

# Ridges in the Dark Energy Survey for cosmic trough identification

Ben Moews,<sup>1\*</sup> Morgan A. Schmitz,<sup>2</sup> Andrew J. Lawler,<sup>3</sup> Joe Zuntz,<sup>1</sup> Alex I. Malz,<sup>4</sup> Rafael S. de Souza,<sup>5</sup> Ricardo Vilalta,<sup>6</sup> Alberto Krone-Martins,<sup>7,8</sup> Emille E. O. Ishida,<sup>9</sup> for the COIN Collaboration

<sup>1</sup>*Institute for Astronomy, University of Edinburgh, Royal Observatory, Edinburgh EH9 3HJ, UK*

<sup>2</sup>*Department of Astrophysical Sciences, Princeton University, 4 Ivy Ln., Princeton, NJ08544, USA*

<sup>3</sup>*Department of Statistics, Baylor University, One Bear Place #97140, Waco, TX 76798, USA*

<sup>4</sup>*Ruhr-University Bochum, Astronomical Institute, German Centre for Cosmological Lensing, Universitätsstr. 150, 44801 Bochum, Germany*

<sup>5</sup>*Key Laboratory for Research in Galaxies and Cosmology, Shanghai Astronomical Observatory, Chinese Academy of Sciences, 80 Nandan Rd., Shanghai 200030, China*

<sup>6</sup>*Department of Computer Science, University of Houston, 4800 Calhoun Rd., Houston TX, USA*

<sup>7</sup>*Donald Bren School of Information and Computer Sciences, University of California, Irvine, CA 92697, USA*

<sup>8</sup>*CENTRA/SIM, Faculdade de Ciências, Universidade de Lisboa, Ed. C8, Campo Grande, 1749-016, Lisboa, Portugal*

<sup>9</sup>*Université Clermont Auvergne, CNRS/IN2P3, LPC, F-63000 Clermont-Ferrand, France*

Accepted XXX Received XXX; in original form XXX

## ABSTRACT

Cosmic voids and their corresponding redshift-aggregated projections of mass densities, known as troughs, play an important role in our attempt to model the large-scale structure of the Universe. Understanding these structures leads to tests comparing the standard model with alternative cosmologies, constraints on the dark energy equation of state, and provides evidence to differentiate among gravitational theories. In this paper, we extend the subspace-constrained mean shift algorithm, a recently introduced method to estimate density ridges, and apply it to 2D weak-lensing mass density maps from the Dark Energy Survey Y1 data release to identify curvilinear filamentary structures. We compare the obtained ridges with previous approaches to extract trough structure in the same data, and apply curvelets as an alternative wavelet-based method to constrain densities. We then invoke the Wasserstein distance between noisy and noiseless simulations to validate the denoising capabilities of our method. Our results demonstrate the viability of ridge estimation as a precursor for denoising weak lensing quantities to recover the large-scale structure, paving the way for a more versatile and effective search for troughs.

**Key words:** cosmology: large-scale structure of Universe – gravitational lensing: weak – methods: statistical – methods: data analysis

## 1 INTRODUCTION

The cosmic web, a latticework structure of enormous proportions, represents one of the largest physical patterns in the Universe (Bond et al. 1996). Its filamentary nature is a byproduct of the hierarchical growth of large-scale structure, and gives rise to four main classes of substructures: Galaxy clusters, filaments, sheets, and the large regions of near-emptiness known as *cosmic voids* (Zeldovich et al. 1982).

These voids are characterised by underdensities in the dark matter distribution, presenting much simpler dynamics than their non-linear and high-density counterparts (Hamaus et al. 2016). They are valuable cosmological probes, since they can encode relevant cosmological information and suffer from fewer sources of

systematical error (Peebles 2001; Lavaux & Wandelt 2012). Their population statistics may be predicted from a given cosmological theory (see Fry 1986; White et al. 1987; Li 2011), which can provide observational constraints to current models (Hoyle & Vogeley 2004; Gruen et al. 2018). However, in spite of their usefulness, the detection of voids is no trivial task.

The challenge starts from the lack of a dominant definition thereof, and their detection remains a focal topic of interest in cosmology (Cai et al. 2015; Gruen et al. 2016; Sánchez et al. 2017; Nadathur et al. 2017; Adermann et al. 2018; Brouwer et al. 2018; Xu et al. 2019). Several established approaches to void detection employ Voronoi tessellation as described by El-Ad & Piran (1997) and Gaite (2005). Galaxy particles can, for example, be enclosed in distance-based Voronoi cells and grouped into larger zones, where a watershed simulation naturally leads to the identification of large basins (low-density regions) corresponding to

\* E-mail: bmoews@roe.ac.uk

voids (Neyrinck 2008). Markedly different approaches to void identification have also been proposed. For instance, Aragón-Calvo et al. (2007) use computer vision techniques to classify void morphology, applying scale-independent morphology filters to identify primary cosmological structures such as walls, filaments, and voids.

Alternatively, a different family of techniques exploits notions of data topology for void identification. Aragón-Calvo et al. (2010), for example, apply topological segmentation, while Xu et al. (2019) apply notions of topological data analysis to find dimensional holes via persistent homology, in which zero, one, and two-dimensional holes are identified as clusters, loops of filaments, and voids, respectively. In the latter case, voids are considered statistically significant if their structure persists as long as data neighbourhoods increase in size. Classifications yielded by these techniques exhibit differences that impact the science case for which each method was developed (Libeskind et al. 2018).

Despite the difficulties posed by the automation of void detection, the rich variety of current methods and techniques have led to important advances in cosmology. As an example, the accurate location and modelling of voids can be exploited to derive clustering and abundance statistics such as void mass, and to constrain dark energy (Pisani et al. 2015). The dynamics of matter flowing away from the centre of a void are instrumental to gain more insight on cosmological parameters, as described by (Dekel & Rees 1994), while other applications include probing alternative dark matter models and tests of general relativity (Yang et al. 2015; Barreira et al. 2015).

Another difficulty inherent to void detection is the need for accurate redshift measurements. While this calls for larger, and ideally complete, galaxy spectroscopic surveys, these tend to cover only small areas of the sky. Photometric galaxy surveys, such as the Panoramic Survey Telescope and Rapid Response System (Pan-STARRS; Chambers et al. 2016), the Dark Energy Survey (DES; Flaugher et al. 2015), and the upcoming Vera C. Rubin Observatory Legacy Survey of Space and Time (LSST; Ivezić et al. 2008), on the other hand, can provide observations covering larger areas.

Beyond these studies of cosmic voids, or the common study of the galaxy clusters that occupy the nodes of the cosmic web, the characterisation of the filaments that connect them can also be of great interest. The filament between clusters Abell 0399 and 0401, for example, has been shown to host quiescent galaxies and hot gas, and to emit in radio (Bonjean et al. 2018; Govoni et al. 2019). Automated detection of filaments in a large volume has recently been carried out by Malavasi et al. (2020), using galaxy samples from the Sloan Digital Sky Survey (SDSS; York et al. 2000), while Galárraga-Espinosa et al. (2020) apply a similar approach to hydrodynamical simulations to derive the expected statistical properties of filaments.

In these studies, the distributions of galaxies from spectroscopic surveys (or simulations thereof) are used to detect filaments. Alternatively, weak gravitational lensing can be used as the observable to extract filaments from (Mead et al. 2010; Maturi & Merten 2013). While a few individual detections have been reported between specific clusters, as described by Dietrich et al. (2012) and Jauzac et al. (2012), the very low signal-to-noise ratio of the filament signal typically requires the stacking of large numbers of pairs of clusters to make detection from weak lensing possible (see, for example, Xia et al. 2020, and references therein). The very low amplitude of the lensing signal of filaments is enough to make detection, from wide photometric surveys and using lensing observables, a very hard problem. Even if this was not the case, any such attempt to characterise many filaments over a large 3D volume

would suffer from the same dependency on redshift accuracy as that discussed in the case of studies focussed on voids.

A way to avoid relying on accurate redshift measurements, while still readily taking advantage of photometric surveys, is to consider the 2D-projected counterparts to the elements that make up large-scale structure instead (though see also Sánchez et al. 2017, for a void finder built to work on photometric surveys, using redshift slices instead of a full 2D projection). In the case of voids, these are known as *troughs*, which represent the most underdense regions on the sky. Since troughs comprise regions of lower density across the line of sight in the projected space only, it is sufficient to use photometric measurements, obviating spectroscopic redshifts (Clampitt et al. 2013; Gruen et al. 2016).

In this work, we propose an algorithm to detect 2D density ridges as a way to denoise cosmic structure in mass density maps. For this purpose, we extend the subspace-constrained mean shift (SCMS) algorithm introduced by Ozertem & Erdogmus (2011) to fit our application case, and apply our method to the DES Y1 data release (Flaugher et al. 2015). The general methodology works by defining a denoised and sparse representation of the filamentary structure, and as a consequence, the locations of regions of emptiness emerge naturally. Considering the aforementioned limitations of spectroscopic surveys, trough finders can be applied as an alternative to void finders to recover and study underdense regions from weak lensing studies, giving particular relevance to the present application to DES Y1.

While previous work demonstrates how to identify filaments as density ridges using the SCMS algorithm (see Chen et al. 2015a,b, 2016, further discussed in Section 2), our method extends past implementations in several ways. We incorporate the haversine distance, a more suitable approach for spherical surfaces, and also customise the mesh size for ridge estimation and optimisation of the bandwidth. Another important difference is that these studies apply the algorithm to galaxy data, either from simulations or SDSS, or directly to dark matter particles from N-body simulations. In this work, we apply the approach to 2D weak lensing mass maps instead; the ridges we recover are thus extracted from the projected matter distribution and not directly comparable to the filaments of the cosmic web. Finally, we compare our ridges to a search based on curvelets, an extension of the wavelet transform more suitable for the study of curvy and filamentary structures (Candès et al. 2006).

This paper is organised as follows. Section 2 explains the SCMS algorithm and implemented extensions; it includes our kernel density estimation technique, experimental data, and simulations. Section 3 describes our experimental results and compares them to a curvelet-based denoising technique. Section 4 comments on our approach and future directions, with an emphasis on the advantages and drawbacks of the proposed methodology. Finally, Section 5 provides a summary and conclusions.

## 2 METHODOLOGY AND DATA

This section provides background information on the SCMS algorithm in Section 2.1, and describes past applications, as well as extensions from both this and prior work, in Section 2.2. Lastly, we introduce the Dark Energy Survey and its mass maps, together with our approach to sample generation and the creation of noisy and noiseless simulations for verification purposes in Section 2.3.

## 2.1 Subspace-constrained mean shift

Introduced by [Ozertem & Erdogmus \(2011\)](#), the SCMS algorithm is a recent addition to statistical methods dealing with the estimation of density ridges. Starting with a mesh of points placed in equidistant steps across the parameter space, the algorithm seeks to establish local principal curves in iterative steps. This can be visualised as a cloud of points shifting closer towards the nearest underlying structure at each iteration, akin to the process in which mass in our universe converges toward better-defined cosmic filaments over time. The latter can be observed in N-body simulations such as the Millennium Simulation by [Springel et al. \(2005\)](#) and its Millennium-II successor by [Boylan-Kolchin et al. \(2009\)](#), as well as the Bolshoi simulation by [Klypin et al. \(2011\)](#) and the MultiDark simulation ([Riebe et al. 2013](#)).

In more formal terms, such a ridge is the maximisation of the local density in the normal direction as given by the Hessian matrix. Let  $\nabla p(x)$  be the gradient of a probability density function  $p$  on a space of dimension  $d$ , and  $H(x)$  its Hessian matrix of second derivatives. We diagonalize  $H$ :

$$H(x) = U(x)\lambda(x)U(x)^\top. \quad (1)$$

with the RHS terms in descending order of  $\lambda$ . We then take the  $d-1$  eigenvectors (columns of  $U$ ) which correspond to the  $d-1$  smallest eigenvalues  $\lambda$ , i.e. omitting the column for the largest eigenvalue, since this is the direction parallel to the ridge.

Taking these eigenvectors  $v_i$  and their linear projection operator:

$$L(x) \propto L(H(x)) = v'v'^\top, \quad (2)$$

then we can project the gradient of the probability onto the eigenvectors as:

$$G(x) = L(x)\nabla p(x). \quad (3)$$

A ridge  $R$  can then be expressed as ([Genovese et al. 2014](#)) the locations  $x$  where  $G$  is zero everywhere and the omitted largest eigenvalue is positive<sup>1</sup>:

$$R = \{x : \|G(x)\| = 0, \lambda_{d+1}(x) < 0\}. \quad (4)$$

Algorithm 1 shows pseudocode describing the SCMS algorithm as it appears in [Moews et al. \(2019\)](#), including thresholding. Line 4 shows a kernel density estimation with a radial basis function (RBF) kernel,  $\mathcal{K}(x) = (1/\sqrt{2\pi})\exp(-0.5x^2)$ , meaning

$$\text{KDE}_{\text{RBF}}(x, \beta) = \frac{1}{|\theta|(2\pi\beta^2)^{\frac{d}{2}}} \sum_{i=1}^{|\theta|} e^{-\frac{\|x-\theta_i\|^2}{2\beta^2}}. \quad (5)$$

## 2.2 Extensions and previous applications

Since its inception a few years ago, the SCMS algorithm was found to be a valuable tool in a variety of fields. This includes the analysis of 3D neuron structures in tissue images by [Bas & Erdogmus \(2011\)](#), as well as the identification of road networks in satellite images by combining the SCMS algorithm with the geodesic method and tensor voting ([Miao et al. 2014](#)).

<sup>1</sup> without this second condition we also locate valleys

### Algorithm 1 SCMS with thresholding as in ([Moews et al. 2019](#))

```

1: Input: Coordinates  $\theta$ , bandwidth  $\beta$ , threshold  $\tau$ , iterations  $N$ 
2: Output: Density ridge point coordinates  $\psi$ 
3: procedure SCMS( $\theta, \beta, \tau, N$ )
4:    $\kappa(x) \leftarrow \text{KDE}_{\text{RBF}}(\theta, \beta)$ , using Eqn. 5
5:    $\xi_1 \leftarrow (\min(\theta_{*,1}), \max(\theta_{*,1}))$ 
6:    $\xi_2 \leftarrow (\min(\theta_{*,2}), \max(\theta_{*,2}))$ 
7:    $\psi \leftarrow \psi \sim U(\xi_1, \xi_2)_{|\theta|}$ 
8:    $\psi \leftarrow \forall y \in \psi : \kappa(y) < \tau$ 
9:   for  $n \leftarrow 1, 2, \dots, N$  do
10:     for  $i \leftarrow 1, 2, \dots, |\psi|$  do
11:       for  $j \leftarrow 1, 2, \dots, |\theta|$  do
12:          $\mu_j = \frac{\psi_i - \theta_j}{\beta^2}$ 
13:          $\sigma_j = \mathcal{K}_{\text{RBF}}\left(\frac{\psi_i - \theta_j}{\beta}\right)$ 
14:       end for
15:        $H(x) = \frac{1}{|\theta|} \sum_{j=1}^{|\theta|} \sigma_j \left( \mu_j \mu_j^\top - \frac{1}{\beta^2} \mathbb{I} \right)$ 
16:        $v, \lambda \leftarrow v, \lambda$  from eigendecomp. eig( $H(x)$ )
17:        $v' \leftarrow$  entries in  $v$  corresp. to sortasc( $\lambda$ )1,2,\dots,d-1
18:        $\psi_i \leftarrow v'v'^\top \frac{\sum_{j=1}^{|\psi|} \sigma_j \theta_j}{\sum_{j=1}^{|\psi|} \sigma_j}$ 
19:     end for
20:   end for
21:   return  $\psi$ 
22: end procedure

```

### 2.2.1 Applications to astronomy and thresholding

The first application of the SCMS algorithm in astronomy is also the one that is most closely related to our work. For the purpose of investigating galaxy evolution, [Chen et al. \(2015a\)](#) employ the algorithm to constrain matter distributions at different redshifts, introducing the concept of thresholding. This refers to the ‘cutting’ of the initial grid of unconverged ridge points according to a kernel density estimate of the data, effectively deleting coordinates falling within low-density areas to avoid the identification of filaments in these areas. More formally, the threshold  $\tau$  is determined by computing the root mean square of the differences between the average density estimate and the grid points’ density estimates,

$$\tau = \sqrt{\frac{\sum (\phi - \bar{\phi})^2}{G}}, \quad (6)$$

for an element-wise subtraction of an array of grid point density estimates  $\phi$  and their average  $\bar{\phi}$ , as well as the number of grid points  $G$ . In a related paper, they also explore the algorithm’s suitability to show that dark matter is traced by baryonic matter across large-scale structure ([Chen et al. 2015b](#)). Additionally, [Chen et al. \(2016\)](#) provide a filament catalogue for SDSS, and show the influence of nearest-filament distances on galaxy properties like size, colour, and stellar mass.

Another example, this time in cosmology, is the application of the SCMS algorithm by [He et al. \(2017\)](#); they investigate non-Gaussianities of the matter density field to provide lensing effects based on filaments. [Hendel et al. \(2018\)](#) use the algorithm to identify morphological substructures in stellar debris, trying to classify the latter to gain a better understanding of galactic mergers through left-over collision disruptions.

The most recent use of the algorithm pertains to the field of quantitative criminology, where multiple extensions are introduced ([Moews et al. 2019](#)). The study investigates the optimisation of police patrols via density ridges, using publicly available data

from the City of Chicago over multiple years to assess the validity and stability of predictive ridges. Apart from thresholding as described above, the study makes use of the haversine formula as described by [Inman \(1835\)](#), which calculates the great-circle (or orthodromic) distance as a way to prevent distorted measures that would result from, for example, using the Euclidean distance. As a more specialised case of the law of haversines, it computes the distance between two points along the surface of a sphere. The formula makes use of the haversine function for a given angle  $\alpha$ ,

$$\text{hav}(\alpha) = \frac{1 - \cos(\alpha)}{2} = \sin^2\left(\frac{\alpha}{2}\right), \quad (7)$$

and can be used to calculate the relative haversine distance between two such points,  $\theta_1$  and  $\theta_2$ , as

$$\delta_{\text{hav}}(\theta_1, \theta_2) = \text{hav}(\theta_{2,1} - \theta_{1,1} + \cos \theta_{1,1} \cos \theta_{2,1} \text{hav}(\theta_{2,2} - \theta_{1,2})). \quad (8)$$

The resulting ridge estimation tool by [Moews et al. \(2019\)](#) is publicly available as a python package called DREDGE<sup>2</sup>, with a corresponding open-source repository<sup>3</sup>. In this paper, we adapt and extend DREDGE in order to apply it to DES Year 1 mass density maps. In addition to the existing thresholding and integrated use of the haversine formula, we parallelise time-consuming parts of the code, enable the manual setting of the mesh size for the ridge estimation, and implement a mathematically grounded optimisation of the required bandwidth. The latter is due to [Moews et al. \(2019\)](#) using a simple automatic bandwidth calculation that was previously proposed for, and is specific to, the field of criminology. Some of the features in the original implementation of DREDGE are redundant for this paper as well, most notably the ability to set a top-percentage threshold to extract only ridges falling within the highest-density areas, as researchers and practitioners in criminal justice are often interested in focussing on ‘hot spots’ ([Braga 2005](#)).

For the purpose of parallelising DREDGE, we make use of the embarrassing parallelism inherent in the updating function of the ridge points at each iteration. As this update is not reliant on other ridge points during the respective iteration, multiprocessing offers an easily accessible option to speed up the algorithm’s runtime.

### 2.2.2 Kernel density estimation

Kernel density estimation is a well-developed non-parametric data smoothing technique that has seen wide use in cosmological research applications. [Park et al. \(2007\)](#), for example, use an adaptive smoothing bandwidth with a spline kernel, and [Mateus et al. \(2007\)](#) apply a  $k$ -nearest neighbours density estimator to estimate the local number density of galaxies in an SDSS sample.

Similar methods have been employed in other surveys, for example by [Scoville et al. \(2007b\)](#), who identify large-scale structure and estimate dimensions, number of galaxies, and mass using an adaptive smoothing technique over a sample in the Cosmic Evolution Survey (COSMOS; [Scoville et al. 2007a](#)). Similarly, [Jang \(2006\)](#) uses a multivariate kernel density estimator with a cross-validated smoothing parameter to estimate galaxy cluster density over a sample of the Edinburgh-Durham Cluster Catalogue (EDCC).

The SCMS algorithm’s bandwidth  $\beta$ , plays a crucial role in determining the bias-variance relationship of the resulting distribution. A larger bandwidth results in a smoother distribution with

less variance and more bias, whereas a smaller bandwidth results in a less smooth distribution with more variance and less bias. Finding an optimal bandwidth for the kernel estimator in the context of DES mass maps is crucial to ensure that dense regions will not be oversmoothed, and higher-density areas of the projected large-scale structure will not be blurred into troughs. Conversely, optimising the bandwidth allows us to properly preserve the properties of low-density and extended structures.

In this application, we use a likelihood cross-validation approach to find the optimal bandwidth parameter. This method provides a density estimate that is close to the actual density in terms of the Kullback-Leibler divergence (KLD; [Kullback & Leibler 1951](#), a measure of relative entropy). Cross-validation is performed using a maximum likelihood estimation of the leave-one-out kernel estimator of  $f_{-i}$ , which is given by

$$f_{-i}(X_i) = \frac{1}{(n-1)h} \sum_{j=1, j \neq i} K_h(X_i, X_j), \quad (9)$$

where  $h$  is the bandwidth parameter and  $K_h$  represents the generalised product kernel estimator.

We use the generalised product kernel estimator on the latitude and longitude coordinates of the DES Y1 data as described by [Li & Racine \(2006\)](#), as

$$K_h(X_i, X_j) = \prod_{s=1}^q h_s^{-1} k\left(\frac{X_{is} - X_{js}}{h_s}\right), \quad (10)$$

where  $q$  is the dimension of  $X_i$ ,  $X_{is}$  is the  $s^{\text{th}}$  component of  $X_i$  ( $s = 1, \dots, q$ ),  $h_s$  is the smoothing parameter for the given component of  $X_i$ , and  $k(\cdot)$  is a univariate kernel function. We select this nonparametric kernel estimator as it doesn’t assume any functional form of the data, only that it satisfies regularity conditions such as smoothness and differentiability. In our case, the  $\beta$  value we use in SCMS is the best  $h$  value obtained by cross-validation, while the variables  $X_i, X_j$  correspond to the sky positions  $\theta$ .

## 2.3 Data and simulations

### 2.3.1 DES mass maps

The Dark Energy Survey (DES) is a six-year photometric survey project to image 5000 square degrees of the sky in *grizY* filters using the DECam camera on the Blanco telescope, Cerro Tololo, Chile ([Flaugher et al. 2015](#)). The primary purpose of the survey is to generate a dataset for cosmology, and in particular one suitable for weak gravitational lensing measurements. DES observations completed in January 2019, and data analysis for the project is ongoing.

Weak lensing measurements use galaxies as a backlight to determine the projected gravitational fields along the paths of their light rays. Gravity bends the light paths, resulting in a shearing of galaxy images that can be measured from the mean ellipticity of a large-enough sample of objects. One application of weak lensing is to generate mass maps. Assuming the general relativity relationship between gravitational convergence  $\kappa$  and mass, we can use ellipticity catalogues to map the projected, weighted overdensity in a given pixel of the survey. The weighting depends on the redshift distribution of the observed galaxies and the redshift-distance relationship which, in turn, depends on the underlying cosmology.

The DES Year 1 (Y1) data release<sup>4</sup>, as described by [Drlica-](#)

<sup>2</sup> <https://pypi.org/project/dredge>

<sup>3</sup> <https://github.com/moews/dredge>

<sup>4</sup> <https://des.ncsa.illinois.edu/releases/y1a1>



Wagner et al. (2018), includes 2D-projected mass maps (see Chang et al. 2018) estimated from cosmic shear measurements from the survey’s first year (Zuntz et al. 2018). The creation of these maps is based on catalogues made using the METACALIBRATION method (see Huff & Mandelbaum 2017; Sheldon & Huff 2017) and the redshifts estimated in Hoyle et al. (2018). They were generated using the spherical Kaiser-Squires method on the galaxy shear catalogues to invert shears to convergences (Kaiser & Squires 1993; Schneider 1996).

Various mass maps were made using different selections of source galaxies, thus having different redshift weight functions. For this initial project, we use only the maps made with the widest range of galaxies, from  $z = 0.2$  to  $z = 1.3$ . Here, we use the E-mode maps and their corresponding masks<sup>5</sup>.

### 2.3.2 Sample generation from DES maps

Although the mass maps are already in the form of a field on which Hessians and gradients could be calculated, in practice, the DES Y1 mask, which has a large number of small excised regions at the locations of bright stars, makes calculating derivatives a very noisy process, even with aggressive masks. It is considerably simpler to instead generate samples from the map and apply the SCMS algorithm described above. In order to generate these samples, we compute a mean value  $\mu_i$  per pixel,

$$\mu_i = \max(1 + \omega \kappa_i, 0), \quad (11)$$

where  $\kappa_i$  is the projected overdensity in pixel  $i$ , and  $\omega$  is a parameter that we can tune as desired. If  $\omega$  is too low, the ridges will not be detectable in the map, and if it is too high, the lower ridges will disappear because the highest density peaks dominate. We find that a value of  $\omega = 50$  works well for DES data to convert the  $\sim 2\%$  map fluctuations to  $\mathcal{O}(1)$  point density variation. We then generate a number  $n_i$  of samples per pixel, using a Poisson distribution  $n_i \sim \text{Poi}(\mu_i)$ , and place these samples uniformly within the pixel.

Note that no interpolation or reconstruction is done in masked regions; we simply omit masked pixels and do not generate samples inside them. This method is simple to apply and causes no numerical difficulties, but it does mean that ridges at the edge of the mask regions may not be detected. This occurs in regions where interpolation would also fail. Methods using these ridge catalogs should account for this, for example in simulations.

### 2.3.3 Flask simulations

To provide a testbed for our method before employing real data, we also build a suite of simulated maps using the FLASK<sup>6</sup> software (see Xavier et al. 2016a,b), which generates tomographic log-normal random fields that approximate large-scale structure distributions. We use the Planck best-fit  $\Lambda$ CDM cosmological parameters in the code,  $\sigma_8 = 0.25$ , and the same redshift distribution as estimated for the DES source galaxies (see Hoyle et al. 2018), normalised to the correct overall density. We then generate true  $\kappa$  maps, which we treat as idealised noiseless simulations, and galaxy ellipticity catalogues, which we use with a spherical Kaiser-Squires map-making method

<sup>5</sup> [http://desdr-server.ncsa.illinois.edu/despublic/y1a1\\_files/mass\\_maps\\_files/y1a1\\_spt\\_mcal\\_0.2\\_1.3\\_kE.fits](http://desdr-server.ncsa.illinois.edu/despublic/y1a1_files/mass_maps_files/y1a1_spt_mcal_0.2_1.3_kE.fits) and [y1a1\\_spt\\_mcal\\_0.2\\_1.3\\_mask.fits](http://desdr-server.ncsa.illinois.edu/despublic/y1a1_spt_mcal_0.2_1.3_mask.fits)

<sup>6</sup> <https://github.com/hsxavier/flask>

to generate a noisy  $\kappa$  map. Lastly, we apply the DES masks to both noisy and noiseless simulated maps.

While this approach works well for the experiments performed in this work, it is insufficient for pseudo-3D extensions of our method discussed in Section 4. This is due to flask not generating features like clusters and filaments.

## 3 EXPERIMENTAL RESULTS

In this section, we discuss and implement a distance-based statistical test to verify our results and test the degree of robustness to noise through noisy and noiseless simulated dark matter density maps in Section 3.1. We explore the properties and specificities of our method’s tracing of the large-scale structure through a quantifiable comparison to the curvelet transforms in Section 3.2. Lastly, we present the extracted ridges, together with a comparison to previous research on trough identification from DES Y1 data, in Section 3.3.

### 3.1 Statistical functionality verification

As is common with studies dealing with cosmic voids, validating our approach is not an easy task, even when we can apply it to simulations. Indeed, while simulations allow us to apply our method to noiseless data, they do not provide us with ‘ground-truth’ ridges or troughs, since these need to be estimated and defined, even in the absence of noise. A viable analysis is the comparison of ridges and troughs recovered when running the proposed approach in two different settings, applying it to both a noiseless simulation and one that contains realistic levels of noise, as described in Section 2.3.3. In other words, we can test for robustness to observational noise. This test, however, requires the choice of a similarity criterion, or distance metric, between sets of ridges.

Optimal transport theory (see Villani 2008), and specifically the Wasserstein distance, provides us with a natural, principled means of computing such distances. Per the Monge-Kantorovich interpretation of optimal transport, the Wasserstein distance can be understood as the minimal possible cost incurred to move a certain amount of mass from one distribution to another. This is precisely what a set of ridges and troughs represent; a certain distribution of mass, projected in two dimensions across (a part of) the sky.

Consider two distributions of mass,  $p_1, p_2 \in \mathbb{R}^N$ , sampled on some discrete space of dimension  $N$ . Let  $C \in \mathbb{R}^{N \times N}$  be the cost matrix for which each entry  $C_{ij}$  contains the cost of moving mass from position  $i$  to position  $j$ . We define the Wasserstein distance as

$$W(p_1, p_2) = \arg\min_{T \in \Pi(p_1, p_2)} \langle T, C \rangle. \quad (12)$$

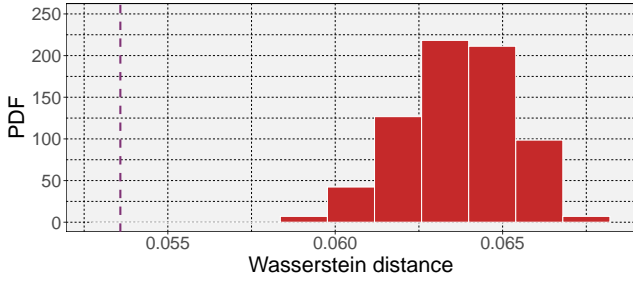
Here,  $\Pi(p_1, p_2)$  is the set of transport plans between  $p_1$  and  $p_2$ . For any such matrix  $T$ , each of its entries  $T_{ij}$  contains the amount of mass of  $p_1$  that is transported from position  $i$  within  $p_1$  (denoted  $p_{1,i}$ ) to position  $j$  within  $p_2$ . By construction, for any row  $i$ , we then have

$$\sum_j T_{ij} = p_{1,i}. \quad (13)$$

Similarly, summing over the columns of  $T$  yields the entries of  $p_2$ . We can then simply express the set of acceptable transport plans as

$$\Pi(p_1, p_2) = \left\{ T \in \mathbb{R}^{N \times N}, \forall i, j, \sum_j T_{ij} = p_{1,i}, \sum_i T_{ij} = p_{2,j} \right\}. \quad (14)$$

In order to solve Eqn. (12), meaning to find the optimal  $T$ , we use the



**Figure 1.** Wasserstein distance between the ridges obtained on the noisy simulation and either its noiseless counterpart, as a vertical dashed line in purple, or a set of 101 random distributions of mass in red.

entropic regularisation scheme proposed by Cuturi (2013), which allows for the Wasserstein distance between discrete measures to be computed by using an iterative scheme (see the textbook by Peyré & Cuturi 2019, for a recent overview of computational schemes for the practical computation of optimal transport quantities).

The raw output of SCMS is a set of shifted ‘mesh points’, each with a 2D position vector, that make up the ridges. While we could compute the Wasserstein distance directly between two sets of mesh points, the objects of interest in our case are the ridges they comprise, and the troughs thus delimited, rather than the points themselves. For this reason, we convert each set of ridges into a binary 2D image with each pixel’s value set to zero if no meshpoint fell within it, and one otherwise. We then compute the Wasserstein distance between the two resulting images. Here,  $N$  is equal to the total number of pixels, in our case 10952, and the cost matrix  $C$  is the Euclidean distance between each pair of pixel positions.

In order to provide a basis for comparison, we generate random maps by projecting the DES mask onto the image plane and uniformly sampling the same number of non-zero pixels as those present in the images that contain our ridges. We then compute the Wasserstein distance between those ridges and the ones obtained from the noisy simulation. The distribution of these distances are represented in Figure 1, along with the distance computed between the two sets of ridges. As can be seen, our ridge-finding method shows robustness to realistic amounts of noise. This test allows for a general confirmation that the ridges we obtain contain physical information about the distribution of matter, obviating any specification of the problem under study, that is, a precise definition of troughs or voids. Further testing tailored to the application of interest, however, is recommended. These tests could likely use the same settings and simulations, where we compare the output in both noiseless and noisy cases, for example the cosmological constraints derived from the troughs delimited by those two sets of ridges).

### 3.2 DES ridges and curvelet comparison

Sparse signal processing (see Starck et al. 2015) provides solutions to many signal retrieval problems including, but not exclusive to, image denoising. The approach relies on finding a representation space in which the signal can be sparsely represented, a typical example being a sinusoidal signal, which can be fully expressed with very few non-zero coefficients in Fourier space. While natural signals are rarely sinusoidal, wavelets are commonly used as a sparse basis of representation (Mallat 1999). They can, however, perform poorly when the features to be recovered are rectilinear or elongated, as is the case for estimating ridges. Because of this

shortcoming, analogous transforms have been designed specifically for these cases, namely ridgelets and curvelets (Candès & Donoho 1999; Candès et al. 2006).

These transforms have led to several applications in astrophysics and cosmology; a wide range of examples are shown by Starck et al. (2003). Starck et al. (2004) use wavelets, ridgelets and curvelets to detect and characterise, on simulations, various sources of CMB anisotropies, which include imprints of inflation, the Sunyaev-Zel’dovich effect, and cosmic strings. The latter have also been studied in a CMB framework by Vafaei Sadr et al. (2017) and Hergt et al. (2017), using the curvelet transform, while Laliberté et al. (2018) use ridgelets to that end within N-body simulations. Gallagher et al. (2011) apply both to solar astrophysics, and Jiang et al. (2019) use curvelets for radio transient detection.

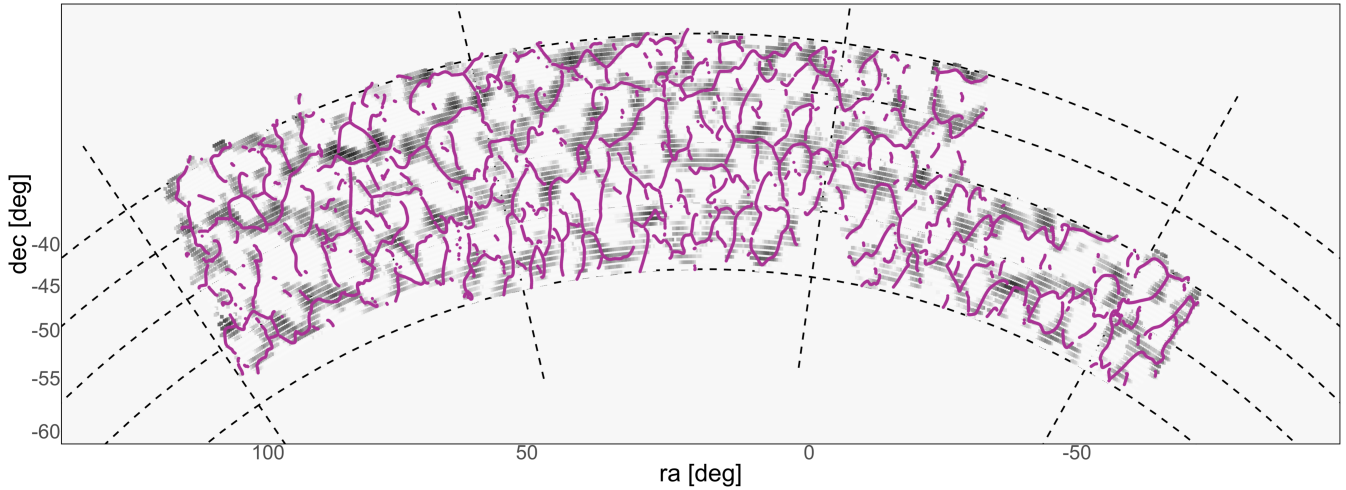
In our case, the application of curvelets allows for a straightforward and entirely independent means of recovering the ridges. We perform a simple denoising, that is, a thresholding of the input mass in curvelet space. Our SCMS-recovered ridges are obtained from the samples described in Section 2.3.2, which themselves rely on the choice of the  $\omega$  parameter. We use these samples to generate a two-dimensional, discrete image by counting their number in each pixel bin. We then apply curvelet denoising to the resulting image, using the freely available Sparse2D package<sup>7</sup>. The thresholds are selected using the False Discovery Rate approach described by Benjamini & Hochberg (1995), and the coarse scale is discarded.

The resulting denoised map, converted back to sky coordinates, is shown in Figure 2. Upon visual inspection, we find an overall good agreement between the uncovered structured overdensities and the ridges obtained by SCMS, overplotted in purple. To get a more quantitative view of the differences in the structures recovered by both approaches, we project the ridges yielded by the SCMS algorithm into the same pixel grid as that used by the curvelet denoising step. For every pixel that then contains part of one of the ridges, we check the corresponding value in the curvelet reconstruction. We consider all pixels where that value is 0 to be a mismatch. Figure 3 shows all such mismatches in orange, while parts of the ridges that match with the curvelet reconstruction are shown in purple.

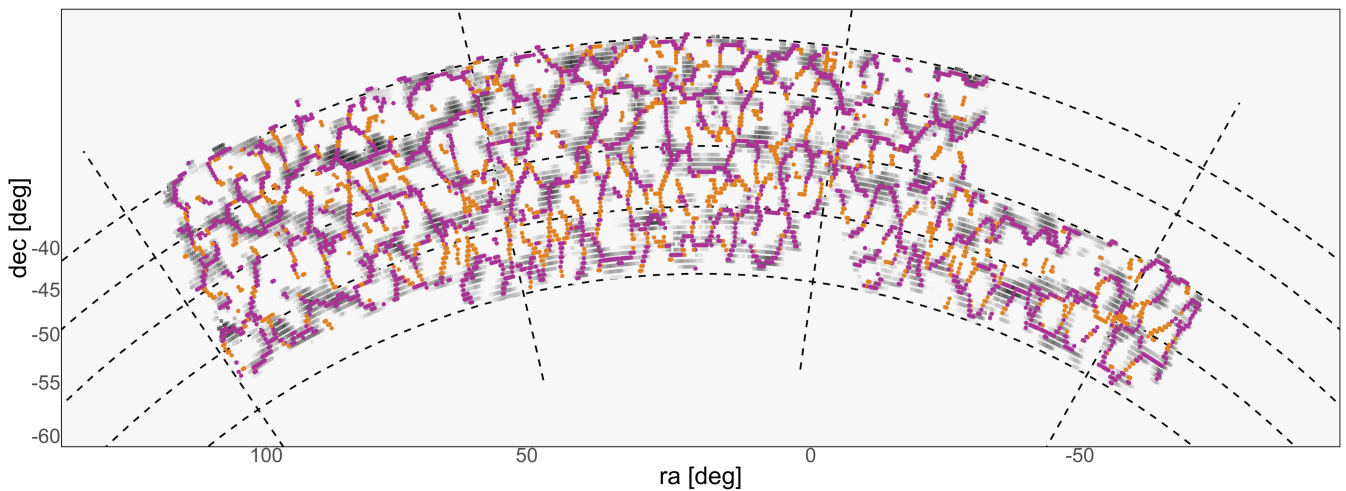
About 31% of the SCMS-derived ridges mismatch the curvelet reconstruction. As can be seen from Figure 3, these mostly correspond to areas where structure is present in two nearby areas of the sky, and where the curvelet reconstruction keeps those two areas disjoint, while SCMS-derived ridges link them together. Although very different in their heuristics, ultimately, both approaches perform some form of denoising of the input mass maps. Those differences between the two denoising approaches could indicate that one of the two fails at properly separating noise from signal. The areas of disagreement could either be due to the curvelets considering signal to be noise, or to SCMS algorithm turning noise into ridges. In order to clarify this issue, we reprocess the DES data while tuning the parameters of each method. In the case of SCMS, we impose a stronger denoising by multiplying the bandwidth,  $\beta$ , obtained as described in Section 2.2.2, by a factor  $\xi \in \{1.25, 1.50, 2.00\}$ , or considering a higher threshold  $\Upsilon > \tau$  (see Section 2.2.1). The percentage of mismatching ridges for both cases is shown in Figure 4.

Artificially increasing the strength of the denoising performed by the SCMS algorithm does not lead to ridges that match the curvelet reconstruction more closely. In fact, the percentage of mismatch between the two approaches is not even monotonous in the case of bandwidth changes. In all cases, we still observe ridges

<sup>7</sup> <https://github.com/CosmoStat/Sparse2D>



**Figure 2.** Comparison of density ridges and a curvelet reconstruction. Ridges in purple are superimposed on structural constraints obtained via curvelet denoising in shades of orange, with higher densities shifting from lighter to darker. Both results are based on DES Y1 weak lensing mass density maps.



**Figure 3.** Similar to Figure 2, but with ridges shown in purple where they match the curvelet reconstruction, and orange otherwise.

that tend to be more connected in the SCMS case. Similarly, we try imposing weaker denoising in the curvelet approach, by using increasingly lower values of  $k$  in  $k\hat{\sigma}$  thresholds instead of the False Detection Rate approach, where  $\hat{\sigma}$  is estimated from the data using the Median Absolute Deviation estimator. Once again, this yields no clear increase in the match between the two outputs.

This shows that the support of each method, in their respective (hyper)parameter spaces, are disjoint. In other words, the differences between curvelet reconstruction and SCMS-derived ridges seen in Figure 3 are due to intrinsic differences between the approaches, as opposed to a failure of either at the denoising task. Implications of those differences for potential applications will be further discussed in Section 4.

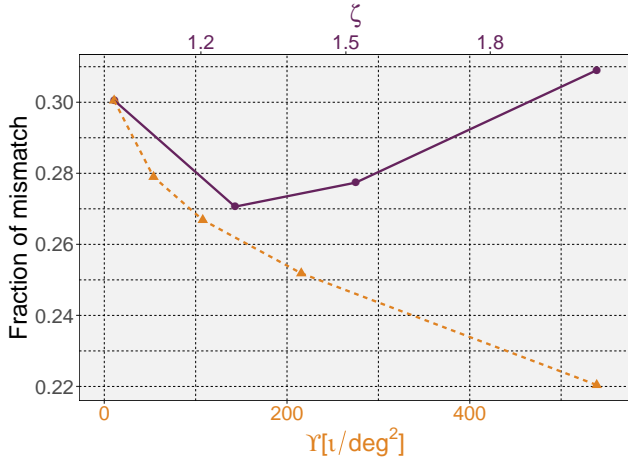
### 3.3 Ridges in the Dark Energy Survey

After the above comparison with curvelet transforms, we perform a second comparison of our results with independent methodology, as

void and trough detection remain a current focus of interest within the cosmology community, as described in Section 1. Specifically, Gruen et al. (2018) derive cosmological constraints via density split statistics, counting tracer cells to split lines of sight and measuring counts-in-cells and gravitational shear in regions of varying density. Using luminous red galaxies (REDMAGIC galaxies at  $0.2 < z < 0.45$ ) to trace the foreground matter density field, they count the number of galaxies that fall within circular top-hat apertures within radii  $\theta_T = [10', 20', 30', 60']$ .

Each line of sight is then assigned to a density quintile based on those counts, the data for which is publicly available<sup>8</sup>, with the highest quintile being of interest to us in terms of high-density ridge analogues. Due to the DES Y1 data being more inhomogeneous in depth than the SDSS DR8 data used for the REDMAGIC

<sup>8</sup> <https://des.ncsa.illinois.edu/releases/y1a1/density>, file `trough_maps.tar.gz`



**Figure 4.** Fraction of ‘mismatch’ between SCMS-derived ridges and the curvelet reconstruction, when varying the  $\tau$  and  $\beta$  hyperparameters of the SCMS algorithm. The bottom axis shows the threshold on the ridges ( $\Upsilon$ ) in units of meshpoints ( $l$ ) per square degree ( $\text{deg}^2$ ). The top axis shows the multiplication factor ( $\zeta$ ) of the bandwidth. The further to the right, the stronger the implicit denoising performed by the algorithm. The absence of a clear decreasing trend shows that the differences between both methods, as illustrated by Figure 3, are not due to hyperparameter choices.

catalogue (see Aihara et al. 2011), tracer galaxies are removed and the area is defined as fully masked if the sample from the catalogue is not complete to  $z = 0.45$ .

Given that Gruen et al. (2018) use the same DES Y1 data we make use of in this work, in combination with SDSS DR8 data, this enables us to compare our findings against results from the DES collaboration without being based on the same underlying mass density map as the curvelet comparison. Section 3.3 shows the same DES Y1-extracted ridges as in Figure 2, overlaid with the highest-quintile density measurements from Gruen et al. (2018).

While these quintile-based comparisons are more ‘spotty’ due to the masking based on depth completeness described above, and qualitatively different when compared to curvilinear structures extracted by the SCMS algorithm, the lines generally trace the same high-density regions shown in the figure. The shown percentiles also use foreground galaxies as tracers of the matter field, as opposed to the mass density maps based on weak lensing that are used as the input for our modified version of DREDGE. In contrast to the previous comparison with curvelets as an alternative method in Section 3.2, the goal of this experiment is to perform a comparison across both methods and underlying data to reach a consensus in terms of high-density regions on the sky.

## 4 DISCUSSION

### 4.1 Overview

Ridges derived from our methodology have several properties that make them interesting for lensing trough studies, compared to simply using local minima. First, since the method generates a point cloud as a starting point, rather than working with the map directly, it does not require that we smooth the convergence map to use it. This means that smaller-scale troughs can potentially be probed. Secondly, since the trough points maximise the distance from local ridge structures, they probe inter-cluster and inter-filament regions directly, rather than by proxy. Finally, masks from bright stars on

small scales (below the typical ridge segment length) do not significantly affect the operation of the algorithm.

We incorporate several extensions from previous research into our implementation of the SCMS algorithm, both for accuracy and performance reasons, and describe these in Section 2.2. Thresholding the initial mesh of points based on density estimates, as introduced by Chen et al. (2015a), ensures that the curvilinear structure is constrained to higher-density areas, which solves potential issues associated with identifying ridges based on sparsely populated regions in the data. In addition, we make use of the haversine formula, as implemented by Moews et al. (2019) for geospatial analysis, to prevent distorted measures on the curved sky.

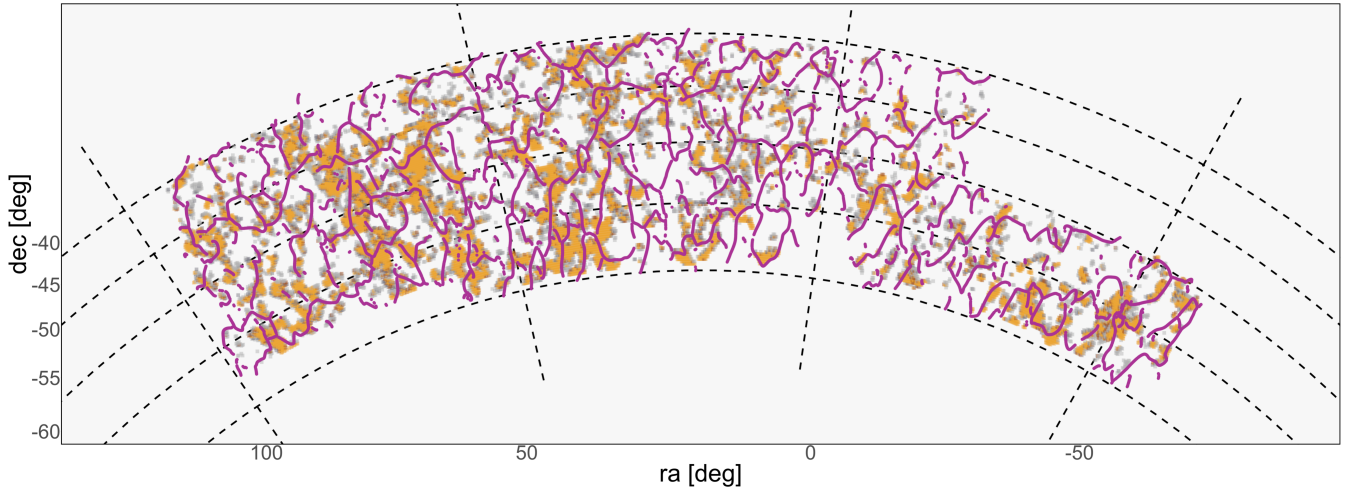
We further extend the algorithm by exploiting the potential for embarrassing parallelism inherent in the updating function of ridge points at each iteration. This allows users to reduce the runtime significantly by using a multiprocessing setup, thus removing a major obstacle when applying the SCMS algorithm to large datasets. The algorithm itself is reliant on the choice of a bandwidth, which plays a crucial role in determining the bias-variance relationship of the distribution, with larger bandwidths resulting in smoother distributions with less variance and more bias, and with smaller bandwidths resulting in a less-smooth distribution with more variance and less bias. For our current application, this means that large-scale structure in dark matter density maps could be blurred into cosmic troughs through bandwidths that are too large, while recovered ridges could show spurious fine-grained structure through bandwidths that are too small. We solve this by introducing a likelihood cross-validation to the SCMS algorithm to automatically find the optimal bandwidth in a data-driven way, providing a density estimate close to the actual density in terms of the Kullback-Leibler divergence.

The experiments performed in the course of this work are designed to test for both noise robustness and the correct tracing of high-density regions. In terms of noise, we generate noisy and noiseless simulations that correspond to the DES Y1 mass density maps in Section 3.1, and calculate the Wasserstein distance between binarised maps of the resulting curvilinear structures. In addition, we calculate the same metric for randomly generated maps with uniform sampling to place the results into context, showing robustness to realistic levels of noise. While this test provides some measure of the robustness of our method to the noise present in real data, as discussed at the end of Section 3.1, additional testing is recommended to avoid introducing biases specific to any application. However, as our goal is to propose a general approach to recover ridges, our tests were chosen to be as general as possible.

The overall agreement we find between our ridges and the structure recovered by curvelet denoising in Section 3.2, despite the independence and vast difference between the two approaches, is a good indication that our ridges successfully capture information contained in the matter distribution. The differences between the two methods lie mostly in the curvelet reconstruction leading to more disconnected patches, while our ridges tend to connect these areas. Our experiments show that this is the result of intrinsic differences between both approaches, and not due to a poor choice of hyperparameters. While clearly different, neither of the two resulting estimates are ‘wrong’. This illustrates an important point: For a given matter distribution, there is no such thing as a set of true ridges. Choosing a proper definition of such ridges is a non-trivial task in itself, which is precisely one of the motivations for the present work.

In this paper, where we aim to present both methods in as general a framework as possible, we will simply point out the difference





**Figure 5.** Comparison of density ridges and previous results from Gruen et al. (2018). Ridges from this work are shown in purple, and are superimposed on mass density probabilities that were obtained by measuring counts-in-cells along lines of sight of the foreground luminous red galaxies REDMAGiC sample derived by .

in the structure they reveal, highlighting once again that neither is more correct than the other. This could change in the context of any given application; tests tailored to a specific application could be used to determine which of the two is more appropriate or useful in that context. It should also be noted that a key difference between the two approaches is that the SCMS algorithm does indeed produce ridges, that is, curvilinear structures. The output of the curvelet denoising, on the other hand, is still a full two-dimensional map. The curvelet reconstruction also still contains information about the amplitude of the signal at each position, while our ridges are binary, meaning that every position either contains a ridge or not. If such a truly curvilinear format is required for a given application, and the patchier nature of the curvelet-recovered ‘ridges’ is better suited for the respective task at hand, it would be straightforward to combine both approaches (see Figure 3).

To further our analysis, we compare ridges extracted from DES Y1 mass density maps, which are based on weak lensing, to trough mass probabilities by Gruen et al. (2018) in Section 3.3, using the highest quintile corresponding to large-scale structure. Since the data is based not on weak lensing, but on luminous red galaxies (REDMAGiC galaxies at  $0.2 < z < 0.45$ ) to trace the foreground matter density field, and is also limited due to combining DES Y1 and SDSS DR8 data in terms of their depth coverage, this serves as a comparison across both methods and sources. Given this combination in Gruen et al. (2018), a stronger masking is present, leading to a more ‘spotty’ nature of the visible structure when compared to the curvelet representation in Section 3.2. Despite these differences, we see both ridges and quintile probabilities falling into the same areas, with DES Y1-extracted ridges tracing this complementary dataset.

Our cross-validated bandwidth optimisation approach to kernel density estimation provides a data-driven and, via the Kullback-Leibler divergence, mathematically motivated way to perform ridge estimation. As with any approach that is data-driven, sanity checks should be performed to ensure physically plausible ridges. These derived ridges are, therefore, quite dependent on both the data quality of the DES Y1 survey and the shape assumptions in the density estimation, as well as on the distribution space being approximately

symmetric and unimodal. In our approach, care should be taken when considering samples derived from heavy-tailed distributions, as the bandwidth optimisation can be prone to overestimation (Silverman 1986).

For an easier dissemination of our results, both for verification and further analysis of this denoising approach to dark matter density maps, we also release a catalogue<sup>9</sup> of the ridges extracted from DES Y1 data.

## 4.2 Applications

- consistency with standard models of structure formation in a specific environment - why interesting - intermediate regime between cluster and wide field lensing - less studied - screening mechanisms
- filaments a poorly understood - if we can understand how ridges relate to filaments

- 3d void catalog - why interesting - why this ridge method is useful to this - lensing along ridges - why interesting - why this ridge method is useful to this - In terms of future directions and follow-ups, we intend to extend our method to the three-dimensional case to identify full voids instead of trough projections. One way to reach this goal is to use tomographic 2D information, effectively identifying voids within a layered pseudo-3D structure. This can be performed using tomographic reconstruction techniques (see Herman 2009), which have been highly enhanced using deep neural networks (Wang et al. 2018).

This planned future application would yield a three-dimensional void catalogue, the statistical properties of which could constitute powerful cosmological probes. To start, the higher-order statistics of the spatial distribution of voids could be compared to results from Hamaus et al. (2014) and Pycke & Russell (2016). Additional insights could also be gained from the spatial correlation of voids as a function of their sizes and shapes.

Another interesting avenue we plan to pursue is the inclusion of lensing information along identified ridges to further bolster the

<sup>9</sup> Link will be added here after acceptance

viability of this approach, as well as the matching of a cosmic web reconstructed from DES data to Planck LSS (Bouchet 2016).

Just as the statistical characterisation of the triaxiality of galaxy clusters has proven fruitful for cosmological analysis (see Simet et al. 2017; Melchior et al. 2017; Chiu et al. 2018), so too might the triaxiality distribution of voids. One could further imagine a hierarchical inference procedure that uses the void size and shape distributions as a prior to iteratively update ambiguous photometric redshift probability density functions, just as those same probability density functions are used to hierarchically infer the void size and shape distributions from the tomographic projections used to derive the triaxial void catalogue. A supplemental application of these distributions could serve as a validation test for the sophisticated simulated catalogues being developed for future LSS surveys (Korytov et al. 2019).

In short, this method brings within reach a number of promising avenues for testing cosmological models, developing novel cosmological probes, and validating cosmological analysis pipelines. In this paper, we demonstrate the utility of curvilinear structures for the denoising of large-scale structure based on weak lensing, extending the available methodology for this purpose. We plan to explore these applications in future papers, but also welcome potential users of the ridge catalogue to experiment with the dataset we provide alongside the presented results.

## 5 CONCLUSION

This work presents a new ridge estimation approach based on an extension of the subspace-constrained mean shift algorithm, a filamentary search method, and releases the corresponding results as a catalogue of curvilinear structure. As an application case of current relevance, we apply the method to dark matter mass density maps from the DES Y1 data release to extract high-density ridges between cosmic troughs. Our results demonstrate the viability of ridge estimation as a precursory step for denoising cosmic filaments, leading to a versatile and effective identification of cosmic troughs.

We extend the SCMS algorithm by including the haversine distance, a customisation of the mesh size for ridge estimation, automatic optimisation of the bandwidth used in the process, and the parallelisation of the updating function for mesh points to scale down the algorithm's runtime. In addition, we also include the thresholding extension of the SCMS algorithm from a previous application to astronomical data.

In order to test the robustness of our method, we recover the ridges from simulations under different noise levels, and use the Wasserstein distance as a comparison metric. We further compare our extracted ridges, which are based on DES Y1 weak lensing data, with curvelet denoising of the same data, and with high-density quintiles derived from both DES Y1 and SDSS DR8 foreground galaxies limited by inhomogeneous depth coverage. This allows us to compare results across both methods and data sources, leading to highly reasonable agreement.

## ACKNOWLEDGEMENTS

This work was developed during the 6<sup>th</sup> COIN Residence Program<sup>10</sup> (CRP#6) held in Chamonix, France in August 2019. This

project is financially supported by CNRS as part of its MOMENTUM programme over the 2018–2020 period. The Cosmostatistics Initiative<sup>11</sup> (COIN) is a non-profit organization whose aim is to nourish the synergy between astrophysics, cosmology, statistics, and machine learning communities. AKM acknowledges the support from the Portuguese Fundação para a Ciência e a Tecnologia (FCT) through grants SFRH/BPD/74697/2010, PTDC/FIS-AST/31546/2017 and from the Portuguese Strategic Programme UID/FIS/00099/2013 for CENTRA.

## REFERENCES

- Adermann E., Elahi P. J., Lewis G. F., Power C., 2018, *MNRAS*, 479, 4861  
 Aihara H., et al., 2011, *ApJS*, 193, 29  
 Aragón-Calvo M. A., Jones B. J. T., van de Weygaert R., van der Hulst J. M., 2007, *A&A*, 474, 315  
 Aragón-Calvo M. A., Platen E., van de Weygaert R., Szalay A. S., 2010, *ApJ*, 723, 364  
 Barreira A., Cautun M., Li B., Baugh C. M., Pascoli S., 2015, *J. Cosmology Astropart. Phys.*, 2015, 028  
 Bas E., Erdogmus D., 2011, *Neuroinformatics*, 9, 181  
 Benjamini Y., Hochberg Y., 1995, *J. Royal Stat. Soc. B*, 57, 289  
 Bond J. R., Kofman L., Pogogyan D., 1996, *Nature*, 380, 603  
 Bonjean V., Aghanim N., Salomé P., Douspis M., Beelen A., 2018, *A&A*, 609, A49  
 Bouchet F., 2016, in *Frontiers of Fundamental Physics 14*. Berlin: Springer, doi:10.22323/1.224.0002  
 Boylan-Kolchin M., Springel V., White S. D. M., Jenkins A., Lemson G., 2009, *MNRAS*, 398, 1150  
 Braga A. A., 2005, *J. Exp. Criminol.*, 1, 317  
 Brouwer M. M., et al., 2018, *MNRAS*, 481, 5189  
 Cai Y.-C., Padilla N., Li B., 2015, *MNRAS*, 451, 1036  
 Candès E. J., Donoho D. L., 1999, *Philos. Trans. R. Soc. A*, 357, 2495  
 Candès E., Demanet L., Donoho D., Ying L., 2006, *Mult. Mod. Sim.*, 5, 861  
 Chambers K. C., et al., 2016, preprint, (arXiv:1612.05560)  
 Chang C., et al., 2018, *MNRAS*, 475, 3165  
 Chen Y.-C., Ho S., Freeman P. E., Genovese C. R., Wasserman L., 2015a, *MNRAS*, 454, 1140  
 Chen Y.-C., et al., 2015b, *MNRAS*, 454, 3341  
 Chen Y.-C., Ho S., Brinkmann J., Freeman P. E. P., Wasserman L., 2016, *MNRAS*, 461, 3896  
 Chiu I.-N., Umetsu K., Sereno M., Ettori S., Meneghetti M., Merten J., Sayers J., Zittrn A., 2018, *ApJ*, 860, 126  
 Clampitt J., Cai Y.-C., Li B., 2013, *MNRAS*, 431, 749  
 Cuturi M., 2013, in *Advances in Neural Information Processing Systems*. Curran Associates Inc., p. 2292  
 Dekel A., Rees M. J., 1994, *ApJ*, 422, L1  
 Dietrich J. P., Werner N., Clowe D., Finoguenov A., Kitching T., Miller L., Simionescu A., 2012, *Nature*, 487, 202  
 Drlica-Wagner A., et al., 2018, *ApJS*, 235, 33  
 El-Ad H., Piran T., 1997, *ApJ*, 491, 421  
 Flaugher B., et al., 2015, *AJ*, 150, 150  
 Fry J. N., 1986, *ApJ*, 306, 358  
 Gaité J., 2005, *Eur. Phys. J. B*, 47, 93  
 Galárraga-Espinosa D., Aghanim N., Langer M., Gouin C., Malavasi N., 2020, arXiv e-prints,  
 Gallagher P. T., Young C. A., Byrne J. P., McAteer R. T. J., 2011, *Adv. Space Res.*, 47, 2118  
 Genovese C. R., Perone-Pacifico M., Verdini I., Wasserman L., 2014, *Ann. Stat.*, 42, 1511  
 Govoni F., et al., 2019, *Science*, 364, 981  
 Gruen D., et al., 2016, *MNRAS*, 455, 3367  
 Gruen D., et al., 2018, *Phys. Rev. D*, 98, 023507

<sup>10</sup> <https://cosmostatistics-initiative.org/residence-programs/crp6>

<sup>11</sup> <https://cosmostatistics-initiative.org>

- Hamaus N., Sutter P. M., Wandelt B. D., 2014, *Proc. Int. Astron. Union*, 11, 538
- Hamaus N., Pisani A., Sutter P. M., Lavaux G., Escoffier S., Wandelt B. D., Weller J., 2016, *Phys. Rev. Lett.*, 117, 091302
- He S., Alam S., Ferraro S., Chen Y.-C., Ho S., 2017, *Nat. Astron.*, 2, 401
- Hendel D., Johnston K. V., Patra R. K., Sen B., 2018, preprint ([arXiv:1811.10613](https://arxiv.org/abs/1811.10613))
- Hergt L., Amara A., Brandenberger R., Kacprzak T., Refregier A., 2017, *J. Cosmology Astropart. Phys.*, 2017, 004
- Herman G. T., 2009, *Fundamentals of computerized tomography: Image reconstruction from projections*. Berlin: Springer
- Hoyle F., Vogeley M. S., 2004, *ApJ*, 607, 751
- Hoyle B., et al., 2018, *MNRAS*, 478, 592
- Huff E., Mandelbaum R., 2017, preprint, ([arXiv:1702.02600](https://arxiv.org/abs/1702.02600))
- Inman J. W., 1835, *Navigation and nautical astronomy for the use of British seamen*, 3 edn. London: W. Woodward, C. & J. Rivington
- Ivezic Z., et al., 2008, preprint ([arXiv:0805.2366](https://arxiv.org/abs/0805.2366))
- Jang W., 2006, *Comput. Stat. Data Anal.*, 50, 760
- Jauzac M., et al., 2012, *MNRAS*, 426, 3369
- Jiang M., Cui B., Yu Y.-F., Cao Z., 2019, *IEEE Access*, 7, 107389
- Kaiser N., Squires G., 1993, *ApJ*, 404, 441
- Klypin A. A., Trujillo-Gomez S., Primack J., 2011, *ApJ*, 740, 102
- Korytov D., et al., 2019, *ApJS*, 245, 26
- Kullback S., Leibler R. A., 1951, *Ann. Math. Stat.*, 22, 79
- Laliberte S., Brandenberger R., Camargo Neves da Cunha D., 2018, arXiv e-prints,
- Lavaux G., Wandelt B. D., 2012, *ApJ*, 754, 109
- Li B., 2011, *MNRAS*, 411, 2615
- Li Q., Racine J. S., 2006, *Nonparametric Econometrics: Theory and Practice*. No. 8355 in Economics Books, Princeton, NJ: Princeton University Press
- Libeskind N. I., et al., 2018, *MNRAS*, 473, 1195
- Malavasi N., Aghanim N., Douspis M., Tanimura H., Bonjean V., 2020, arXiv e-prints,
- Mallat S., 1999, *A Wavelet Tour of Signal Processing - The Sparse Way*. Amsterdam: Elsevier
- Mateus A., Sodré L., Cid Fernandes R., Stasińska G., 2007, *MNRAS*, 374, 1457
- Maturi M., Merten J., 2013, *A&A*, 559, A112
- Mead J. M., King L. J., McCarthy I. G., 2010, *MNRAS*, 401, 2257
- Melchior P., et al., 2017, *MNRAS*, 469, 4899
- Miao Z., Wang B., Shi W., Wu H., 2014, *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, 7, 4762
- Moews B., Argueta Jaime R. J., Gieschen A., 2019, preprint (ArXiv:1907.03206)
- Nadathur S., Hotchkiss S., Crittenden R., 2017, *MNRAS*, 467, 4067
- Neyrinck M. C., 2008, *MNRAS*, 386, 2101
- Ozertem U., Erdogmus D., 2011, *J. Mach. Learn. Res.*, 12, 1249
- Park C., Choi Y.-Y., Vogeley M. S., Gott J. Richard I., Blanton M. R., SDSS Collaboration 2007, *ApJ*, 658, 898
- Peebles P. J. E., 2001, *ApJ*, 557, 495
- Peyré G., Cuturi M., 2019, *Found. Trends Mach. Learn.*, 11, 355
- Pisani A., Sutter P., Hamaus N., Alizadeh E., Biswas R., Wandelt B. D., Hirata C. M., 2015, *Phys. Rev. D*, 92, 083531
- Pycke J. R., Russell E., 2016, *ApJ*, 821, 110
- Riebe K., et al., 2013, *Astron. Nachr.*, 334, 691
- Sánchez C., et al., 2017, *MNRAS*, 465, 746
- Schneider P., 1996, *MNRAS*, 283, 837
- Scoville N., et al., 2007a, *ApJS*, 172, 1
- Scoville N., et al., 2007b, *ApJS*, 172, 150
- Sheldon E. S., Huff E. M., 2017, *ApJ*, 841, 24
- Silverman B. W., 1986, *Density estimation for statistics and data analysis*. Vol. 26, Cleveland, OH: CRC press
- Simet M., McClintock T., Mandelbaum R., Rozo E., Rykoff E., Sheldon E., Wechsler R. H., 2017, *MNRAS*, 466, 3103
- Springel V., et al., 2005, *Nature*, 435, 629
- Starck J.-L., Donoho D. L., Candès E. J., 2003, *A&A*, 398, 785
- Starck J.-L., Aghanim N., Forni O., 2004, *A&A*, 416, 9
- Starck J.-L., Murtagh F., Fadili J., 2015, *Sparse image and signal processing: Wavelets and related geometric multiscale analysis*. Cambridge: Cambridge University Press
- Vafaei Sadr A., Movahed S., Farhang M., Ringeval C., Bouchet F., 2017, *MNRAS*, 475, 1010
- Villani C., 2008, *Optimal transport: Old and new*. Vol. 338, Springer Science & Business Media
- Wang G., Ye J. C., Mueller K., Fessler J. A., 2018, *IEEE Trans. Med. Imaging*, 37, 1289
- White S. D. M., Frenk C. S., Davis M., Efstathiou G., 1987, *ApJ*, 313, 505
- Xavier H. S., Abdalla F. B., Joachimi B., 2016a, FLASK: Full-sky Lognormal Astro-fields Simulation Kit (ascl:1606.015)
- Xavier H. S., Abdalla F. B., Joachimi B., 2016b, *MNRAS*, 459, 3693
- Xia Q., et al., 2020, *A&A*, 633, A89
- Xu X., Cisewski-Kehe J., Green S. B., Nagai D., 2019, *Astron. Comput.*, 27, 34
- Yang L. F., Neyrinck M. C., Aragón-Calvo M. A., Falck B., Silk J., 2015, *MNRAS*, 451, 3606
- York D. G., et al., 2000, *AJ*, 120, 1579
- Zeldovich I. B., Einasto J., Shandarin S. F., 1982, *Nature*, 300, 407
- Zuntz J., et al., 2018, *MNRAS*, 481, 1149

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