest of the parameter β , which is the region around $\beta = 1$.

For example, if we choose

$$L_1 = 10^{-3}; \quad L_2 = 10^3$$

it arises that, from defining the the correction factor κ :

$$\kappa = \int_{L_1}^{L_2} f_B(t) dt = \frac{1}{\sqrt{1+L_1^2}} - \frac{1}{\sqrt{1+L_2^2}} = \frac{1}{\sqrt{1+10^{-6}}} - \frac{1}{\sqrt{1+10^6}}$$

we can define an approximation to the distribution function for B defined as:

$$f_{\bar{B}}(\beta) = \begin{cases} \frac{1}{\kappa} f_B(\beta) & \beta \in [L_1, L_2] \\ 0 & \text{otherwise,} \end{cases}$$

The mean of this function is in fact defined, and its expression is as follows:

$$\begin{aligned}
 E[\bar{B}] &= \frac{1}{\kappa} \left[\frac{L_1}{\sqrt{L_1^2+1}} - \frac{L_1^2 \operatorname{asinh}(L_1)}{L_1^2+1} + \frac{L_2^2 \operatorname{asinh}(L_2)}{L_2^2+1} \right. \\
 &\quad \left. - \frac{L_2}{\sqrt{L_2^2+1}} + \frac{\operatorname{asinh}(L_2)}{L_2^2+1} - \frac{\operatorname{asinh}(L_1)}{L_1^2+1} \right]
 \end{aligned}$$

If we take $L_1 = 1/L_2$:

$$\begin{aligned}
 E[\bar{B}] &= \frac{1}{\kappa} \left[\frac{1/L_2}{\sqrt{L_2^{-2}+1}} - \frac{L_2^{-2} \operatorname{asinh}(1/L_2)}{L_2^{-2}+1} + \frac{L_2^2 \operatorname{asinh}(L_2)}{L_2^2+1} \right. \\
 &\quad \left. - \frac{L_2}{\sqrt{L_2^2+1}} + \frac{\operatorname{asinh}(L_2)}{L_2^2+1} - \frac{\operatorname{asinh}(1/L_2)}{L_2^{-2}+1} \right]
 \end{aligned}$$

However, the result is strongly dependant on the value of L_2 , and, to a lesser extent, the value of L_1 . Let

$$A(1/L_2, L_2) = \int_{1/L_2}^{L_2} t f_B(t) dt$$

It is straightforward to see that $A(1/L_2, L_2)$ depends on L_2 , where we took different values for L_2 and $L_1 = 1/L_2$.

With this we prove that it is not possible to obtain a distribution for $\bar{\beta}$. Furthermore, not only is it not possible to solve analytically, but it is also not possible to make a formal bootstrap estimation of the error. However, as we show in the next section, β can be used to define a new parameter with better statistical properties.

2.3 Test for the fraction of vectors in excess to random w.r.t. a reference direction

To find a robust estimator we consider the fraction of values of β that are greater than some critical value. Given that when the perpendicular and parallel components are equal there is no preference in any of the two directions, we can posit that said critical value be $\beta = 1$. Therefore, we define the parameter:

$$\hat{\eta} = \frac{N(\beta > 1)}{N(\beta < 1)} \tag{15}$$

where N is the number of observations of a sample that fulfills the conditions indicated between parentheses. Under H_0 , based on the probability density function, one expects that

$$\eta_0 = \frac{P(\beta > 1)}{P(\beta < 1)} \tag{16}$$

To calculate the value of η_0 , we take into account that, using the probability function, F_B (see Eq. 12):

$$P(\beta > 1) = 1 - F_B(1) = 1 - \left(1 - \frac{1}{\sqrt{1+\beta^2}} \right) \Big|_{\beta=1} = \frac{1}{\sqrt{2}}$$

$$P(\beta < 1) = F_B(1) = 1 - \frac{1}{\sqrt{1+\beta^2}} \Big|_{\beta=1} = 1 - \frac{1}{\sqrt{2}}$$

i.e.:

$$\begin{aligned}
 \eta_0 &= \frac{P(\beta > 1)}{P(\beta < 1)} = \frac{\int_0^1 f_B(t) dt}{\int_1^\infty f_B(t) dt} = \frac{\frac{1}{\sqrt{2}}}{1 - \frac{1}{\sqrt{2}}} = \frac{1}{\sqrt{2} - 1} \\
 &\cong 2.4142
 \end{aligned} \tag{17}$$

We propose that $\hat{\eta}$ is an estimator of η_0 , i.e., we need to check if:

$$E(\hat{\eta}) = \eta$$

Let there be a random sample of N values of β , we define:

$$n = N(\beta > 1)$$

Given that the probability of obtaining a value of $\beta > 1$ is $P(\beta > 1) = 1/\sqrt{2}$, the variable n has a binomial distribution,

$$f_n(n) = \operatorname{Bin}(p, N) = \binom{N}{n} p^n (1-p)^{N-n}$$

with $p = 1/\sqrt{2} \approx 0.707$. Therefore,

$$\eta = \frac{n}{N-n}$$

and the distribution of η can then be calculated from the distribution of n , taking into account that:

$$P_\eta \left(\eta = \frac{k}{N-k} \right) = P_n(n = k)$$

where $k = \frac{\eta N}{\eta - 1}$.

It is equivalent then, although much more efficient, to generate random variables of the distribution of η with this method than with

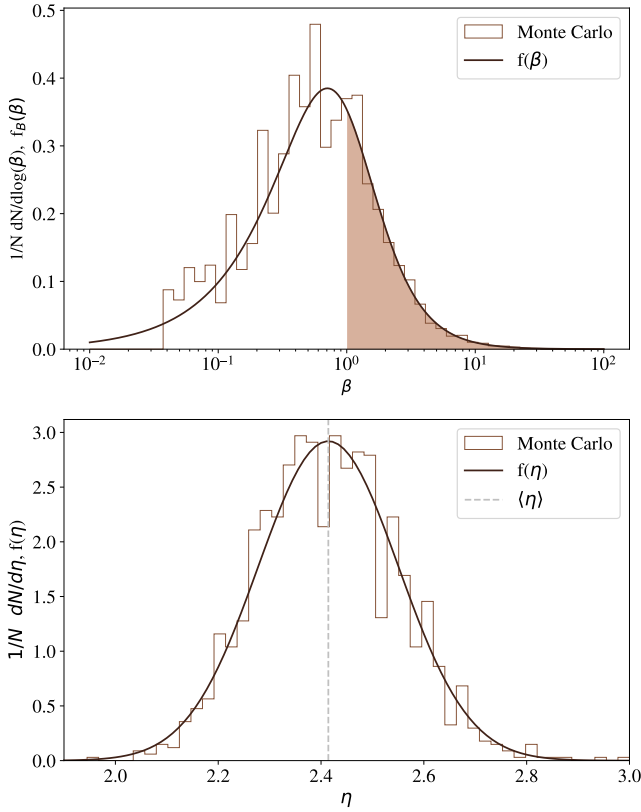


Figure 2. (upper:) Distribution of β obtained from the theoretical derivation (solid line) and from Monte Carlo simulations (histogram), and definition of η (fraction of samples with $\beta > 1$). (bottom:) Histogram of the η variables sorted with the Monte Carlo method (from samples of β) and with the theoretical distribution approximation. The mean theoretical value ($1/\sqrt{2}$) is shown in the dashed line and the histogram corresponds to a Monte Carlo realization of eta values, using

a Monte Carlo method. The comparison between the two random samples can be seen in the bottom panel of Figure 2.

Then, taking into account the fact that the expectation value of the variable $n \sim \text{Bin}(N, p)$ is np , we have to calculate the expectation value of the ratio. This problem is generally not well defined, but it can be solved approximately.

Let X and Y be R.V. defined as $X = n$, $Y = N - n$. If $q = 1 - p$, the expectation values of these variables are:

$$\mu_X = Np; \quad \mu_Y = N - Np = N - \mu_X$$

and the variances:

$$\sigma_X^2 = \sigma_Y^2 = NP(1 - p) = Npq$$

with $q = 1 - p$.

In the Appendix we show that the variance of this estimator is given by

$$\text{Var}(\eta) \approx \frac{28}{N}, \quad (18)$$

and Fig. 3 shows that for sample sizes larger than approximately 100, using the theoretical value is equivalent to using Monte Carlo simulations, with the advantage of needing comparatively no computation time.

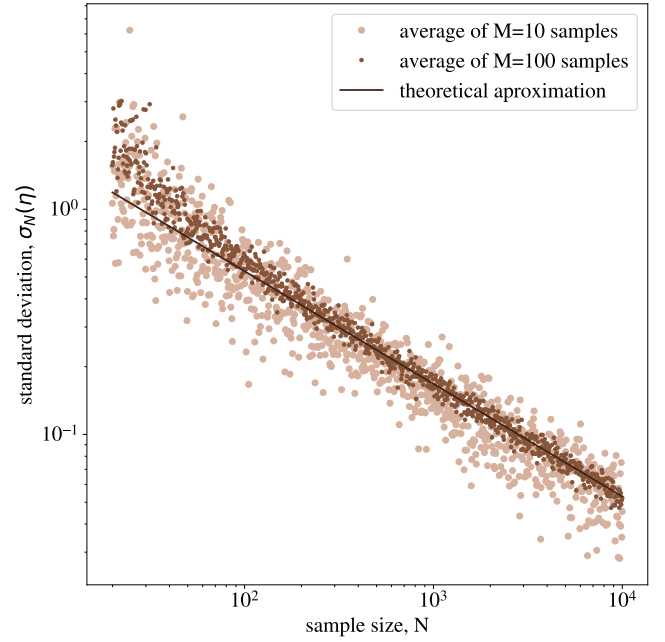


Figure 3. Variation of Monte Carlo estimations of the variance of M samples of values of η , calculated from N samples of values of β (N , in the X axis, simulated), for $M=10$ (large dots) and $M=100$ (small dots). The theoretical variance is shown to be an appropriate estimation for samples with size $M > 100$.

2.4 A test for the OLS coefficients of the cosine distribution

Another option is to analyze the distribution of $\cos(\lambda)$ to determine if it is distinguishable from the expected distribution of a sample of random orientations of vectors \vec{S} . As discussed, said distribution is uniform under H_0 . Working with samples of limited size, the statistical fluctuations can generate differences between the two sets of data, even when they arise from the same distribution. Therefore, we want to compare the two distributions and establish whether their difference is sufficient to discard H_0 .

The comparison between two observed distributions is performed in a more robust way from the empirical cumulative distribution. If one has a random sample of a variable X , $\{x_i\}_{i=1}^n$, where the values are sorted, the empirical cumulative distribution function (ECDF) is:

$$F_e[x] = \frac{|\{X/X < x\}|}{|\{X\}|} \quad (19)$$

This function is used in the Kolmogorov-Smirnov test, where the statistic D is defined as the maximum difference between the two cumulative distributions. Following this idea, to describe the difference between an observed distribution and a control distribution $F_c(x)$, we consider:

$$\Delta(x) = F_e[X](x) - F_c(x) \quad (20)$$

In the case of the distribution of the λ parameter, we know that under H_0 it has to be uniform between 0 and $\pi/2$, which yields $F_c(x) = 2x/\pi$. Then,

$$\Delta(x) = F_e[X](x) - \frac{2x}{\pi} \quad (21)$$

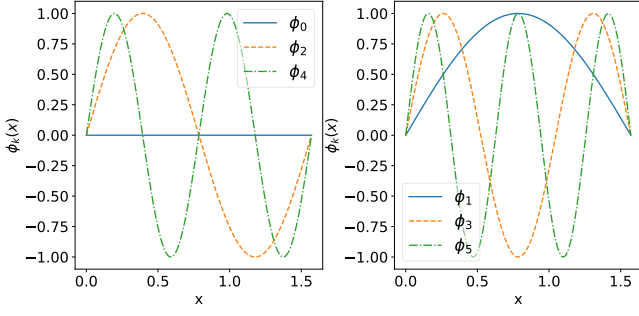


Figure 4. First even (left panel) and odd (right panel) elements of the base of functions $\phi(x)$.

This function, by definition, begins and ends in zero, i.e.,

$$\Delta(x) = 0 \text{ for } x = 0; x = \pi/2.$$

In order to represent the function $\Delta(x)$ one can use a base of orthogonal functions. If $f(x)$ is a continuous function in the interval $[0, \pi/2]$, then it can be written as:

$$\Delta(x) = \sum_{k=0}^{\infty} a_k \phi_k(x) \quad (22)$$

where

$$\phi_k(x) = \sin(2kx).$$

The first elements of this base can be seen in Figure 4. And by defining orthogonality of the base we can finally write our orthonormal system as:

$$\phi_k(x) = \frac{4}{\pi} \sin(2kx). \quad (23)$$

That said, in the case of data sets, there is no continuous function, but a discrete sampling. Fourier analysis allows for an expansion in terms of a finite sum of sines in the case of a set of equally distanced points. If the points are not equally distanced we must resort to alternative strategies. One might make a binning, but this alternates the information available. Another option is to fit a model to the observed points. For example, we might take as a model:

$$\Phi(x) = \sum_{k=0}^M a_k \phi_k(x), \quad (24)$$

which is a model that is linear in the parameters that we want to calculate: the coefficients a_k .

Even though this method is enough to represent a distribution, it is not useful as a statistic to characterize the parameters of the distribution. However, it is possible to take advantage of the conditions of symmetry to characterize the distribution.

For example, if we assume symmetry, i.e. consider only the cosines of acute angles, we expect the even parameters to be zero. Ec. 24 is modelling the difference between a cumulative distribution of the cosine distribution calculated from data and a control function, so it is expected that when studying the effects of preferential alignments in a vector population this function will take values that are either mostly positive or negative, indicating a net alignment perpendicular or parallel to the preferred direction, respectively. Therefore, in the case of symmetry where one is interested in the acute angle λ , it is

noticeable from the right panel of Fig. 4 that the first term, represented in a blue curve, will be the dominant term in Ec. 24 and is proportional to the parameter a_1 ; the area under the curve, whether positive or negative, will represent a net alignment in the perpendicular or parallel direction, respectively.

The coefficients a_k can be obtained from ordinary least squares from data. Given Ec. 24 and the residuals of the ECDF of the data $y_i = i/n - i$ of size n , we intend to minimize:

$$\chi^2 = \sum_{i=1}^n \left(\frac{y_i}{\sigma_i} - \frac{1}{\sigma_i} \Phi(x) \right)^2 = |\mathbf{A}a_k - \mathbf{B}|, \quad (25)$$

where \mathbf{A} and \mathbf{B} are in matrix notation and x is the data we want to fit, i.e. the cosines calculated from the sample. So by minimizing this expression it follows that $\mathbf{A}a_k$ equals \mathbf{B} (Hastie et al. 2001, Chapter 3), so:

$$a_k = \frac{(\mathbf{A}^T \mathbf{B})_k}{(\mathbf{A}^T \mathbf{A})_k}. \quad (26)$$

The base of functions is the set of all harmonic functions $\phi_k(x) = \sin(k\pi x)$, which we truncate at $k=4$, and assuming $\sigma_i = 1$, we have:

$$(\mathbf{A}^T \mathbf{A})_k = \sum_{i=1}^n \phi_k^2(x_i). \quad (27)$$

On the other hand, we have $(\mathbf{A}^T \mathbf{B})_k = \sum_{i=1}^n \phi_k(x_i) y_i$, so, after replacing we have:

$$a_k = \frac{\sum_{i=1}^n \phi_k(x_i) y_i}{\sum_{i=1}^n \phi_k^2(x_i)}. \quad (28)$$

In our case of study we have $x = \cos(\lambda)$, and assuming sorted data: $y_i = i/n - i$, we finally arrive at an analytical expression for the OLS coefficients:

$$a_k = \frac{\sum_{i=1}^n i \sin[k\pi \cos(\lambda_i)]}{\sum_{i=1}^n \sin^2[k\pi \cos(\lambda_i)]} \left(\frac{n-1}{n} \right). \quad (29)$$

This is an expression that directly relates the parameters to the data. As previously indicated, the coefficient a_1 (Ec. 30) gives the first order approximation for the residual function (Ec. 24):

$$a_1 = \frac{\sum_{i=1}^n i \sin[\pi \cos(\lambda_i)]}{\sum_{i=1}^n \sin^2[\pi \cos(\lambda_i)]} \left(\frac{n-1}{n} \right). \quad (30)$$

If the data consists of a vector population with a net alignment perpendicular to the preferred direction, Ec. 24 will resemble the blue curve of the right panel of Fig. 4 with the coefficient a_1 taking positive values. This is due to the data presenting an excess in the lower values of cosines and, as a consequence, the ECDF taking values larger than the control function so that the residues are positive. On the other hand, if the data presents a net alignment in the parallel direction, a_1 will take negative values.

3 APPLICATION TO SYNTHETIC DATA AND COMPARISON WITH OTHER METHODS

To test the efficiency of the methods presented above with regards to usual methods, such as the average cosine, we apply them to 3 sets

of synthetic data. These data are generated sorting random points on the surface of a 3-dimensional ellipsoid with axis a , b , and c , with various eccentricities defined in the usual manner: $e^2 = 1 - c^2/a^2$, where $c < a$ and $a = b$. It is worth noting that by varying the c axis, we are defining this vertical z direction as the preferential direction for spherical symmetry. We chose to establish three different eccentricities to test the methods: $e^2 = 0.6$, 0.4 , and 0 , going from elongated to isotropic, respectively. In this manner, we are simulating a population of vectors with no preferential orientation for the isotropic case, to one with a strong alignment trend for the largest eccentricity.

First, we studied the stability of the estimators with various sample sizes. Figure 5 shows the mean and standard deviation calculated for 50 random realizations of samples of size N_{ran} . For a sample size of over a few hundred the estimation of the parameters a_1 and η is reliable, and for sample sizes of over $\sim 10^4$ the relative error is sufficiently small to distinguish between little variations in eccentricities. This is promising in the sense that, for present and future large scale surveys with large samples, even a small effect of alignment would be detectable with these methods.

In Fig. 6 we test the statistical significance of the η parameter when compared to the average of cosines, $\langle \cos(\lambda) \rangle$. As explained above, we generate random points along the surface of three ellipsoids with eccentricity values of 0.6 , 0.4 , and 0 . These populations yield the three cosine histograms showed in panel a), where we include the mean values along with the standard deviation of the distribution. We note that the standard deviation of the cosine distribution is of the order of its mean.

Panel c) shows the logarithm of the β parameter defined as the ratio of perpendicular and parallel components of the vector, the parallel direction being that of the preferential direction for spherical symmetry. The mean value of β is similar for every population, so this is not an ideal parameter to study. However, the cumulative number of vectors that have a larger perpendicular component, a.k.a. the η parameter, is noticeable.

In order to account for sample variance, we perform bootstrap sampling of our observable β to obtain a distribution of η from which we can define a mean and a standard deviation. Such bootstrap distributions for each eccentricity are shown in panel d). We include the mean value and standard deviation of the distributions, as well as that of the isotropic case in grey colour, which has been derived theoretically in Sec. 2.3.

It is noticeable that the η distributions have larger standard deviation for larger mean values. This is a consequence of the definition of β , where values for larger parallel components are limited between 0 and 1 , while values for larger perpendicular values have no theoretical upper limit. If one were to define β in the inverse manner, the same divergent behaviour happens for larger parallel alignments. This is a feature of the parameter to keep in mind. However, while the upper limit is infinite in theory, it is not in practice. On one hand one would have to find vectors of infinite norm for this to be a problem. Furthermore, alignment corresponding to an eccentricity of 0.6 is unlikely to be found in observables such as galactic orientations, much less higher values of eccentricity. In other words, we are testing these parameters in the limits of practical situations.

To finally assess the efficacy and significance of the η parameter, we repeat the above procedure by generating random points along the surface of the different ellipsoids with 50 random seeds, therefore yielding 50 values of η and average cosines. Furthermore, to study the

statistical significance with respect to random behaviours we define the variable ζ as:

$$\zeta_X = \frac{X - \bar{X}_{\text{ran}}}{\sigma_{X,\text{ran}}}, \quad (31)$$

where X is the random variable we want to test, in this case: η and $\langle \cos(\lambda) \rangle$. Given its definition, the ζ_X variable contains information, not only about how much does the variable X deviates from isotropic behaviour, but also how statistically significant this deviance is. The isotropic values for the mean and standard deviation for η have been theoretically derived in Sec. 2.3. The mean and standard deviation for the cosines in the isotropic case can be calculated as those of a uniform distribution. The probability distribution function of a uniform distribution is:

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases} \quad (32)$$

where, in the case of the cosines, $a=0$ and $b=1$. So the mean and standard deviation would be:

$$E(\cos_{\text{ran}}) = 0.5, \quad (33)$$

and

$$\sigma_{\cos,\text{ran}} = \sqrt{\frac{1}{12}} \approx 0.289. \quad (34)$$

Panel e) shows the computation of ζ_η and ζ_{\cos} . We observe that the statistical significance for deviations from isotropic behaviour as measured with η is much higher than with the cosines. For an eccentricity of $e^2 = 0.6$ we have a significance of around 12.5 for η and 0.3 for the cosines.

In Fig. 7 we perform an equivalent analysis for the OLS coefficients method using the same synthetic data, as can be seen by comparing the cosine distributions in panels a) of both Fig. 7 and 6. For this method we first calculate the ECDF of the cosines (panel c) of the data corresponding to the three cases of varying anisotropy. The residues are calculated by subtracting from the ECDF of the data the one corresponding to an isotropic distribution which is the straight line of $\text{ECDF}(\cos(\lambda)) = \cos(\lambda)$. We generate the data and perform this calculation 50 times with different random seeds, as shown in panel d). We fit each of this curves and plot the mean of the fit with solid lines and their corresponding 3σ with the shadowed bands in panel f).

The linear regression for each curve yields a set of coefficients a_k , where the one that determines the basic shape of the fit is a_1 (see section 2.4). We perform bootstrap resampling of the data in order to estimate the mean and standard deviation of this coefficient. Panel e) shows the bootstrap distribution of the a_1 coefficient corresponding to the cosine distributions shown in the same colors, where the dotted vertical lines correspond to the mean. Furthermore, in this panel we indicate in text the mean values along with the standard deviation of the distribution. It is readily noticeable that the mean of this coefficient is more robustly determined than the mean of the cosines.

Finally, for each of the random realizations we plot the normalized parameters ζ_{a_1} and $\zeta_{\langle \cos \rangle}$, where we find that the OLS coefficient detects alignment with higher statistical significance.

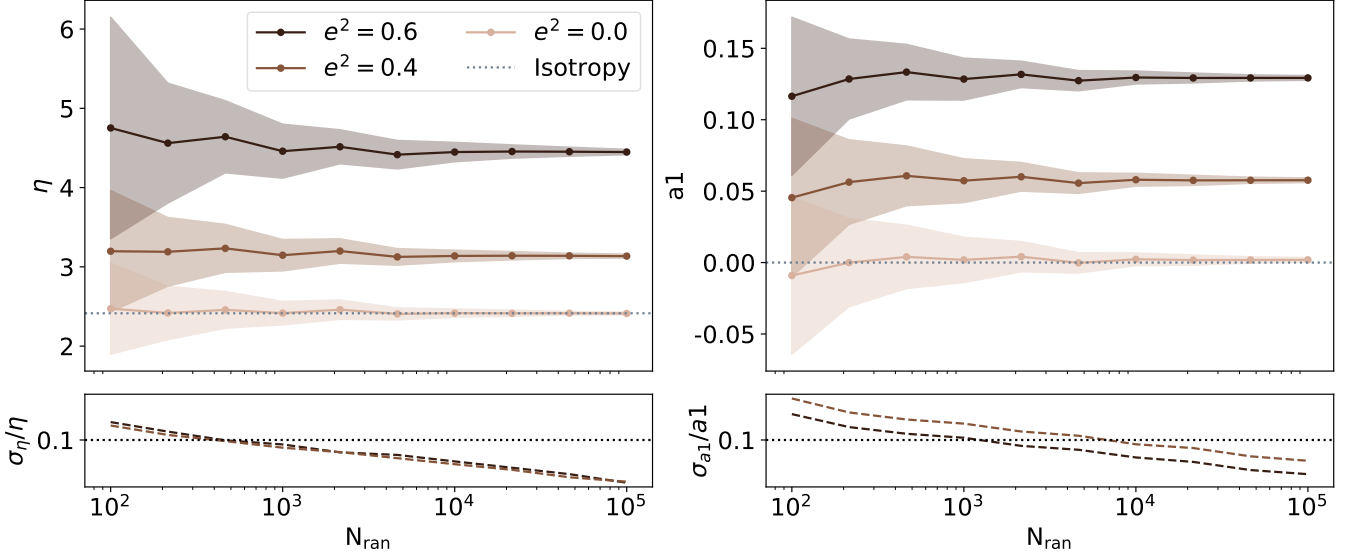


Figure 5. Stability of the methods with respect to the sample size. We applied the methods to 50 random realizations of synthetic data with alignments corresponding to three values of increasing eccentricities: 0, 0.4, and 0.6. The mean and standard deviations, represented by the solid lines and shadowed regions respectively, were calculated with these 50 independent results. We find that the mean of the parameters η and a_1 is stable even with a sample size of a few hundred. A relative error of 10% is achieved with a sample size of $\sim 10^3$ for the η parameter. The same sample size for the same relative error is achieved when applying the OLS coefficient method to the data with the larger eccentricity.

4 SUMMARY

In this paper we present two methods to detect and quantify alignment signals and test their statistical robustness. The first method uses the deviance of the relative fraction of vector components in the plane perpendicular to a reference direction with respect to an isotropic distribution. We have derived the first and second moments of the distribution of the resulting estimator, η , and can thus reliably assess its statistical significance. The second method is based on a fit over the residuals of the ECDF of the data with respect to the one expected for a uniform distribution. The fit uses a small set of harmonic orthogonal functions and does not rely on any binning scheme. The amplitude of the fit, i.e. the amplitude of the alignment signal, can be described by the first odd OLS parameter, a_1 .

For the first method, we derive the distribution of the test statistic β , defined as the ratio between the perpendicular and parallel vector components (Eq. 4), under the null hypothesis. We find that the probability distribution $f_B(\beta)$ is a pathological distribution where the moments are undefined, and as such, is not a robust statistic. However, using this statistic we consider the fraction of its values that are greater than one, given that when the perpendicular and parallel components are equal there is no preference in any of the two directions and so $\beta = 1$ can be taken as a critical value. Therefore, we define this ratio as the parameter η in Eq. 15. We find that, for β defined as in Eq. 4, the parameter η has an expectation value and variance given by $\eta_0 \approx 2.4142$ and $\text{Var}(\eta) \approx 28.1421/N$, respectively. The gaussian behaviour of this parameter allows for the first two moments to be sufficient to describe its distribution. The advantage of knowing the theoretical distribution of the parameters is the ability to accurately determine the statistical significance of any signal detected without investing computation resources and time into Monte Carlo simulations.

The second method of alignment analysis we presented yields OLS coefficients of a fit of the residues of the ECDF of the cosines of the data with respect to a random sample. The first odd coefficient,

a_1 , of the harmonic expansion of the residual function is sufficient to characterize the amplitude of the alignment signal, with zero being consistent with isotropic orientations. Positive values of a_1 indicate perpendicular alignment while negative values indicate parallel alignment with respect to the preferred direction.

We have compared these methods with others commonly used in the literature, mainly based on the average of the cosine distribution and using Monte Carlo simulations. This comparison was achieved by simulating a population of vectors with three different degrees of alignment (from no alignment, to intermediate, to greatly aligned), and testing how well the alignment signal was recovered by both the new and traditional parameters.

We find that the proposed statistics allow the detection of alignment signals with a larger significance. For a deviation of approximately 0.25σ from an isotropic distribution of cosines, we obtain a significance of 10– and 12 σ for the OLS coefficient a_1 and the η parameter respectively. In a forthcoming paper (Dávila-Kurbán et al., in prep) we apply the first method presented in this paper, i.e. the fraction of vectors parameter η , to data in a cosmological simulation.

We have assessed the effects of uncertainties in the measurement of parallel and perpendicular vector components. To that end, we modelled observational errors by introducing gaussian noise into the components. The standard deviation was chosen to be 10 per cent of the mean vector norm. There is a linear relation between the "real" parameters, η and a_1 , obtained with the raw synthetic data, and the "observed" parameters that account for the introduced mock observational error. We performed 100 realizations of the calculations with and without mock errors. The y–intercept b corresponding to the maximum divergence from the real values can be expressed, for a sample size N of 1000 and 5000 respectively, as $\Delta b_\eta = .7$ and $.3$ for the first method, and $\Delta b_{a_1} = .03$ and $.01$ for the second method. For reference, the range of values of the two parameters for the three eccentricities tested are: $2.0 < \eta < 5.5$ and $-.05 < a_1 < 0.16$ for $N=1000$; and $2.2 < \eta < 5.0$ and $-.02 < a_1 < 0.15$ for $N=5000$.

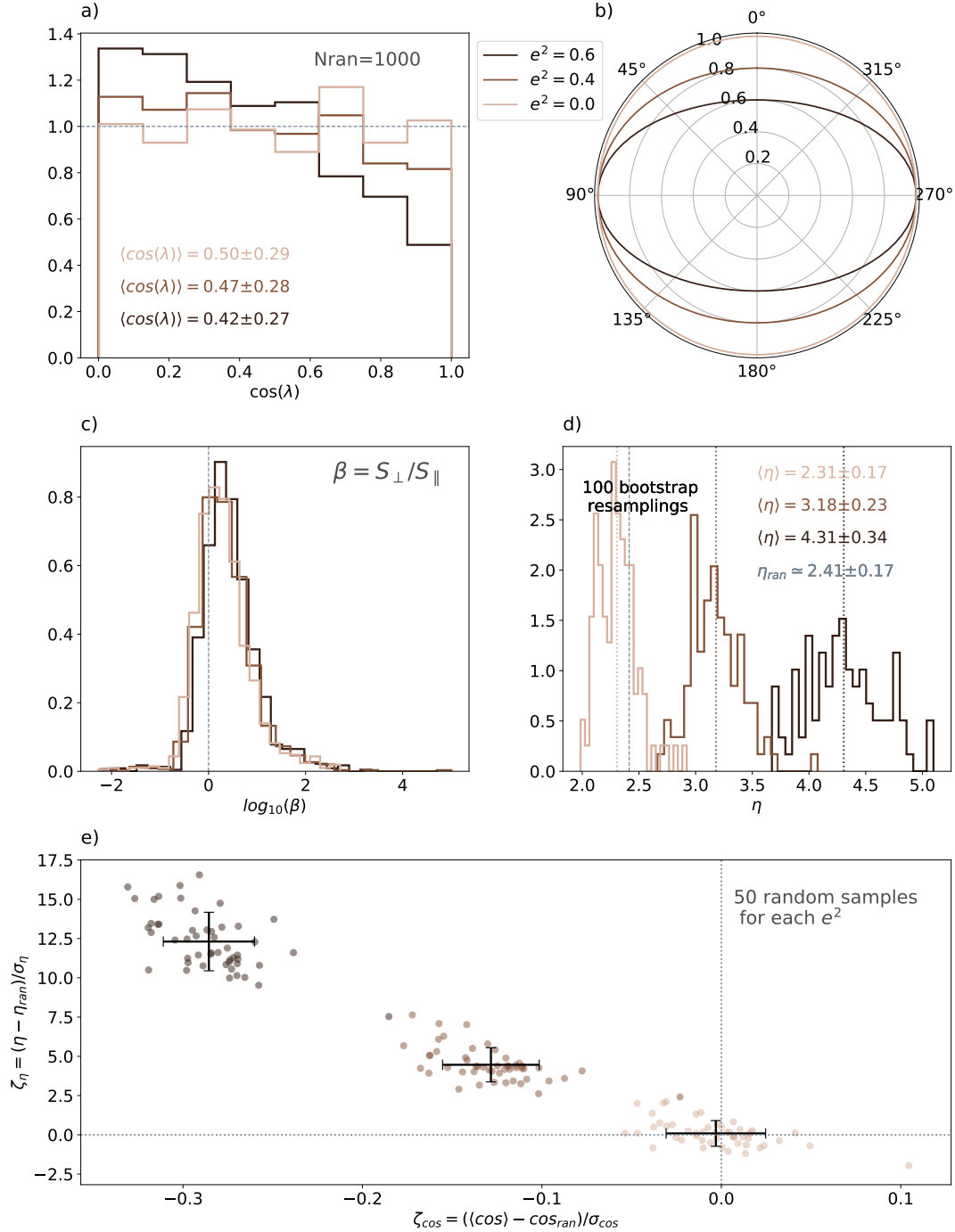


Figure 6. Representation of the method starting from a distribution of alignments of a vector population arriving to the “fraction of vectors” parameter η , along with its significance compared to the average of cosines. Panels a) and b) show the histograms of cosines of angles, which is the usual manner of studying alignments, corresponding to directions scattered along the surface of ellipsoids with eccentricities 0.6, 0.4, and 0, i.e. from strong alignment to random behaviour. The first panel includes the mean and standard deviation of the distribution of cosines. The method then consists of calculating the parameter β as the ratio of perpendicular to parallel components, whose distributions for the three alignments are shown on panel c), and then obtaining η : the number of vectors with $\beta > 1$. Sample variance is taken into account by bootstrap resampling the data, and thus a mean and standard deviation for η can be estimated (panel d). To assess the stability of the parameter we perform this calculation several times by varying the random seed of the initial samples and thus obtaining several estimations of the bootstrap mean of η . And finally, in order to quantify the significance of the estimated alignment signal, we define a variable ζ_{η} normalized with respect to the behaviour of the estimator η under the null hypothesis. Panel e) shows ζ_{η} against $\zeta_{(cos)}$, where we find that η can detect alignment signal with higher statistical significance.

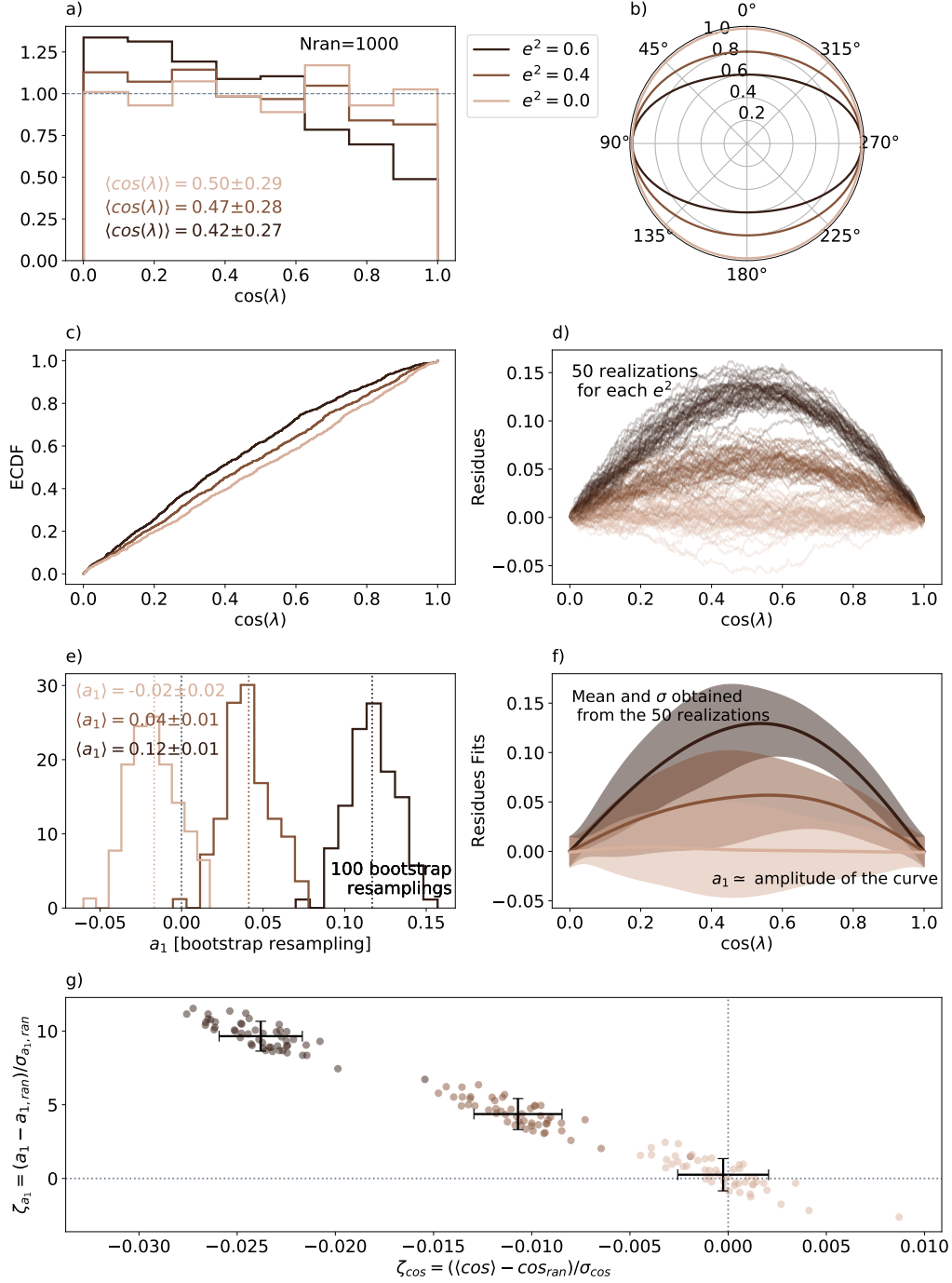


Figure 7. Representation of the method starting from a distribution of alignments of a vector population and deriving the the OLS coefficient a_1 , along with its significance compared to the average of cosines. Panels a) and b) show the histograms of cosines of angles, which is the usual manner of studying alignments, corresponding to directions scattered along the surface of ellipsoids with eccentricities 0.6, 0.4, and 0, i.e. from strong alignment to random behaviour. The first panel includes the mean and standard deviation of the distribution of cosines. The method then consists of calculating the residues of the ECDF of the cosine distribution, shown in panel c), and that of an isotropic behaviour. To assess the stability and significance of the parameter we perform this calculation several times by varying the random seed of the initial samples and thus obtaining the residual curves shown in panel d). Fits are performed over these curves. The amplitudes of said fits are characterized by the parameter a_1 , whose distributions along with their mean and standard deviation are shown in panel e). A representation of the fits is plotted in panel f), with a solid line representing the mean and a colored band showing their standard deviation. And finally, in order to quantify the significance of the estimated alignment signal, we define a variable ζ_{a_1} normalized with respect to the behaviour of the estimator a_1 under the null hypothesis. Panel g) shows ζ_{a_1} against ζ_{cos} , where we find that the coefficient a_1 can detect alignment signal with higher statistical significance.

This can be used to roughly estimate the size of samples required to achieve the detection of small signal in the data. With these tools, large forthcoming surveys and simulations can provide new insights on small amplitude signals of alignments unseen in current surveys with lower number of galaxies.

ACKNOWLEDGEMENTS

This work was partially supported by the Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET, Argentina), the Secretaría de Ciencia y Tecnología, Universidad Nacional de Córdoba, Argentina, and the Agencia Nacional de Promoción de la Investigación, el Desarrollo Tecnológico y la Innovación, Ministerio de Ciencia, Tecnología e Innovación, Argentina. This research has made use of NASA's Astrophysics Data System. Visualizations made use of python packages and inkscape software.

DATA AVAILABILITY

The data underlying in this article are available on request to the corresponding author.

REFERENCES

- Aragón-Calvo M. A., Yang L. F., 2014, *MNRAS*, **440**, L46
- Aragón-Calvo M. A., van de Weygaert R., Jones B. J. T., van der Hulst J. M., 2007, *ApJ*, **655**, L5
- Aryal B., Saurer W., 2000, *A&A*, **364**, L97
- Betancort-Rijo J. E., Trujillo I., 2009, arXiv e-prints, p. arXiv:0912.1051
- Brainerd T. G., 2005, *ApJ*, **628**, L101
- Brunino R., Trujillo I., Pearce F. R., Thomas P. A., 2007, *MNRAS*, **375**, 184
- Codis S., Jindal A., Chisari N. E., Vibert D., Dubois Y., Pichon C., Devriendt J., 2018, *Monthly Notices of the Royal Astronomical Society*, Volume 481, Issue 4, p.4753-4774, 481, 4753
- Colless M., et al., 2001, *MNRAS*, **328**, 1039
- Cuesta A. J., Betancort-Rijo J. E., Gottlöber S., Patiri S. G., Yepes G., Prada F., 2008, *MNRAS*, **385**, 867
- Doroshkevich A. G., 1970, *Astrophysics*, **6**, 320
- Duris F., Gazdarica J., Gazdaricova I., Strieskova L., Budis J., Turna J., Szemes T., 2018, *Journal of Statistical Distributions and Applications*, **5**
- Fall S. M., Efstathiou G., 1980, *Monthly Notices of the Royal Astronomical Society*, **193**, 189
- Flin P., Godłowski W., 1986, *MNRAS*, **222**, 525
- Forero-Romero J. E., Contreras S., Padilla N., 2014, *Monthly Notices of the Royal Astronomical Society*, **443**, 1090
- Gillespie D. T., 1983, *American Journal of Physics*, **51**, 520
- Godłowski W., 1993, *MNRAS*, **265**, 874
- Godłowski W., 1994, *MNRAS*, **271**, 19
- Godłowski W., Ostrowski M., 1999, *MNRAS*, **303**, 50
- Hahn O., Carollo C. M., Porciani C., Dekel A., 2007, *MNRAS*, **381**, 41
- Hastie T., Tibshirani R., Friedman J., 2001, *The Elements of Statistical Learning*. Springer Series in Statistics, Springer New York Inc., New York, NY, USA
- Hu F. X., Wu G. X., Su H. J., Liu Y. Z., 1995, *A&A*, **302**, 45
- Hu F. X., Yuan Q. R., Su H. J., Wu G. X., Liu Y. Z., 1998, *ApJ*, **495**, 179
- Jaaniste J., Saar E., 1978, in Longair M. S., Einasto J., eds. Vol. 79, *Large Scale Structures in the Universe*. p. 448
- Joachim B., et al., 2015, *Space Science Reviews*, Volume 193, Issue 1-4, pp. 1-65, 193, 1
- Kashikawa N., Okamura S., 1992, *PASJ*, **44**, 493
- Kiessling A., et al., 2015, *Space Sci. Rev.*, **193**, 67
- Koopman P. A. R., 1984, *Biometrics*, **40**, 513
- Kraljic K., Dave R., Pichon C., 2019, *Monthly Notices of the Royal Astronomical Society*, **493**, 362
- Lee J., 2004, *ApJ*, **614**, L1
- Lee J., Pen U.-L., 2000, *The Astrophysical Journal*, Volume 532, Issue 1, pp. L5-L8., 532, L5
- Libeskind N. I., Hoffman Y., Forero-Romero J., Gottlöber S., Knebe A., Steinmetz M., Klypin A., 2013, *MNRAS*, **428**, 2489
- Mo H., van den Bosch F. C., White S., 2010, *Galaxy Formation and Evolution*
- Panko E., Piwowarska P., Godłowska J., Godłowski W., Flin P., 2013, *Astrophysics*, **56**, 322
- Peebles P. J. E., 1969, *ApJ*, **155**, 393
- R.M. P., Bonett D. G., 2008, *Statistics in Medicine*, **27**, 5497
- Riebe K., et al., 2011, arXiv e-prints, p. arXiv:1109.0003
- Slosar A., White M., 2009, *J. Cosmology Astropart. Phys.*, **2009**, 009
- Trujillo I., Carretero C., Patiri S. G., 2006, *ApJ*, **640**, L111
- Varela J., Betancort-Rijo J., Trujillo I., Ricciardelli E., 2012, *ApJ*, **744**, 82
- White S. D. M., 1984, *Astrophysical Journal*, Part 1 (ISSN 0004-637X), vol. **286**, Nov. 1, 1984, p. 38-41. *NASA-supported research.*, 286, 38
- Wu G. X., Hu F. X., Su H. J., Liu Y. Z., 1997, *A&A*, **323**, 317
- Yang X., van den Bosch F. C., Mo H. J., Mao S., Kang X., Weinmann S. M., Guo Y., Jing Y. P., 2006, *MNRAS*, **369**, 1293
- Yuan Q. R., Hu F. X., Su H. J., Huang K. L., 1997, *AJ*, **114**, 1308
- Zhang Y., Yang X., Wang H., Wang L., Luo W., Mo H. J., van den Bosch F. C., 2015, *ApJ*, **798**, 17

APPENDIX A: VARIANCE OF $\hat{\eta}$

We want to calculate the expectation value for the ratio $Q = X/Y$. From the definition of the expectation value it stems that it is not possible to derive a simple expression for $Q[R]$. Another reason that prevents from analytically solving the distribution of Q is that the denominator can be zero. A way of solving this problem is to rewrite the function in a way that avoids having a singularity. It is possible to make such an approximation from developing the Taylor series of $\bar{Q}(X, Y) = X/Y$ around $(X, Y) = (\mu_X, \mu_Y)$, i.e. $\bar{Q} = Q + R$:

$$\begin{aligned}\bar{Q}(X, Y) &= \bar{Q}(\mu_X, \mu_Y) + \frac{\partial \bar{Q}}{\partial X}(\mu_X, \mu_Y)(X - \mu_X) \\ &\quad + \frac{\partial \bar{Q}}{\partial Y}(\mu_X, \mu_Y)(Y - \mu_Y) + R,\end{aligned}$$

where R is the order 2 error given by the Taylor theorem (Duris et al. 2018; Koopman 1984; R.M. & Bonett 2008). Then, we can estimate the expectation value of \bar{Q} , which is approximately $E[\bar{Q}] \approx E[Q]$, where:

$$\begin{aligned}E[Q] &= E\left[Q(\mu_X, \mu_Y) + \frac{\partial Q}{\partial X}(\mu_X, \mu_Y)(X - \mu_X) + \frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)(Y - \mu_Y)\right] \\ &= E\left[Q(\mu_X, \mu_Y)\right] + E\left[\frac{\partial Q}{\partial X}(\mu_X, \mu_Y)(X - \mu_X)\right] + E\left[\frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)(Y - \mu_Y)\right] \\ &= Q(\mu_X, \mu_Y) + \frac{\partial Q}{\partial X}(\mu_X, \mu_Y)E\left[X - \mu_X\right] + \frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)E\left[Y - \mu_Y\right] \\ &= Q(\mu_X, \mu_Y)\end{aligned}$$

Then, to a first order approximation and considering $Q = \eta$, $X = n$, $Y = N - n$:

$$E(\hat{\eta}) \approx \frac{Np}{N - Np}$$

and

$$\frac{Np}{N - Np} = \frac{p}{1 - p} = \frac{\frac{1}{\sqrt{2}}}{1 - \frac{1}{\sqrt{2}}} = \frac{P(\beta > 1)}{P(\beta < 1)} = \eta_0$$

therefore

$$E(\hat{\eta}) \approx \eta_0$$

The second order moment of $\hat{\eta}$ can be obtained calculating the variance of the first order approximation of $Q(X, Y)$:

$$\begin{aligned}Var(Q) &= Var\left[Q(\mu_X, \mu_Y) + \frac{\partial Q}{\partial X}(\mu_X, \mu_Y)(X - \mu_X) + \frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)(Y - \mu_Y)\right] \\ &= \left(\frac{\partial Q}{\partial X}(\mu_X, \mu_Y)\right)^2 Var[X - \mu_X] + \left(\frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)\right)^2 Var[Y - \mu_Y] + \\ &\quad + \left(2\frac{\partial Q}{\partial X}(\mu_X, \mu_Y)\frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)\right) Cov[X, Y] \\ &= \left(\frac{\partial Q}{\partial X}(\mu_X, \mu_Y)\right)^2 \sigma_X^2 + \left(\frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)\right)^2 \sigma_Y^2 + \left(2\frac{\partial Q}{\partial X}(\mu_X, \mu_Y)\frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)\right) Cov[X, Y] \\ &= \left[\left(\frac{\partial Q}{\partial X}(\mu_X, \mu_Y)\right)^2 + \left(\frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)\right)^2\right] \sigma_X^2 + \left(2\frac{\partial Q}{\partial X}(\mu_X, \mu_Y)\frac{\partial Q}{\partial Y}(\mu_X, \mu_Y)\right) Cov[X, Y]\end{aligned}$$

To evaluate the derivatives, we have:

$$\begin{aligned}\frac{\partial \eta}{\partial X}(\mu_X, \mu_Y) &= \frac{\partial(X/Y)}{\partial X}(\mu_X, \mu_Y) = \frac{1}{Y} \Big|_{(X,Y)=(\mu_X,\mu_Y)} \\ &= \frac{1}{\mu_Y} = \frac{1}{N(1-p)} = \\ &= \frac{1}{Nq}\end{aligned}\tag{A1}$$

and

$$\begin{aligned}\frac{\partial \eta}{\partial Y}(\mu_X, \mu_Y) &= \frac{\partial X/Y}{\partial Y}(\mu_X, \mu_Y) = -\frac{X}{Y^2} \Big|_{(X,Y)=(\mu_X,\mu_Y)} \\ &= -\frac{\mu_X}{\mu_Y^2} = -\frac{Np}{N^2(1-p)^2} \\ &= -\frac{p}{Nq^2}\end{aligned}\tag{A2}$$

The covariance between X and Y is given by:

$$\begin{aligned}
 \text{Cov}(X, Y) &= \frac{1}{N} \sum_{i=1}^N (x_i - E(X))(y_i - E(Y)) \\
 &= \frac{1}{N} \sum_{i=1}^N (x_i - \mu_X) \left((N - x_i) - (N - \mu_X) \right) \\
 &= \frac{1}{N} \sum_{i=1}^N (x_i - \mu_X) (-x_i + \mu_X) \\
 &= -\frac{1}{N} \sum_{i=1}^N (x_i - \mu_X)^2 \\
 &= -\sigma_X^2
 \end{aligned} \tag{A3}$$

Then, using the expressions A1, A2 y A3

$$\begin{aligned}
 \text{Var}(\eta) &= \left[\left(\frac{1}{Nq} \right)^2 + \left(\frac{p}{Nq^2} \right)^2 \right] Npq + 2 \frac{1}{Nq} \frac{p}{Nq^2} Npq \\
 &= \frac{1}{N} \left(\left(\frac{1}{q^2} + \frac{p^2}{q^4} \right) pq + 2 \frac{p^2}{q^2} \right) \\
 &= \frac{1}{N} \left[\frac{p}{q} + 2 \left(\frac{p}{q} \right)^2 + \left(\frac{p}{q} \right)^3 \right]
 \end{aligned}$$

Evaluating in $p = 1/\sqrt{2}$ we have that

$$\frac{p}{q} = \frac{1/\sqrt{2}}{1 - 1/\sqrt{2}} = \sqrt{2} + 1$$

Then,

$$\text{Var}(\eta) \approx \frac{28.14214}{N} \tag{A4}$$

It can be seen in Figure 3 that the theoretical variance of the expression A4 and values of the variance calculated from 10 or 100 samples of values of η , obtained in equivalent simulated samples of β with the method of the binomial distribution.