

# Science-driven 3D data compression

David Alonso<sup>1</sup>

<sup>1</sup>*Oxford Astrophysics, Department of Physics, Keble Road, Oxford, OX1 3RH, UK*

Photometric redshift surveys map the distribution of matter in the universe through the positions and shapes of galaxies with poorly resolved measurements of their radial positions. While a tomographic analysis can be used to recover some of the large-scale radial modes present in the data, this approach suffers from a number of practical shortcomings, and the criteria to decide on a particular binning scheme are commonly blind to the ultimate science goals. We present a method designed to separate and compress the data into a small number of uncorrelated radial modes, circumventing the main problems of standard tomographic analyses. The method is based on the Karhunen-Loève transform, and is connected to other 3D data compression bases advocated in the literature, such as the Fourier-Bessel or Fourier-Laguerre bases. We apply this method to the specific cases of constraining dark energy with weak lensing and primordial non-Gaussianity with galaxy clustering. In the case of weak lensing we show that the method is able to compress the vast majority of the information into a single mode, considerably simplifying the analysis with respect to a traditional tomographic approach.

## I. INTRODUCTION

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

## II. METHOD

### A. The Karhunen-Loève transform

The idea behind the Karhunen-Loève transform, as developed within the field of cosmological data analysis in e.g. CITES, is to compress a given data vector into a small set of modes containing most of the useful information on a particular parameter (or set of parameters). Let  $\mathbf{a}$  be a data vector of dimension  $N_s$ , and let  $\theta$  be a particular parameter we want to measure. Under the assumption that  $\mathbf{x}$  is Gaussianly distributed with mean 0 and covariance  $\mathbf{X}$ , a set of linear combinations  $y_p \equiv \mathbf{e}_p^\dagger \mathbf{x}$  can be found such that the  $y_p$  are white and uncorrelated ( $\langle y_p y_q^* \rangle = \delta_{pq}$ ) and such that the first  $m < N_s$  combinations contain most of the information about  $\theta$ . This is done by solving the generalized eigenvalue problem CITE:

$$\partial_\theta \mathbf{C} \mathbf{e}_p = \lambda_p \mathbf{C} \mathbf{e}_p, \quad (1)$$

where  $\partial_\theta \equiv \partial/\partial\theta$ .

#### 1. The K-L transform for the signal-to-noise

Let us decompose the data vector  $\mathbf{x}$  into uncorrelated signal and noise components  $\mathbf{x} = \mathbf{s} + \mathbf{n}$  where, in this context, the signal is the part of the data containing any information of cosmological interest, and the noise is any contaminant preventing us to access it. In this particular case, the data covariance matrix can be split into their independent contributions  $\mathbf{X} = \mathbf{S} + \mathbf{N}$ .

The K-L transform has traditionally been used to design an eigenbasis that maximizes the overall signal-to-noise ratio (e.g. CITES). This can be done by defining a fictitious parameter  $\alpha$  multiplying the signal part of the data with fiducial value  $\alpha = 1$  (i.e.  $\mathbf{x} = \alpha \mathbf{s} + \mathbf{n}$ ). In this case, the eigenvalue equation (Eq. 1) takes the form:

$$(\mathbf{S} + \mathbf{N})\mathbf{e}_p = \lambda_p \mathbf{N} \mathbf{e}_p, \quad (2)$$

where we have redefined  $2/(2 - \lambda_p) \rightarrow \lambda_p$ . This can be cast into a standard eigenvalue equation using the Cholesky decomposition of the noise covariance matrix  $\mathbf{N} = \mathbf{L} \mathbf{L}^\dagger$ :

$$[\mathbf{L}^{-1} \mathbf{C} (\mathbf{L}^{-1})^\dagger] \tilde{\mathbf{e}}_p = \lambda_p \tilde{\mathbf{e}}_p, \quad (3)$$

where  $\tilde{\mathbf{e}}_p \equiv \mathbf{L}^\dagger \mathbf{e}_p$ .

At this point it is worth noting that the generalized eigenvalue problem in Eq. 2 can be understood as the problem diagonalizing  $\mathbf{C}$  under a non-standard dot product  $\circ$  given by the inverse noise covariance matrix (i.e.  $\mathbf{a} \circ \mathbf{b} \equiv \mathbf{a}^\dagger \mathbf{N}^{-1} \mathbf{b}$ ). Under this dot product, an eigenbasis  $\mathbf{F} \equiv (\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_{N_s})$  can be found such that  $\mathbf{F}$  is orthonormal  $\mathbf{F} \circ \mathbf{F} = \mathbf{I}$ , and the covariance of the transformed data vector  $\mathbf{y} \equiv \mathbf{F} \circ \mathbf{x}$  is diagonal:

$$\langle \mathbf{y} \mathbf{y}^\dagger \rangle = \mathbf{F}^\dagger \mathbf{N}^{-1} \mathbf{C} \mathbf{N}^{-1} \mathbf{F} = \mathbf{\Lambda} \equiv \text{diag}(\lambda_1, \dots, \lambda_{N_s}). \quad (4)$$

Using the orthonormality of  $\mathbf{F}$ , this can be cast into the same form as Eq. 3, where  $\mathbf{f}_p = \mathbf{L} \tilde{\mathbf{e}}_p = \mathbf{N} \mathbf{e}_p$ .

Finally, note that, because both  $\mathbf{S}$  and  $\mathbf{N}$  are positive-definite matrices, their eigenvalues will also be positive. Since the eigenvalues of  $\mathbf{N}$  under the K-L transform are, by construction, 1, the elements of  $\Lambda$  above will all be greater than 1, and converging to 1 for the noise-dominated modes.

## 2. The K-L transform with correlated contaminants

Let us now consider a more general case, where we further split the noise into two parts  $\mathbf{n} \rightarrow \mathbf{n} + \mathbf{m}$ , where  $\mathbf{m}$  is a contaminant with a non-zero correlation with the signal. The covariance matrix of the data is then given by:

$$\langle \mathbf{x} \mathbf{x}^\dagger \rangle = \alpha^2 \mathbf{S} + 2\alpha \mathbf{M}_s + \mathbf{M} + \mathbf{N}, \quad (5)$$

where  $\mathbf{M}_s \equiv (\langle \mathbf{m} \mathbf{s}^\dagger \rangle + \langle \mathbf{s} \mathbf{m}^\dagger \rangle)/2$ ,  $\mathbf{M} \equiv \langle \mathbf{m} \mathbf{m}^\dagger \rangle$  and we have kept the fictitious parameter  $\alpha$  defined in the previous section.

Following the same steps as before, we obtain a similar generalized eigenvalue problem to that in Eq. 2:

$$\mathbf{C} \mathbf{e}_p = \lambda_p \tilde{\mathbf{N}} \mathbf{e}_p \quad (6)$$

with a modified noise covariance matrix  $\tilde{\mathbf{N}} \equiv \mathbf{M}_s + \mathbf{M} + \mathbf{N}$ .

## B. Application to tomographic datasets

The standard method to draw cosmological constraints from photometric redshift surveys is to divide the galaxy sample into bins in photo- $z$  space and use the information encoded in all the relevant auto- and cross-correlations between different bins, making use of various calibration methods in order to estimate the true redshift distribution of each bin. Several criteria can be followed in order to select these redshift bins, such as minimising the correlation between non-neighbouring bins or preserving a roughly constant number density on all bins. Other approaches (CITES) involve projecting the main observable (e.g. galaxy overdensity or shear) onto the Fourier-Bessel eigenbasis. None of these schemes are manifestly optimal from the point of view of  $S/N$ , final cosmological constraints or contaminant deprojection. This section presents an alternative slicing scheme addressing these shortcomings, based on the K-L transform.

### 1. Tomographic analyses

Let us start by assuming that we have split the galaxy sample into  $N_s$  subsamples. As mentioned above, we will think of each of these subsamples as some kind of redshift binning (e.g. binning galaxies in terms of their maximum-likelihood redshift), but the formalism applies to any set of subsamples. Let  $f^\alpha(\hat{\mathbf{n}})$  be the a field on the sphere at the angular position  $\hat{\mathbf{n}}$  and defined in terms of

the properties of the sources in the  $\alpha$ -th sample (e.g. the cosmic shear field  $\gamma^\alpha$  or the galaxy overdensity  $\delta^\alpha$ ), and let  $\phi^\alpha(z)$  be the redshift distribution of these sources. Finally, let  $a_{\ell m}^\alpha$  be the spherical harmonic coefficients of  $f^\alpha$ <sup>1</sup>. The power spectrum for our set of subsamples is defined as the two-point correlator of  $a^\alpha$ :

$$\langle \mathbf{a}_{\ell m} \mathbf{a}_{\ell' m'}^\dagger \rangle \equiv \delta_{\ell \ell'} \delta_{m m'} \mathbf{C}_\ell, \quad (7)$$

where we have packaged  $a^\alpha$  as a vector for each  $(\ell, m)$ :  $\mathbf{a}_{\ell m} \equiv (a_{\ell m}^1, \dots, a_{\ell m}^{N_s})$ . In general, the observed field will receive contributions from the true cosmological signal ( $\mathbf{s}$ ) and measurement noise ( $\mathbf{n}$ ).

Once the choice of subsamples  $\alpha$  is made, the standard analysis method would proceed by performing a likelihood evaluation of the two-point statistics of these subsamples. While this procedure is relatively simple, it suffers from a number of drawbacks, an incomplete list of which is:

1. It is not clear what the optimal strategy should be to define the sub-samples. One could make sure to exploit all of the information present in the data by using a large number of very narrow redshift bins, and let the likelihood evaluation pick up the information encoded in them.
2.  $\mathbf{C}_\ell^{\alpha\beta}$  is a  $N_s \times N_s \times N_\ell$  data vector. Thus increasing  $N_s$  will increase the computational time required for each likelihood evaluation like  $N_s^2$  and number of elements of the covariance matrix of  $\mathbf{C}_\ell^{\alpha\beta}$  like  $N_s^4$ , with the corresponding increase in complexity needed to estimate this covariance. Although this can be partially alleviated by considering only correlations between neighbouring redshift shells, the amount of information lost by neglecting all correlations beyond a given neighbouring index is not clear a priori.
3. Estimating the redshift distribution for a large number of subsamples can be inaccurate, depending on the method used to do so, on the quality of the photometric redshift posterior information and on the statistics of the available spectroscopic sample.

### 2. Optimal radial eigenbasis

Following the description in Section II A 1, it is straightforward to derive an optimal set of radial, uncorrelated eigenmodes.

<sup>1</sup> Spin-2 fields, such as the cosmic shear, will be decomposed in spin-2 spherical harmonics, however the discussion below holds for fields of arbitrary spin.

1. We start by assuming that the field  $\mathbf{a}$  has been measured in a number of narrow redshift bins, and by defining the inverse-variance weighted field  $\tilde{\mathbf{a}}_{\ell m} \equiv \mathbf{N}_\ell^{-1} \mathbf{a}_{\ell m}$ .
2. Let us consider a set of linear combinations of the weighted field measured on narrow redshift bins:

$$\mathbf{b}_{\ell m} = \mathbf{F}_\ell^\dagger \cdot \tilde{\mathbf{a}}_{\ell m} \equiv \mathbf{F}_\ell \circ \mathbf{a}, \quad (8)$$

where  $\mathbf{F}_\ell$  is a yet-unspecified matrix and, as in Section II A 1, we have let  $\mathbf{N}^{-1}$  define the non-standard dot product  $\mathbf{v}_\ell \circ \mathbf{w}_\ell \equiv \mathbf{v}_\ell^\dagger \cdot \mathbf{N}_\ell^{-1} \cdot \mathbf{w}_\ell$ . The power spectrum for this new observable would then simply be given by:

$$\mathbf{D}_\ell \equiv \langle \mathbf{b}_{\ell m} \mathbf{b}_{\ell m}^\dagger \rangle = \mathbf{F}_\ell^\dagger \circ \mathbf{C}_\ell \circ \mathbf{F}_\ell. \quad (9)$$

3. Requiring that the new modes be uncorrelated, we can identify Eq. 9 with the generalized eigenvalue equation 4, which defines the K-L eigenbasis  $\mathbf{F}_\ell$  by additionally requiring that it be orthonormal ( $\mathbf{F}_\ell \circ \mathbf{F}_\ell = \mathbf{I}$ ). Note that, after this transformation and without any further optimization, some of the practicalities of the original problem the original problem are already simplified, since we can now focus on the diagonal elements of the new power spectrum and its covariance.
4. The data can be further compressed if by assuming that we are interested in measuring a set of cosmological parameters  $\Theta \equiv \{\theta_1, \dots\}$ . The information regarding this set of parameters encoded in a given data vector  $\mathbf{x}$  can be quantified in terms of its Fisher matrix (the expectation value of the Hessian of the log-likelihood with respect to  $\Theta$ ), which assuming  $\langle \mathbf{x} \rangle = 0$  reads

$$\mathcal{F}_{ij} \equiv \langle \partial_i \partial_j \mathcal{L} \rangle = \frac{1}{2} \text{Tr} \left( \partial_i \mathbf{X} \mathbf{X}^{-1} \partial_j \mathbf{X} \mathbf{X}^{-1} \right), \quad (10)$$

where  $\mathbf{X} \equiv \langle \mathbf{x} \mathbf{x}^\dagger \rangle$  is the covariance matrix of the data. Since the power spectrum of  $\mathbf{b}_{\ell m}$  defined above is diagonal, this expression gets simplified further, and the Fisher matrix can be decomposed into the independent contributions of each mode:  $\mathcal{F}_{ij} = \sum_p \mathcal{F}_{ij}^p$ , where

$$\mathcal{F}_{ij}^p \equiv \sum_\ell \frac{2\ell+1}{2} (\partial_i \log D_\ell^p) (\partial_j \log D_\ell^p). \quad (11)$$

We can thus rank the eigenvectors  $(\mathbf{F}_\ell)_\alpha^p$  in terms of their information content (in a Fisher-matrix sense).

5. The final set of uncorrelated modes can then be truncated to the first  $M$  defined by this procedure, which will contain the bulk of the information needed to constrain  $\Theta$ .

This strategy therefore allows one to reliably and significantly reduce the dimensionality of the data vector from  $N_s^2 \times N_\ell$  to  $M \times N_\ell$  while minimising the loss of information. Note that, although the method is based on an initial thin-slicing of the galaxy distribution, the fact that the final dataset comprises only a small set of samples means that the method is not penalized in terms of photometric redshift uncertainties. Once the K-L eigenmodes  $\mathbf{F}_\ell$  are found for a fiducial cosmological model, they can be directly applied as weights to all the objects in the survey to generate the  $b_\ell^p$  modes. Furthermore, using  $\mathbf{F}_\ell$  for the fiducial cosmology as model-agnostic weights and inserting them in Eq. 9, the theoretical prediction for the power spectrum of each mode  $D_\ell^p$  can be computed in a model-independent way.

The same methods used to calibrate photo- $z$  uncertainties in the standard tomographic analysis hold in this case with slight modifications (e.g. weighed and  $\ell$ -dependent stacking of photo- $z$  pdfs, or cross-correlations of the weighed maps with a spectroscopic survey in the case of clustering redshifts).

### III. PERFORMANCE AND PARTICULAR EXAMPLES

This Section explores the performance of the K-L decomposition in a number of specific science cases.

#### A. Special case: the harmonic-bessel basis

Let us consider a simplified case where  $f$  is the overdensity field of a non-evolving galaxy population for which we neglect the effect of redshift-space distortions. Let us further assume that we have perfect redshift information, such that we can split the sample into thin radial slices of equal width  $\delta r$ , which we label by their comoving radius  $r$ . The noise in the measurement of  $f$  is given purely by shot noise, and since (as per our initial assumptions) the number density of sources does not change with  $r$ , the noise power spectrum is diagonal and scales like

$$N_\ell(r, r') \propto \frac{\delta_{r, r'}}{r^2}. \quad (12)$$

Thus, the dot product is just given by:

$$\mathbf{b}^\dagger \circ \mathbf{c} \propto \int dr r^2 b(r)^* c(r). \quad (13)$$

In this case, the cross-shell power spectrum is given by,

$$C_\ell^{rr'} = \frac{2}{\pi} \int_0^\infty dk k^2 P_k j_\ell(kr) j_\ell(kr'), \quad (14)$$

and it is trivial to show that the K-L eigenmodes are simply given by the spherical Bessel functions:  $(\mathbf{F}_\ell)_r^k \propto \sqrt{2/\pi} j_\ell(kr)$ :

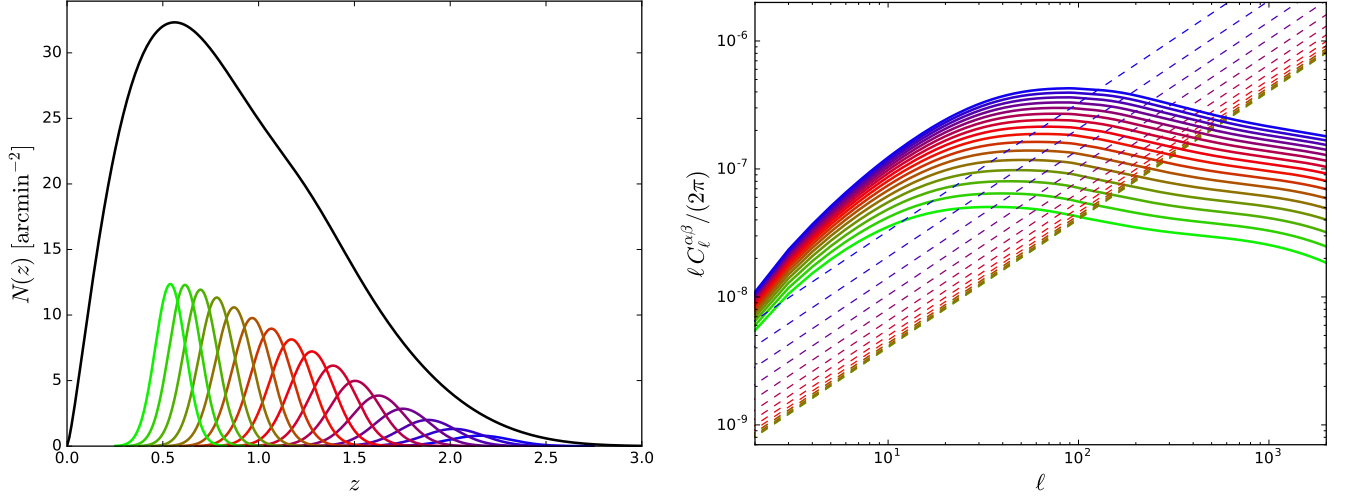


FIG. 1. *Left*: redshift distribution and bins considered for the K-L analysis of a strawman weak-lensing survey. *Right*: shear auto-power spectra of the redshift bins shown in the left panel. The signal and noise power spectra are shown as thick solid and thin dashed lines respectively.

$$D_{\ell}^{kk'} \propto \frac{2}{\pi} \int dr r^2 \int dr' r'^2 j_{\ell}(kr) j_{\ell}(k'r') C_{\ell}^{rr'} \quad (15)$$

$$= \int dq q^2 P_q \left[ \frac{2}{\pi} \int dr r^2 j_{\ell}(qr) j_{\ell}(kr) \right] \left[ \frac{2}{\pi} \int dr' r'^2 j_{\ell}(qr') j_{\ell}(k'r') \right] \quad (16)$$

$$= \int dq q^2 P_q \frac{\delta(k-q)}{q^2} \frac{\delta(k'-q)}{q^2} = P_k \frac{\delta(k-k')}{k^2} = \frac{P_k}{k^2 \Delta k} \delta_{k,k'} \quad (17)$$

This choice of basis defines the so-called harmonic-Bessel (or Fourier-Bessel) decomposition, and has been postulated as a possible data-compression method for the analysis of photometric redshift data (CITES here). In any realistic scenario however (e.g. in the presence of redshift uncertainties, RSDs or for the analysis of weak lensing data), this basis is non-optimal (e.g. different  $k$ s will be correlated), as opposed to the K-L basis described above.

### B. Weak lensing - K-L basis for dark energy

To quantify the performance of the K-L modes for weak lensing we study the case of an LSST-like survey. The survey specifications and the characteristics of the galaxy sample are described in detail in CITE. In summary, we assume a sample with  $\sim 29$  objects per arcmin<sup>2</sup> with the redshift distribution shown in the left panel of Figure 1. We also approximate the photo- $z$  distribution as Gaussian with a scatter  $\sigma_z = 0.05(1+z)$ .

The signal part of the cross-power spectrum between the cosmic shear measurements made in two different red-

shift shells is given by:

$$S_{\ell}^{\alpha\beta} = \frac{2}{\pi} \int_0^{\infty} dk k^2 \Delta_{\ell}^{\alpha}(k) \Delta_{\ell}^{\beta}(k), \quad (18)$$

where the transfer functions  $\Delta_{\ell}^{\alpha}$  take the form:

$$\Delta_{\ell}^{\gamma,\alpha}(k) \equiv \sqrt{\frac{(\ell+2)!}{(\ell-2)!}} \int d\chi W^{\alpha}(\chi) \frac{j_{\ell}(k\chi)}{k^2 a(\chi)} \sqrt{P(k, z(\chi))},$$

$$W^{\alpha}(\chi) \equiv \frac{3H_0^2 \Omega_M}{2} \int_{z(\chi)}^{\infty} dz \phi^{\alpha}(z') \frac{\chi(z') - \chi}{\chi(z') \chi}. \quad (19)$$

Here  $\phi^{\alpha}(z)$  is the redshift distribution of sources in the  $\alpha$ -th bin. The noise power spectrum is white and simply given by the intrinsic ellipticity scatter weighed by the number density of sources in each redshift bin  $\bar{n}^{\alpha}$ :

$$N_{\ell}^{\alpha\beta} = \delta_{\alpha\beta} \frac{\sigma_{\gamma}^2}{\bar{n}^{\alpha}}, \quad (20)$$

with  $\bar{n}^{\alpha}$  in units of srad<sup>-1</sup>. We use  $\sigma_{\gamma} = 0.28$ .

As our initial set of narrow redshift bins, we select top-hat bins in photo- $z$  space for  $z_{\text{ph}} > 0.5$  with a width given by the value of  $\sigma_z$  at the center of the bin. The resulting set of 16 bins is shown in the left panel of Figure 1. The

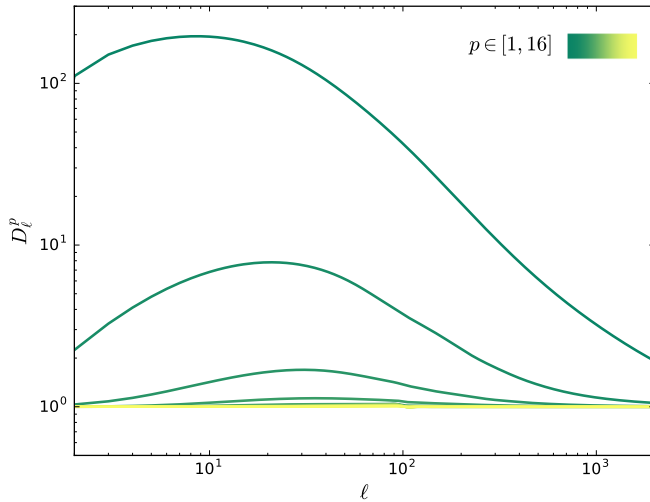


FIG. 2. Power spectra of the K-L eigenmodes for the strawman weak lensing survey. All but the first three modes are noise-dominated, and most of the information is encoded in the first mode.

large overlap between bins implies that a choice of thinner slices is unlikely to unveil significantly more information, and we have verified that the results shown below do not change after doubling the number of bins. The lensing auto-power spectra (both signal and noise) for these bins are shown in the right panel of Figure 3. The elements of  $C_\ell^{\alpha\beta}$  were estimated using a modified version of the code presented in CITE.

We compute the K-L modes for this setup and rank them according to their information content on the dark energy equation of state  $w$ . The power spectra of the resulting set of modes are shown in the left panel of Figure 2. Comparing against the right panel of Fig. 1 we can see that the K-L decomposition effectively separates the signal-dominated and noise-dominated modes, with all modes  $p > 3$  dominated by noise (as we mentioned in Section II A 1, the noise power spectrum gets mapped into 1 under the K-L transform). The fractional contribution of each mode to the total constraint on  $w$  (i.e. its contribution to the corresponding Fisher matrix element) is shown in the left panel of Figure 3. Most of the information ( $\sim 95\%$ ) is contained within a single mode, and the first two modes are able to recover more than 99% of the total. The eigenvectors corresponding to the first and second modes for different values of  $\ell$  are shown in the right panel of the same figure. Firstly, we observe that the eigenvectors preserve roughly the same shape for all  $\ell$ , and converge to the same shape at large  $\ell$ . The first eigenvector upweights the parts of the redshift range with the highest signal-to-noise, penalising the low- $z$  regime due to its poor lensing signal and the high- $z$  bins due to their high shot noise. The second eigenmode then recovers part of this information by marginally upweighting these regions.

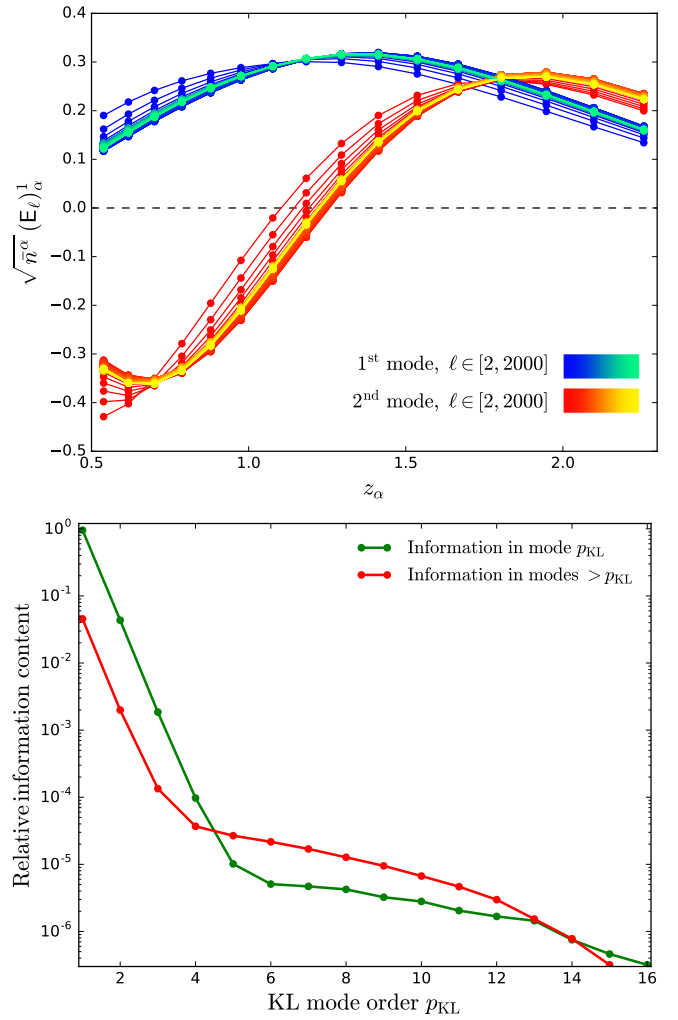


FIG. 3. *Left*: the first (blue-cyan) and second (red-yellow) eigenmodes of the strawman weak-lensing survey for different  $\ell$ . The redshift dependence of the modes stays roughly constant across  $\ell$  and converges to a fixed shape for large  $\ell$ . *Right*: information content of the different eigenmodes. Most of the information  $\sim 95\%$  is encoded in the first bin, and more than 99% of it can be recovered considering only the first 2 modes.

### C. Galaxy clustering - measuring $f_{\text{NL}}$

It is expected that future large-scale photometric surveys will make the search for primordial non-Gaussianity one of the main science cases. This can be achieved through the excess power a non-zero value of  $f_{\text{NL}}$  generates in the two-point statistics of biased tracers of the matter distribution on large scales CITES. Since the signal is most relevant on large scales, one can expect most of the information on  $f_{\text{NL}}$  to be concentrated in a small number of radial modes, which makes the K-L decomposition described above an ideal analysis method.

To explore this possibility we have considered a strawman photometric survey targeting a sample of red galax-



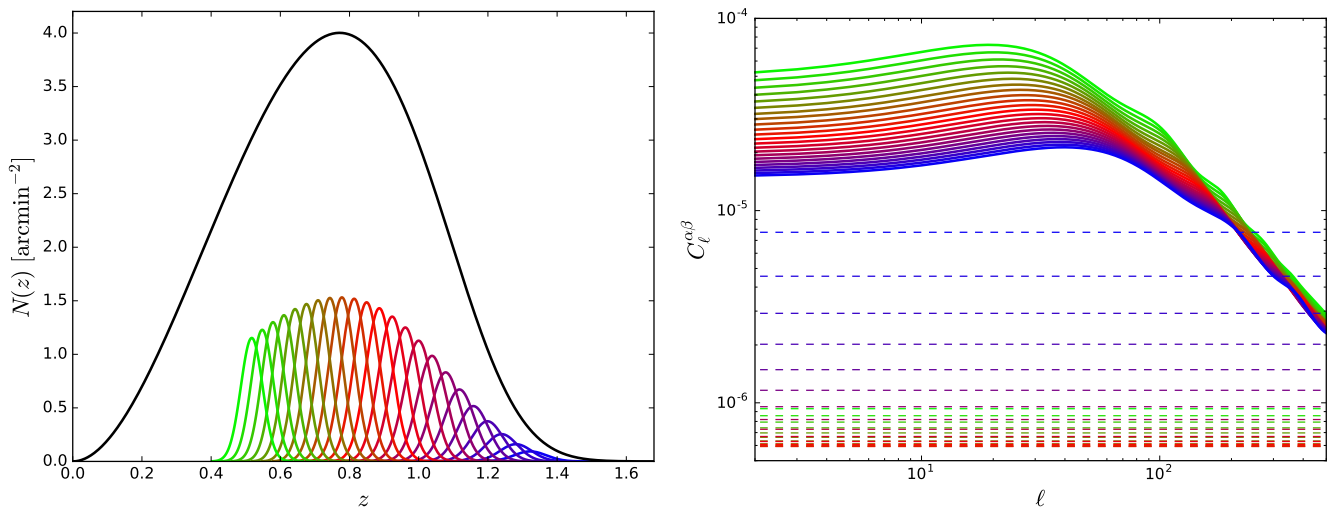


FIG. 4. *Left*: redshift distribution and bins considered for the K-L analysis of a strawman large-scale-structure survey targeting a sample of red galaxies. *Right*: clustering auto-power spectra of the redshift bins shown in the left panel. The signal and noise power spectra are shown as thick solid and thin dashed lines respectively.

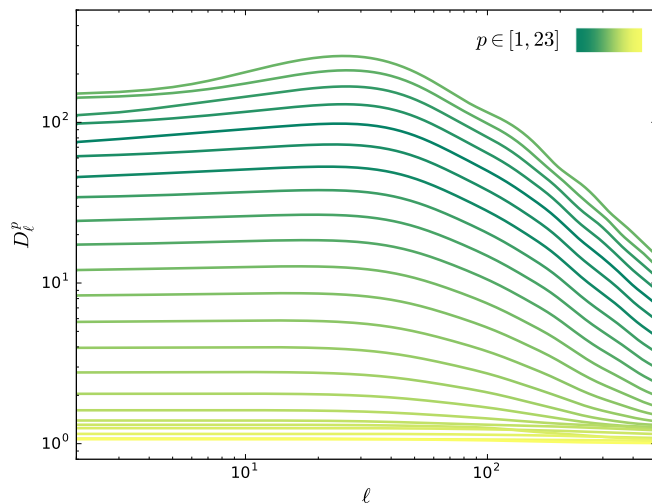


FIG. 5. Power spectra of the K-L eigenmodes for the strawman weak large-scale-structure survey. Unlike in the case of weak lensing, a large number of eigenmodes are signal-dominated. This is due to the overall higher signal-to-noise ratio of galaxy clustering with respect to galaxy shear as well as to the smaller correlations between distant bins.

ies, characterized by their higher bias and better photo- $z$  uncertainties than their blue counterparts. The sample we consider is compatible with what could be observed by LSST, characterized by the redshift distribution shown in the left panel of Fig. 4 (full details can be found in CITE). We assume a photo- $z$  uncertainty of  $\sigma_z = 0.02(1+z)$  and split the sample into redshift bins in photo- $z$  space with  $z_{\text{ph}} > 0.5$  and a width given by the photo- $z$  uncertainty at the bin centre. The resulting set of 23 bins is shown in the left panel of Fig. 4. In the case

of galaxy clustering, the transfer functions entering the signal power spectrum (see Eq. 18) are given by:

$$\Delta_{\ell}^{\delta, \alpha}(k) \equiv \int dz \phi^{\alpha}(z) \Psi_{\ell}(k, z) \sqrt{P(k, z)},$$

$$\Psi_{\ell}(k, z) = b^{\alpha}(z) j_{\ell}(k \chi(z)) - f(z) j_{\ell}''(k \chi(z)) \quad (21)$$

where  $b^{\alpha}(z)$  is the galaxy bias and  $f(z) \equiv d \log \delta / d \log a$  is the growth rate of structure (we have kept the contribution from redshift-space distortions at linear order but ignored the effect of magnification bias). The auto-power spectra for our set of 23 bins are shown in the right panel of Fig. 4.

Using the prescription described in Section II A, we find the K-L eigenmodes and associated power spectra, and rank them according to their contribution to the total Fisher matrix element of  $f_{\text{NL}}$ . The power spectra of the K-L modes are shown in Figure 5. Unlike the case of weak lensing, the information encoded in the galaxy overdensity is local in redshift (as opposed to integrated). Thus the correlation between different bins decays rapidly with redshift separation, and we find a much larger number of signal-dominated modes (e.g. compare with the analogous figure for weak lensing, Fig. 2). The information on  $f_{\text{NL}}$  contained in each mode is shown in the left panel of Fig. 6. In this case, the information is evenly spread across the first  $\sim 10$  modes, and 90% of the total constraining power can be achieved by considering the first 12 eigenvectors. The form of the first 5 of these eigenvectors for  $\ell = 30$  are shown in the right panel of Fig. 6. The eigenmodes are sinusoids with varying frequencies, in agreement with the expectation that, in the limit of  $\sigma_z \rightarrow 0$  and no background redshift dependence, the K-L decomposition is achieved by the spherical Bessel functions (see Section III A).

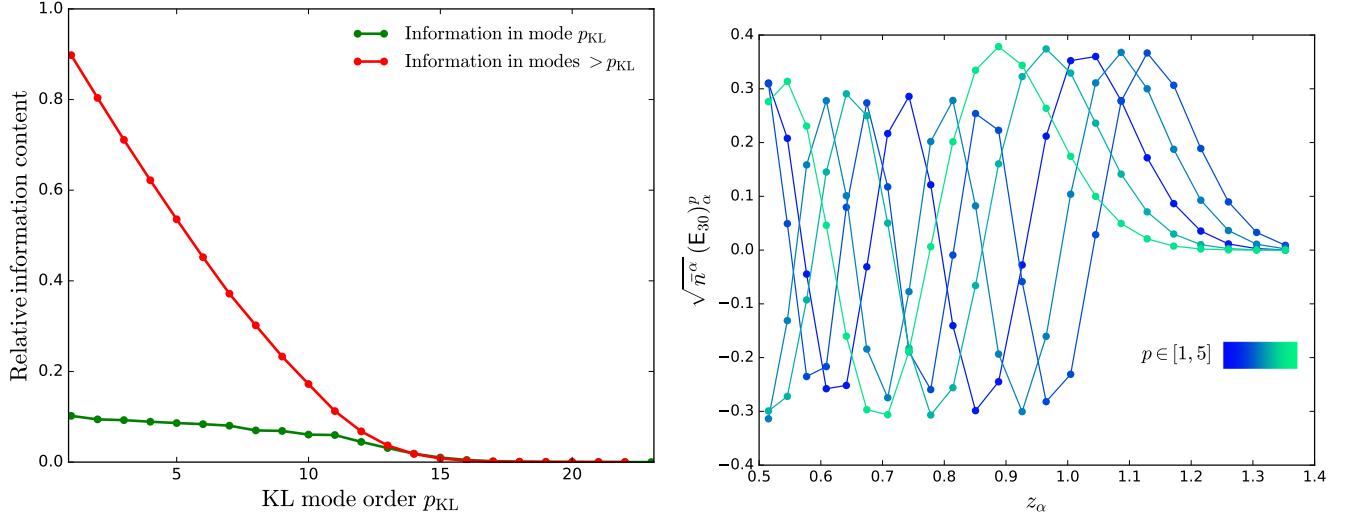


FIG. 6. *Left*: information content of the different K-L eigenmodes for the strawman galaxy clustering survey. The bulk of the information ( $> 90\%$ ) on  $f_{NL}$  is encoded in the first 12 modes. *Right*: the first 5 K-L modes for  $\ell = 30$ . The sinusoidal shape of the modes agrees with the expectation that, in the limit of  $\sigma_z \rightarrow 0$  and no background redshift dependence, the K-L modes should be given by the spherical Bessel functions.

## D. Correlated contaminants

### 1. Intrinsic alignments

Contamination of the weak lensing signal by locally correlated galaxy shapes, commonly known as “intrinsic alignments” (IA), is one of the main concerns for shear surveys. This contaminant cannot be described as an uncorrelated source of irreducible noise, since the causes of IA (e.g. the local tidal field in the non-linear alignment model CITE) are likely correlated the lensing potential.

The K-L decomposition, as described in Section II A 2, can therefore be used to design an optimal basis that maximizes the true lensing signal over any correlated component. The method is similar to the so-called “nulling” approach of CITE.

### 2. Magnification bias

Gravitational lensing of the galaxy positions alters the clustering pattern of galaxies through the so-called magnification bias effect. This appears as an extra term in the galaxy clustering transfer function (Eq. 21):

$$\Delta_\ell^{M,\alpha}(k) = -2\ell(\ell+1) \int d\chi W^{M,\alpha}(\chi) \frac{j_\ell(k\chi)}{k^2 a(\chi)} \sqrt{P(k, z(\chi))},$$

$$W^{M,\alpha}(\chi) = \frac{3H_0^2 \Omega_M}{2} \int_{z(\chi)}^\infty dz' \phi^\alpha(z') \frac{2-5s}{2} \frac{\chi(z') - \chi}{\chi(z')\chi}, \quad (22)$$

where  $s$  is the tilt in the number counts of sources as a function of magnitude limit.

With the aim of measuring this effect, one can think of the density and RSD terms as correlated contaminants on top of it (i.e. taking the place of intrinsic alignments in the case of shear measurements), and thus develop an optimal K-L eigenbasis to measure it CITE.

## IV. PRACTICAL EXAMPLE: WEAK LENSING

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

## V. DISCUSSION

Nulla malesuada porttitor diam. Donec felis erat, congue non, volutpat at, tincidunt tristique, libero. Vivamus viverra fermentum felis. Donec nonummy pellentesque ante. Phasellus adipiscing semper elit. Proin fermentum massa ac quam. Sed diam turpis, molestie vitae, placerat a, molestie nec, leo. Maecenas lacinia. Nam ipsum ligula, eleifend at, accumsan nec, suscipit a, ipsum. Morbi blandit ligula feugiat magna. Nunc eleifend consequat lorem. Sed lacinia nulla vitae enim. Pellentesque tincidunt purus vel magna. Integer non enim. Praesent euismod nunc eu purus. Donec bibendum quam in tel-

lus. Nullam cursus pulvinar lectus. Donec et mi. Nam vulputate metus eu enim. Vestibulum pellentesque felis eu massa.

### ACKNOWLEDGEMENTS

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel justo vitae lacus tincidunt ultrices. Lorem ip-

sum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

---

### Appendix A: Pseudo- $C_\ell$ estimation of the K-L modes

One of the standard methods to estimate the angular power spectrum of any two quantities in the cut sky is the so-called pseudo- $C_\ell$  estimator. This section adapts this method to the modes resulting from the K-L decomposition described before.

The standard pseudo- $C_\ell$  method is based on computing the spherical harmonic coefficients of the mask field:

$$\tilde{a}_{\ell m}^\alpha = \int d\hat{\mathbf{n}} a^\alpha(\hat{\mathbf{n}}) w^\alpha(\hat{\mathbf{n}}), \quad (\text{A1})$$

where  $w^\alpha$  is the weights map characterizing the mask of the field  $a^\alpha$ . One then estimates the power spectrum of this object by averaging over  $m$  for each  $\ell$ :

$$\tilde{C}_\ell^{\alpha\beta} \equiv \frac{\sum_m \tilde{a}_{\ell m}^\alpha \tilde{a}_{\ell m}^{\beta*}}{2\ell + 1}. \quad (\text{A2})$$

This object is then related to the true underlying power spectrum through a mode-coupling matrix  $M_{\ell\ell'}^{\alpha\beta}$  such that

$$\tilde{C}_\ell^{\alpha\beta} = \sum_{\ell'} M_{\ell\ell'}^{\alpha\beta} C_{\ell'}^{\alpha\beta}, \quad M_{\ell\ell'}^{\alpha\beta} \equiv \sum_{\ell''} \frac{(2\ell' + 1)(2\ell'' + 1)}{4\pi} W_{\ell''}^{\alpha\beta} \begin{pmatrix} \ell & \ell' & \ell'' \\ 0 & 0 & 0 \end{pmatrix}^2 \quad (\text{A3})$$

where the coupling matrix  $M$  depends solely on the power spectrum of the masks  $W_\ell^{\alpha\beta} \equiv (2\ell + 1)^{-1} \sum_m w_{\ell m}^\alpha w_{\ell m}^{\beta*}$ .

The extension of this estimator to the power spectrum of the K-L modes is straightforward: we project the masked harmonic coefficients  $\tilde{a}^\alpha$  over the K-L eigenvectors  $\mathbf{E}$  (i.e.  $\tilde{\mathbf{b}}_{\ell m} \equiv \mathbf{E}_\ell \circ \tilde{\mathbf{a}}_{\ell m}$ ) and compute their power spectra by averaging over  $m$ . The resulting estimator takes the form  $\tilde{D}_\ell^p = \sum_{\ell'} M_{\ell\ell'}^{pp'} D_{\ell'}^{p'}$ , where the new mode-coupling matrix is given by:

$$M_{\ell\ell'}^{pp'} \equiv M_{\ell\ell'}^{\alpha\beta} \left[ (\mathbf{E}_\ell)_\alpha^p (\mathbf{N}^{-1})_{\alpha\alpha'} (\mathbf{E}_{\ell'})_{\alpha'}^{p'} \right] \left[ (\mathbf{E}_\ell)_\beta^p (\mathbf{N}^{-1})_{\beta\beta'} (\mathbf{E}_{\ell'})_{\beta'}^{p'} \right] = M_{\ell\ell'} \left[ (\mathbf{E}_\ell)_\alpha^p (\mathbf{N}_\ell^{-1})_{\alpha\beta} (\mathbf{E}_{\ell'})_{\beta}^{p'} \right]^2 \quad (\text{A4})$$

where the second equality holds only if all the maps  $a_\ell^\alpha$  share the same mask  $w$ .<sup>2</sup>

---

<sup>2</sup> Note that, for full-sky coverage  $M_{\ell\ell'} = \delta_{\ell\ell'}$  and using the or-

thonormality of  $\mathbf{E}$  we get  $M_{\ell\ell'}^{pp'} = \delta_{\ell\ell'} \delta_{pp'}$ .