

# Improving Galaxy Clustering Measurements with Deep Learning: analysis of the DECaLS DR7 data

Mehdi Rezaie<sup>1</sup><sup>★</sup>, Hee-Jong Seo<sup>1</sup><sup>†</sup>, Ashley J. Ross<sup>2</sup>, and Razvan C. Bunescu<sup>3</sup>

<sup>1</sup>*Department of Physics and Astronomy, Ohio University, Athens, OH 45701, USA*

<sup>2</sup>*The Center of Cosmology and Astro Particle Physics, the Ohio State University, Columbus, OH 43210, USA*

<sup>3</sup>*School of Electrical Engineering and Computer Science, Ohio University, Athens, OH 45701, USA*

Accepted XXX. Received YYY; in original form ZZZ

## ABSTRACT

Robust measurements of cosmological parameters from galaxy surveys rely on our understanding of systematic effects that impact the observed galaxy density field. In this paper we present, validate, and implement the idea of adopting the systematics mitigation method of Artificial Neural Networks for modeling the relationship between the target galaxy density field and various observational realities including but not limited to Galactic extinction, seeing, and stellar density. Our method by construction allows a wide class of models and alleviates over-training by performing k-fold cross-validation and dimensionality reduction via backward feature elimination. By permuting the choice of the training, validation, and test sets, we construct a selection mask for the entire footprint. We apply our method on the extended Baryon Oscillation Spectroscopic Survey (eBOSS) Emission Line Galaxies (ELGs) selection from the Dark Energy Camera Legacy Survey (DECaLS) Data Release 7 and show that the spurious large-scale contamination due to imaging systematics can be significantly reduced by up-weighting the observed galaxy density using the selection mask from the neural network and that our method is more effective than the conventional linear and quadratic polynomial functions. We perform extensive analyses on simulated mock datasets with and without systematic effects. Our analyses indicate that our methodology is more robust to overfitting compared to the conventional methods. This method can be utilized in the catalog generation of future spectroscopic galaxy surveys such as eBOSS and Dark Energy Spectroscopic Instrument (DESI) to better mitigate observational systematics.

**Key words:** editorials, notices — miscellaneous — catalogs — surveys

## 1 INTRODUCTION

Our current understanding of the Universe is founded upon statistical analyses of cosmological observables, such as the large-scale structure of galaxies, the cosmic microwave background, and the Hubble diagram of distant type Ia supernovae (e.g., Efstathiou et al. 1988; Fisher et al. 1993; Smoot et al. 1992; Mather et al. 1994; Riess et al. 1998; Perlmutter et al. 1999; Ata et al. 2017; Alam et al. 2017; Jones et al. 2018; Akrami et al. 2018; Elvin-Poole et al. 2018). Among these probes, galaxy surveys aim at constructing clustering statistics of galaxies with which we can investigate the dynamic of the cosmic expansion due to dark energy, test Einstein’s theory of gravity, and constrain the total mass of neutrinos and statistical properties of

the primordial fluctuations, etc (Peebles 1973; Kaiser 1987; Mukhanov et al. 1992; Hamilton 1998; Eisenstein et al. 1998; Seo & Eisenstein 2003; Eisenstein 2005; Sánchez et al. 2008; Dalal et al. 2008).

The field of cosmology has been substantially advanced by the recent torrents of spectroscopic and imaging datasets from galaxy surveys such as the Sloan Digital Sky Survey (SDSS), Two Degree Field Galaxy Redshift Survey, and WiggleZ Dark Energy Survey (York et al. 2000; Colless et al. 2001; Drinkwater et al. 2010). The SDSS has been gathering data through different phases SDSS-I (2000-2005), SDSS-II (2005-2008), SDSS-III (2008-2014), and SDSS-IV (2014-2019). In order to derive more robust constraints with higher statistical confidence along with advancements in the technology of spectrographs, software, and computing machines, future large galaxy surveys aim at not only a wider area but also fainter galaxies out to higher redshifts.

<sup>★</sup> E-mail: mr095415@ohio.edu

<sup>†</sup> E-mail: seoh@ohio.edu

As an upcoming ground-based survey, the Dark Energy Spectroscopic Instrument (DESI) experiment will gather spectra of thirty million galaxies over 14,000 deg<sup>2</sup> starting in late 2019; this is approximately a factor of ten increase in the number of galaxy spectra compared to those observed in SDSS I–IV. This massive amount of spectroscopic data will lead to groundbreaking measurements of cosmological parameters through statistical data analyses of the clustering measurements of the 3D distribution of galaxies and quasars (Aghamousa et al. 2016).

The Large Synoptic Survey Telescope (LSST) is another ground-based survey currently being constructed. It will gather 20 Terabytes of imaging data every night and will cover 18,000 deg<sup>2</sup> of the sky for ten years with a sample size of  $2 \times 10^9$  galaxies. The LSST would provide enough imaging data to address many puzzling astrophysical problems from the nature of dark matter, the growth of structure to our own Milky Way. Given such a data volume, many anticipate that the LSST will revolutionize the way astronomers do research and data analysis (Ivezic et al. 2008; LSST Science Collaborations et al. 2017).

The enormous increased data volume provided by DESI and the LSST will significantly improve statistical confidences but will require analyses that are more complex and sensitive to the unknown systematic effects. A particular area of concern is the systematic effects due to imaging attributes such as atmospheric conditions, foreground stellar density, and/or inaccurate calibrations of magnitudes. These systematic effects can affect the target galaxy selections and therefore induce the non-cosmological perturbations into the galaxy density field, leading to excess clustering amplitudes, especially on large scales (see e.g. Myers et al. 2007; Thomas et al. 2011b,a; Ross et al. 2011; Ross et al. 2012; Ho et al. 2012; Huterer et al. 2013; Pullen & Hirata 2013).

Robust and precise measurements of cosmological parameters from the large-scale galaxy clustering are contingent upon thorough treatment of such systematic effects. Many techniques have been developed to mitigate the effects. One can generally classify these methods into the mode projection, regression, and Monte Carlo simulation of fake objects.

The mode-projection based techniques attribute a large variance to the spatial modes that strongly correlate with the potential systematic maps such as imaging attributes, thereby effectively removing those modes from the estimation of power spectrum (see e.g. Rybicki & Press 1992; Tegmark 1997; Tegmark et al. 1998; Slosar et al. 2004; Ho et al. 2008; Pullen & Hirata 2013; Leistedt et al. 2013; Leistedt & Peiris 2014). In detail, the basic mode projection (Leistedt et al. 2013) produces an unbiased power spectrum and is equivalent to a marginalization over a free amplitude for the contamination produced by a given map. The caveat is that the variance of the estimated clustering increases by projecting out more modes for more imaging attributes. The *extended* mode projection technique (Leistedt & Peiris 2014) resolves this issue by selecting a subset of the imaging maps using a  $\chi^2$  threshold

to determine the significance of a potential map. The limitation of the mode-projection based methods is that they are only applicable to the two-dimensional clustering measurements, and they reintroduce a small bias (Elsner et al. 2015). Kalus et al. (2016) extended the idea to the 3D clustering statistics and developed a new step to unbiased the measurements (for an application on SDSS-III BOSS data see e.g., Kalus et al. 2018)

The regression-based techniques model the dependency of the galaxy density on the potential systematic fluctuations and estimate the parameters of the proposed function by solving a least-squares problem, or by cross-correlating the galaxy density map with the potential systematic maps (see e.g. Ross et al. 2011; Ross et al. 2012, 2017; Ho et al. 2012; Delubac et al. 2016; Prakash et al. 2016; Raichoor et al. 2017; Laurent et al. 2017; Elvin-Poole et al. 2018; Bautista et al. 2018). The best fit model produces a *selection mask (function)* or a set of *photometric weights* that quantifies the systematic effects in the galaxy density fluctuation induced by the imaging pipeline, survey depth, and other observational attributes. The selection mask is then used to up-weight the observed galaxy density map to mitigate the systematic effects. The regression-based methods often assume a linear model (with linear or quadratic polynomial terms), and use all of the data to estimate the parameters of the given regression model; however, the assumption that the systematic effects are linear might not necessarily hold for strong contamination, e.g., close to the Galactic plane. Ho et al. (2012) analyzed photometric Luminous Red Galaxies in SDSS-III Data Release 8 and showed that the excess clustering due to the stellar contamination on large scales (e.g., roughly greater than twenty degrees) cannot be removed with a linear approximation. Ross et al. (2013) investigated the local non-Gaussianity ( $f_{NL}$ ) using the BOSS Data Release 9 “CMASS” sample of galaxies (Ahn et al. 2012) and found that a robust cosmological measurement on very large scales is essentially limited by the systematic effects. Their analysis indicated that a more effective systematics correction is preferred relative to the selection mask based on the linear modeling of the stellar density contamination. Recently, Elvin-Poole et al. (2018) developed a methodology based on  $\chi^2$  statistics to rank the imaging maps based on their significance, and derived the selection mask by regressing against the significant maps.

Another promising, yet computationally expensive, approach injects artificial sources into real imaging in order to forward-model the galaxy survey selection mask introduced by real imaging systematics (see e.g. Bergé et al. 2013; Suchyta et al. 2016). Rapid developments of multi-core processors and efficient compilers will pave the path for the application of these methods on big galaxy surveys.

In this paper, we develop a systematics mitigation method based on artificial neural networks. Our methodology models the galaxy density dependence on observational imaging attributes to construct the selection mask, without making any prior assumption of the linearity of the fitting model. Most importantly, this methodology is less prone to over-training and the resulting removal of the clustering

signal by performing  $k$ -fold cross-validation (i.e., splitting the data into  $k$  number of groups/partitions from which one constructs the training, validation, and the test sets) and dimensionality reduction through backward feature elimination (i.e., removing redundant and irrelevant imaging attributes)(see e.g., Devijver & Kittler 1982; John et al. 1994; Koller & Sahami 1996; Kohavi & John 1997; Ramaswamy et al. 2001; Guyon & Elisseeff 2003). By permutation of the training, validation, and test sets, the selection mask for the entire footprint is constructed. We apply our method on galaxies in the Legacy Surveys Data Release 7 (DR7) (Dey et al. 2018) that are chosen with the eBOSS-ELG color-magnitude selection criteria (Raichoor et al. 2017) and compare its performance with that of the conventional, linear and quadratic polynomial regression methods. While the effect of mitigation on the data will be estimated qualitatively as well as quantitatively based on cross-correlating the observed galaxy density field and imaging maps, the data does not allow an absolute comparison to the unknown underlying cosmology. We therefore simulate two sets of 100 mock datasets, without and with the systematic effects that mimic those of DR7, apply various mitigation techniques in the same way we treat the real data, and test the resulting clustering signals against the ground truth.

This paper is organized as follows. Section 2 presents the imaging dataset from the Legacy Surveys DR7 used for our analysis. In Section 3, we describe our method of Artificial Feed Forward Neural Network in detail as well as the conventional multivariate regression approaches. In this section, we also explain the angular clustering statistics employed to assess the level of systematic effects and the mitigation efficiency. We further describe the procedure of producing the survey mocks with and without simulated contaminations. In Section 4, we present the results of mitigation for both DR7 and the mocks. Finally, we conclude with a summary of our findings and a discussion of the benefits of our methodology for future galaxy surveys in Section 5.

## 2 LEGACY SURVEYS DR7

We use the seventh release of data from the Legacy Surveys (Dey et al. 2018). The Legacy Surveys are a group of imaging surveys in three optical (r, g, z) and four Wide-field Infrared Survey Explorer (W1, W2, W3, W4; Wright et al. (2010)) passbands that will provide an inference model catalog amassing 14,000 deg<sup>2</sup> of the sky in order to pre-select the targets for the DESI survey (Lang et al. 2016; Aghamousa et al. 2016). Identification and mitigation of the systematic effects in the selection of galaxy samples from this imaging dataset are of vital importance to DESI, as spurious fluctuations in the target density will likely present as fluctuations in the transverse modes of the 3D field and/or changes in the shape of the redshift distribution. Both effects will need to be modeled in order to isolate the cosmological clustering of DESI galaxies. The ground-based surveys that probe the sky in the optical bands are the Beijing-Arizona Sky Survey (BASS) (Zou

**Table 1.** The Northern Galactic Cap color-magnitude selection of the eBOSS Emission Line Galaxies (Raichoor et al. 2017). We enforce the same selection for the entire sky. Note that our selection is slightly different from Raichoor et al. (2017) in the clean photometry criteria as explained in the main text.

Criterion	eBOSS ELG
Clean Photometry	0 mag < V < 11.5 mag Tycho2 stars mask BRICK_PRIMARY==True brightstarinblob==False
[OII] emitters	21.825 < g < 22.9
Redshift range	-0.068(r-z) + 0.457 < g-r < 0.112 (r-z) + 0.773 0.637(g-r) + 0.399 < r-z < -0.555 (g-r) + 1.901

et al. 2017), DECam Legacy Survey (DECaLS) and Mayall z-band Legacy Survey (MzLS)(see e.g., Dey et al. 2018). Additionally, the Legacy Surveys program takes advantage of another imaging survey, the Dark Energy Survey, for about 1,130 deg<sup>2</sup> of their southern sky footprint (Dark Energy Survey Collaboration: Fermilab & Flaugher 2005). DR7 is data only from DECaLS, and we refer to this data interchangeably as DECaLS DR7 or DR7 hereafter.

We construct the ELG catalog by adopting the Northern Galactic Cap eBOSS ELG color-magnitude selection criteria from Raichoor et al. (2017) on the DR7 sweep files (Dey et al. 2018) with a few differences in the clean photometry criteria (see Table 1). In detail, the original eBOSS ELG selection is based on DR3 while ours is based on DR7. Since the data structure changed from DR3 to DR7, we use `brightstarinblob` instead of `tycho2inblob` to eliminate objects that are near bright stars. In contrast to the original selection criteria, we do not apply the `decam_anymask[grz]=0` cut, as any effect from this cut will be encapsulated by the imaging attributes used in this analysis. Also, we drop the SDSS bright star mask from the criteria, as this mask is essentially replaced by the `brightstarinblob` mask. After constructing the galaxy catalog, we pixelize the galaxies into a HEALPix map (Gorski et al. 2005) with the resolution of 13.7 arcmin ( $N_{\text{side}} = 256$ ) in *ring* ordering format to create the observed galaxy density map.

We consider a total of 18 imaging attributes as potential sources of the systematic error since each of these attributes can affect the completeness and purity with which galaxies can be detected in the imaging data. We produce the HEALPix maps (Gorski et al. 2005) with  $N_{\text{side}} = 256$  and oversampling of four<sup>1</sup> for these attributes based on the DR7 ccds-annotated file using the `validationtests` pipeline<sup>2</sup> and the code that uses the methods described in

<sup>1</sup> In this context, ‘oversampling’ means dividing a pixel into sub-pixels in order to derive the given pixelized quantity more accurately. For example, oversampling of four means subdividing each pixel into 4<sup>2</sup> sub-pixels. If the target resolution is  $N_{\text{side}} = 256$ , the attributes will be derived based on a map with the resolution of 4×256 when oversampling is four.

<sup>2</sup> <https://github.com/legacysurvey/legacypipe/tree/master/validationtests>

Leistedt et al. (2016). These include three maps of Galactic structure: Galactic extinction (Schlegel et al. 1998), stellar density from Gaia DR2 (Brown et al. 2018), and Galactic neutral atomic hydrogen (HI) column density (Bekhti et al. 2016). We further pixelize quantities associated with the Legacy Surveys observations, including the total depth, mean seeing, mean sky brightness, minimum modified Julian date, and total exposure time in three passbands (r, g, and z). For clarity, we list each attribute below:

- **Galactic extinction** ( $EBV$ ), measured in magnitudes, is the infrared radiation of the dust particles in the Milky Way. We use the SFD map (Schlegel et al. 1998) as the estimator of the E(B-V) reddening. The reddening is the process in which the dust particles in the Galactic plane absorb and scatter the optical light in the infrared. This reddening effect affects the measured brightness of the objects, i.e., the detectability of the targets. We correct the magnitudes of the objects for the Milky Way extinction prior to the galaxy (*target*) selection using the extinction coefficients of 2.165, 3.214, and 1.211 respectively for r, g, and z bands based on Schlafly & Finkbeiner (2011).
- **Galaxy depth** ( $depth$ ) defines the brightness of the faintest detectable galaxy at  $5 - \sigma$  confidence, measured in AB magnitudes. The measured depth in the catalogs does not include the effect of Galactic extinction (described above), so we apply the extinction corrections to the depth maps in the same manner.
- **Stellar density** ( $nstar$ ), measured in  $\text{deg}^{-2}$ , is constructed by pixelization of the Gaia DR2 star catalog (Brown et al. 2018) with the g-magnitude cut of  $12 < gmag < 17$ . The stellar foreground affects the galaxy density in two ways. First, the colors of stars overlap with those of galaxies, and consequently stars can be mis-identified as galaxies and included in the sample, which will result in a positive correlation between the stellar and galaxy distribution. Second, the foreground light from stars impacts the ability to detect the galaxies that are behind them, e.g., by directly obscuring their light or by altering the sky background, which will cause a negative correlation between the two distributions. The second effect may reduce the completeness with which galaxies are selected and was the dominant systematic effect on the BOSS galaxies (Ross et al. 2012). The Gaia-based stellar map is a biased set of the underlying stars that actually impact the data. Assuming that there exists a non-linear mapping between the Gaia stellar map and the truth stellar population, linear models might be insufficient to fully describe the stellar contamination. This motivates the application of non-linear models.
- **Hydrogen atom column density** ( $HI$ ), measured in  $\text{cm}^{-2}$ , is another useful tracer of the Galactic structure, which increases at regions closer to the Milky Way plane. The hydrogen column density

map is based on data from the Effelsberg-Bonn HI Survey (EBHIS) and the third revision of the Galactic All-Sky Survey (GASS). EBHIS and GASS have identical angular resolution and sensitivity, and provide a full-sky map of the neutral hydrogen column density (Bekhti et al. 2016). This map provides complementary information to the Galactic extinction and stellar density maps. Hereafter,  $\ln HI$  refers to the natural logarithm of the HI column density.

**Sky brightness** ( $skymag$ ) relates to the background level that is estimated and subtracted from the images as part of the photometric processing. It thus alters the depth of the imaging. It is measured in AB  $\text{mag}/\text{arcsec}^2$ .

- **Seeing** ( $seeing$ ) is the full width at half maximum of the point spread function (PSF), i.e., the sharpness of a telescope image, measured in arcseconds. It quantifies the turbulence in the atmosphere at the time of the observation and is sensitive to the optical system of the telescope, e.g., whether or not it is out of focus. Bad seeing conditions can make stars that are point sources appear as extended objects, therefore falsely being selected as galaxies. The seeing in the catalogs is measured in CCD ‘pixel’. We use a multiplicative factor of 0.262 to transform the seeing unit to arcseconds.
- **Modified Julian Date** ( $MJD$ ) is the traditional dating method used by astronomers, measured in days. If a portion of data taken during a specific period is affected by observational conditions during that period, regressing against MJD could mitigate that effect.
- **Exposure time** ( $exptime$ ) is the length of time, measured in seconds, during which the CCD was exposed to the object light. Longer exposures are needed to observe fainter objects. The Legacy Surveys data is built up from many overlapping images, and we map the total exposure time, per band, in any given area. A longer exposure time thus corresponds to a greater depth, all else being equal.

As part of the process of producing the maps, we determine the fractional CCD coverage per passband,  $fracgood$  ( $f_{pix}$ ), within each pixel with oversampling of four. We define the minimum of  $f_{pix}$  in r, g, and z passbands as the *completeness* weight of each pixel,

$$\text{completeness } f_{pix} = \min(f_{pix,r}, f_{pix,g}, f_{pix,z}). \quad (1)$$

We apply the following arbitrary cuts, somewhat motivated by the eBOSS target selection, on the depth and  $f_{pix}$  values to eliminate the regions with shallow depth and low pixel completeness due to insufficient available information:

$$\begin{aligned}
 \text{depth}_r &\geq 22.0, \\
 \text{depth}_g &\geq 21.4, \\
 \text{depth}_z &\geq 20.5, \\
 \text{and } f_{\text{pix}} &\geq 0.2,
 \end{aligned} \tag{2}$$

which results in 187,257 pixels and an effective total area of 9,459 deg<sup>2</sup> after taking  $f_{\text{pix}}$  into account. We report the mean, 15.9-, and 84.1-th percentiles of the imaging attributes on the masked footprint in Tab. 2.

As an exploratory analysis, we use the Pearson correlation coefficient (PCC) to assess the linear correlation between the data attributes. For two variables  $X$  and  $Y$ , PCC is defined as,

$$\rho_{X,Y} = \frac{\text{cov}(X,Y)}{\sqrt{\text{cov}(X,X)\text{cov}(Y,Y)}}, \tag{3}$$

where  $\text{cov}(X,Y)$  is the covariance between  $X$  and  $Y$  across all pixels. In Fig. 1, we show the observed galaxy density after the pixel completeness (i.e., fraggood  $f_{\text{pix}}$ ) correction in the top panel and the correlation (PCC) matrix between the DR7 attributes as well as the galaxy density ( $ngal$ , the bottom row) in the bottom panel. These statistics indicate that Galactic foregrounds, such as stellar density  $nstar$ , neutral hydrogen column density  $lnHI$ , and Galactic extinction  $EBV$ , are moderately anti-correlated with the observed galaxy density. Each of these maps traces the structure of the Milky Way and the anti-correlation with  $ngal$  implies that, for example, closer to the Galactic plane where the extinction and stellar density are high, there is a systematic decline in the density of galaxies we selected in our sample. The top-left corner of Fig. 1 shows that these three imaging attributes are strongly correlated with each other. Likewise, the negative correlation of  $ngal$  with  $seeing$  indicates that as  $seeing$  increases the detection of ELGs becomes more challenging. On the other hand, we find a positive correlation between  $ngal$  and  $depths$ , which can be explained by the fact that as the depth decreases, e.g., we cannot observe fainter objects, the number of galaxies decreases as well.

This matrix overall demonstrates that the correlation among the imaging variables is not negligible. For instance, in addition to the aforementioned correlation among the Galactic attributes, there is an anti-correlation between the MJD and depth values. Likewise, there is an anti-correlation between the seeing and depth values. The complex correlation between the imaging attributes causes degeneracies, and therefore, complicates the modeling of systematic effects, which cannot be ignored and needs careful treatment.

### 3 METHODOLOGY

#### 3.1 Observed galaxy density

In our methodology, we treat the mitigation of imaging systematics as a regression problem, in which we aim to model the relationship between the observed galaxy density (*label*) and the imaging attributes (*features*) that are the potential sources of the systematic error. Note that we do

**Table 2.** The statistics of the DR7 imaging attributes used in this paper. Due to the non-Gaussian nature of the attributes, we report the mean, 15.9-, and 84.1-th percentile points of the imaging attributes.

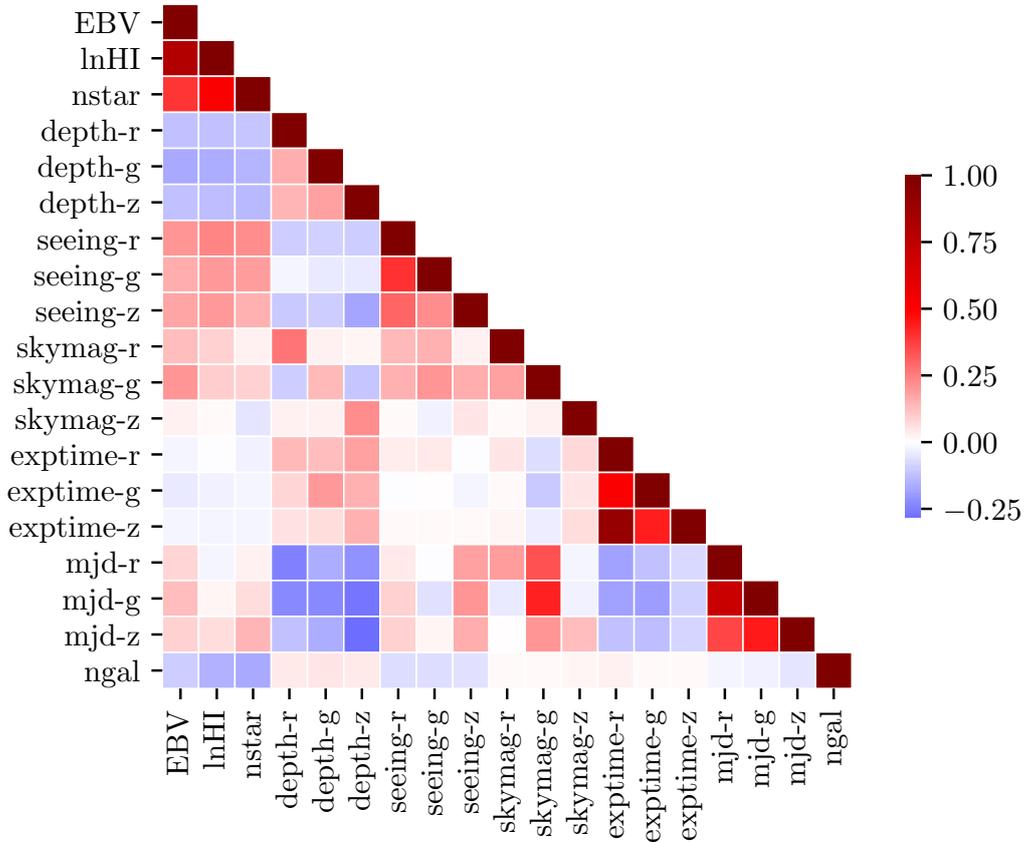
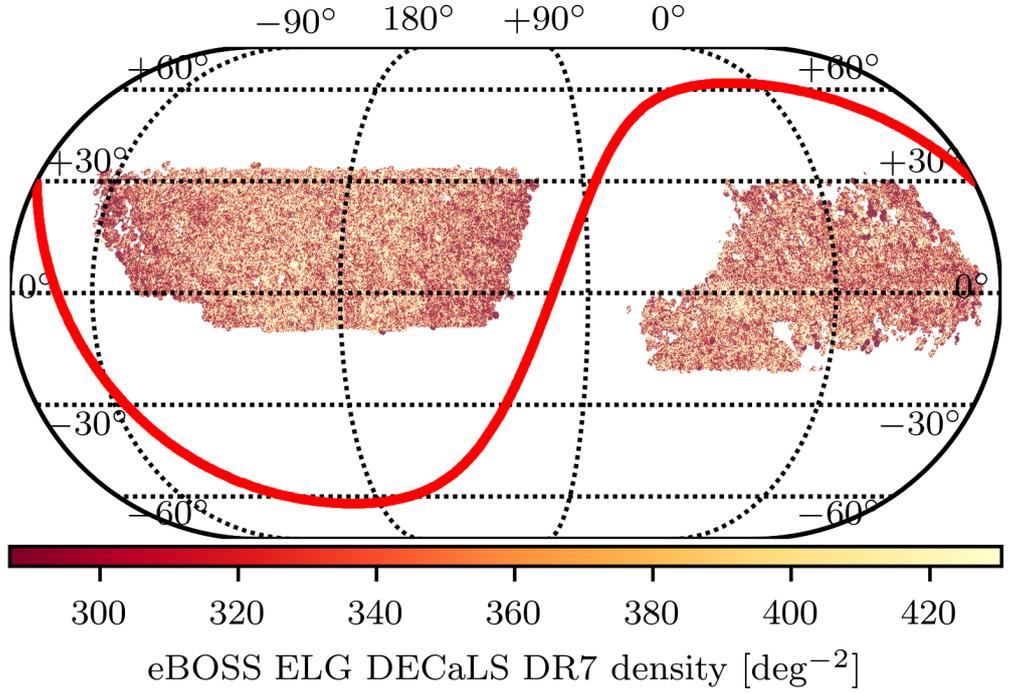
Imaging map	15.9%	mean	84.1%
EBV [mag]	0.023	0.048	0.075
ln(HI/cm <sup>2</sup> )	46.67	47.21	47.71
depth-r [mag]	23.46	23.96	24.33
depth-g [mag]	23.90	24.34	24.55
depth-z [mag]	22.57	22.93	23.23
seeing-r [arcsec]	1.19	1.41	1.61
seeing-g [arcsec]	1.32	1.56	1.78
seeing-z [arcsec]	1.12	1.31	1.51
skymag-r [mag/arcsec <sup>2</sup> ]	23.57	23.96	24.39
skymag-g [mag/arcsec <sup>2</sup> ]	25.06	25.39	25.80
skymag-z [mag/arcsec <sup>2</sup> ]	21.72	22.04	22.38
exptime-r [sec]	138.8	480.7	551.2
exptime-g [sec]	213.3	680.6	642.2
exptime-z [sec]	261.4	651.6	658.1
mjd-r [day]	56599.3	57232.7	57953.3
mjd-g [day]	56856.3	57358.1	57956.3
mjd-z [day]	56402.4	57005.0	57447.3

not include the positional information as *input* features since we do not want the mitigation to fit the cosmological clustering pattern. The solution of the regression then would provide the predicted mean galaxy density (i.e., in the absence of clustering or shot-noise) solely based on the imaging attributes of that location. We use this predicted galaxy density as the *survey* selection mask to be applied to any observed galaxy map in the attempt to eliminate the systematic effects and therefore isolate the cosmological fluctuation. Below we describe our procedure.

In this paper, we focus on the multiplicative systematic effects. The observed number of galaxies within pixel  $i$  can be expressed in terms of the true number of galaxies  $n_i$  and the contamination model  $\mathcal{F}(\mathbf{s}_i)$  as

$$n_i^o(\mathbf{s}_i) = n_i \mathcal{F}(\mathbf{s}_i), \tag{4}$$

where  $\mathbf{s}_i$  is a vector representing the imaging attributes  $\mathbf{s}$  of pixel  $i$ , and the contamination model  $\mathcal{F}(\mathbf{s}_i)$  is an unknown function representing the systematic effects which could be either a linear, non-linear, or a more complex combination of the imaging attributes. Multiplicative systematics are associated with obscuration and area-loss due to foreground stellar density, Galactic extinction, etc. On the other hand, additive systematics are associated mostly with stellar contamination, as described in Myers et al. (2007); Ross et al. (2011); Ho et al. (2012); Prakash et al. (2016); Crocchi et al. (2016). When averaged over many pixels, the effect of additive systematics can be absorbed into the constant term of the multiplicative model  $\mathcal{F}$ , assuming there is no correlation between the imaging maps and the true galaxy density field. The modeling of  $\mathcal{F}(\mathbf{s}_i)$  can be approached by a wide variety of techniques, ranging from the traditional methods based on multivariate functions to non-parametric and non-linear models based on machine learning or deep learning such as random tree forests and neural networks



**Figure 1.** *Top panel:* the pixelated density map of the eBOSS-like ELGs from DR7 after correcting for the completeness of each pixel (see eq., 1) and masking based on the survey depth and completeness cuts, see eq., 2. The solid red curve represents the Galactic plane. This figure is generated by the code described in <https://nbviewer.jupyter.org/github/desiutil/blob/master/doc/nb/SkyMapExamples.ipynb>. *Bottom panel:* the color-coded Pearson correlation matrix between each pair of the DR7 imaging attributes.

(Breiman 2001; Geurts et al. 2006).

The cosmological information is contained in the true overdensity that is given by

$$\delta_i = n_i / (f_{\text{pix},i} \bar{n}) - 1, \quad (5)$$

accounting for the pixel completeness where  $\bar{n}$  is the ‘true’ average number of galaxies. Then,

$$n_i^o = f_{\text{pix},i} \bar{n} (1 + \delta_i) \mathcal{F}(s_i). \quad (6)$$

This  $n_i^o / f_{\text{pix},i}$  is equivalent to the observed  $ngal$  aforementioned. Since we do not know the true average number density  $\bar{n}$  of the data, we estimate  $\bar{n}$  from the average of the observed galaxy field,

$$\hat{\bar{n}} = \frac{\sum_i n_i^o}{\sum_i f_{\text{pix},i}}, \quad (7)$$

and treat  $\hat{\bar{n}} \equiv \bar{n}$ . Due to the finite volume of our sample,  $\hat{\bar{n}} \neq \bar{n}$  even in the absence of systematic effects. This imposes the well-known integral constraint effect on any clustering analysis. We further ignore any systematic effect on  $\hat{\bar{n}}$  due to the fact we use  $n_i^o$ ; that is, Eq. 7 converges to  $\bar{n}$  only when  $\sum_i f_{\text{pix},i} = \sum_i f_{\text{pix},i} \mathcal{F}_i$ . In this sense we are modeling the relative systematic effect without necessarily determining the accurate ‘true’  $\bar{n}$ . We will use simulated results to test our methodology, and the analysis applied to the simulations with a limited footprint will be subject to the similar finite-volume and systematic effects on  $\hat{\bar{n}}$ , thus providing a fair comparison and means to catch any obvious problem with this approximation.

Finally, we define the normalized galaxy density per pixel  $t_i$ ,

$$t_i(s_i) \equiv \frac{n_i^o(s_i)}{f_{\text{pix},i} \hat{\bar{n}}} = (1 + \delta_i) \mathcal{F}(s_i). \quad (8)$$

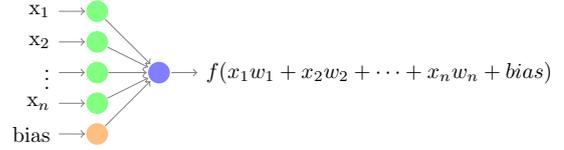
With this definition, we can estimate the unknown contamination model  $\hat{\mathcal{F}}$  (or selection mask) by modeling the dependence of  $t_i$  on  $s_i$ . When averaged over many spatial positions, the cosmological fluctuations  $\delta_i$  will be averaged out and therefore the observed  $t_i$  averaged across many pixels with the same imaging attribute, should only be a function of  $s$  and return  $\mathcal{F}$ :

$$\langle t_i(s_i) \rangle_i \approx \langle \mathcal{F}(s_i) \rangle_i = \mathcal{F}(s). \quad (9)$$

The inverse of the selection mask which is equivalent to the photometric weights ( $w_i^{\text{sys}}$ ) in other studies can therefore be used to remove the systematic effects from the observed galaxy number map (cf. Eq. 4),

$$\hat{n}_i = \frac{n_i^o}{\hat{\mathcal{F}}} = n_i^o w_i^{\text{sys}}. \quad (10)$$

In the following we describe how we obtain  $\hat{\mathcal{F}}$  using different regression approaches, e.g., neural networks and multivariate linear functions. From now on, the terms *features* and *label* associated with each data point refer to  $s$  and  $t$  of each HEALPix pixel, respectively.



**Figure 2.** A schematic diagram of a single neuron with the activation function  $f$ . The neuron takes a set of inputs  $\mathbf{x}=(x_1, x_2, \dots, x_n)$ , multiplies each of them by its associated weight  $\mathbf{w}=(w_1, w_2, \dots, w_n)$ , and sums the weighted values and a threshold (or a constant offset) which is called *bias*, to form a pre-activation value,  $z = \sum_{i=1}^n x_i w_i + \text{bias}$ , which is a linear process. The neuron then transforms the pre-activation  $z$  to the output using the activation function  $f(z)$ , which is where the nonlinear process can enter.

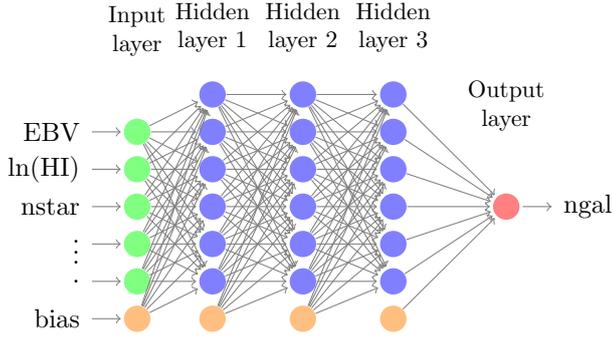
### 3.2 Mitigation with Neural Networks

We will apply *fully connected feed forward neural networks* in order to tackle our regression problem. Fig. 2 illustrates a schematic diagram of a neuron, the building block of a neural network, which generates the output based on a linear combination of the inputs followed by a nonlinear transformation, the activation function  $f$ .

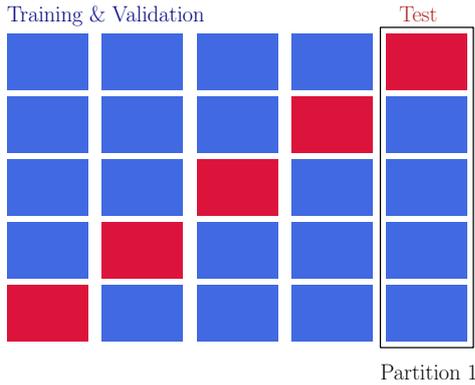
Fig. 3 illustrates the architecture of a fully connected feed forward neural network with the imaging attributes in the input layer, three hidden layers of six non-linear neurons in the middle, and a single neuron without any activation function in the output layer, as an example. The *bias* neuron in each layer is shown in orange and is analogous to the intercept in linear regression. The output of the neural network will be an estimation of the contamination model  $\hat{\mathcal{F}}$  (see Eq. 8). If we use the identity function as the activation function (e.g.,  $f(z) = z$ ), regardless of the number of hidden layers, the neural network is equivalent to a linear model. This means that our methodology is a generalization or an extension of the conventional linear mitigation methods. The modeling capabilities of neural networks depend on the number of hidden layers, type of non-linear activation function and the number of neurons in each hidden layer (see e.g. Cybenko 1989; Hornik et al. 1989; Funahashi 1989; Tamura & Tateishi 1997; Huang 2003; Lin et al. 2017; Rolnick & Tegmark 2017).

We use the rectifier  $f(z) = \max(0, z)$  as the activation function for the hidden layer neurons, which alleviates the ‘vanishing gradient’ problem (see e.g., Nair & Hinton 2010; Glorot et al. 2011; Krizhevsky et al. 2012; Dahl et al. 2013; Montufar et al. 2014).

We utilize  $k$ -fold cross-validation with  $k = 5$  folds/sub-groups to train the parameters, tune the hyper-parameters, and to estimate the predictive performance of the neural network. As illustrated in Fig. 4, we randomly split the entire pixel data (187,257 pixels) into five folds and construct the training, validation and test data sets out of these five folds; three folds are assigned to the training set, one fold is assigned to the validation set, and the remaining one fold is assigned to the test set. A specific assignment of the five folds to these three sets forms one ‘partition’. We construct five different partitions such that each fold is used once as test fold. This  $k$ -fold cross-validation scheme en-



**Figure 3.** A schematic illustration of a fully connected feed forward neural network with the imaging attributes in the input layer, three hidden layers of six neurons in the middle, and a single neuron on the output layer, as an example. The blue-colored neurons have non-linear activation functions, while the red-colored neuron lacks any activation function. In reality, we employ the validation procedure to choose the best number of hidden layers while keeping the total number of hidden layer neurons fixed at 40 (i.e., approximately twice the number of imaging attributes in this study).



**Figure 4.** A visualization of the five-fold permutations of data-split. The data is randomly split into five equal-size folds, and by permutation of the folds we construct five partitions of data. For each partition/permutation, three folds are assigned to the training set, one fold for the validation set, and the remaining one fold for the test set: therefore, 60% training, 20% validation, and 20% test sets. Each column represents a partition. The test folds are shown in red, while the training and validation folds are shown in blue. The key points are : 1) The folds are not contiguous (in RA, DEC) as may be implied by this cartoon. 2) There is no overlap between the training, validation, and test folds within a partition. 3) One can reconstruct the entire data by merging the test folds from the five partitions.

sures that a test example is never used for training or tuning.

We standardize (i.e., renormalize) the label and features of the training, validation, and test folds using the mean and standard deviation of the label and features of the corresponding training fold (i.e., similar to Tab. 2, but for the training set of each partition). We initialize the biases to zero and sample the weights of each layer randomly from a normal distribution whose variance is inversely proportional to the number of neurons of the previous layer (see e.g., He et al. 2015). Using the training fold, we utilize the adaptive gradient descent with momentum (*Adam*,

Kingma & Ba 2014) to update the parameters of the neural network with batches of size  $N_{\text{batch}}$ . Thus, the entire training set is split into  $N_{\text{batch}}$  batches and a gradient update is applied for each batch. One training epoch corresponds to processing the  $N_{\text{batch}}$  batches once (for more details on the training procedure, we refer the reader to see e.g., Ruder 2016). The hyper-parameters of *Adam*, specifically the moments decay rates and the tolerance, are fixed as follows:  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon = 10^{-8}$ . The default learning rate of 0.001 will be tuned using the validation data.

The network is trained to minimize the following cost function:

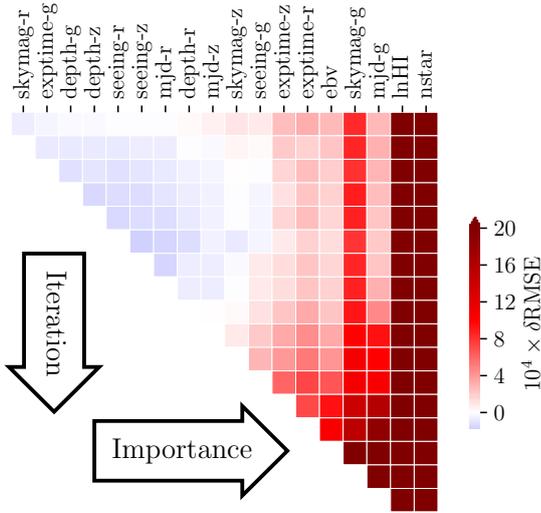
$$J = \frac{1}{N_{\text{batch}}} \sum_i^{N_{\text{batch}}} f_{\text{pix},i} [t_i - \hat{t}_i]^2 + \frac{\lambda}{2} \|w\|^2, \quad (11)$$

where the first term is the Mean Squared Error (MSE) weighted with  $f_{\text{pix},i}$ , and the second term is the L2 regularization term, used to penalize higher weight magnitudes and a larger number of neurons (Hoerl & Kennard 1970). The strength of the L2 penalty term is controlled via the regularization scale  $\lambda$ . The network is trained for a number of training epochs,  $N_{\text{epoch}}$ , although to avoid unnecessary training, we implement the early stopping technique with the tolerance of 1.e-4 and patience of 10, i.e., the training terminates if the validation MSE does not improve more than the tolerance within the last 10 epochs.

### 3.2.1 Backward feature elimination

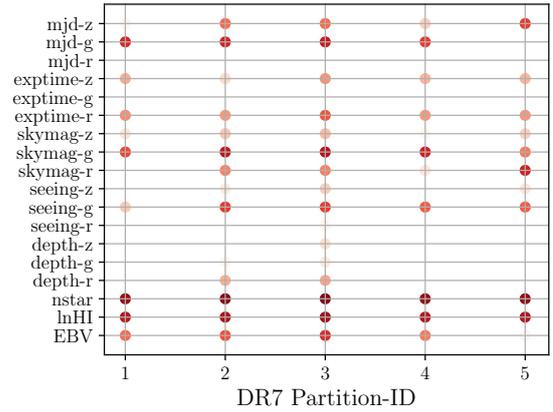
The input features are highly correlated as shown in Fig. 1, and therefore the 18 maps probably contain redundant information. We apply *backward feature elimination* (feature selection) to remove the redundant input features in order to reduce the noise in the prediction as well as to protect the cosmological information by avoiding too much freedom in modeling. We find that reducing the dimension of the input features, i.e., the imaging attributes, is an essential step to avoid over-fitting and regressing out the cosmological clustering.

We perform the feature selection for each partition separately. Initially, we train a linear model on all 18 input features with the following hyper-parameters: the initial learning rate of 0.001, batch-size of 1024, L2 regularization scale of zero, for 500 training epochs with early stopping. We record the validation MSE as a baseline criterion. Then we eliminate one input feature and train the linear model on the remaining 17 features. This trained linear model is applied to the validation set. The input feature whose removal has produced the highest decrease in validation MSE (i.e., the highest improvement in fitting) is permanently eliminated, leaving only 17 features for training. Note that if the feature contained useful information on the systematic effects, removing the feature would have made the fit worse. We repeat the regression using 17 input features, and so on, removing one feature at each iteration until either the validation MSE does not decrease relative to the baseline or all input features are removed. Fig. 5 shows the result of the backward feature elimination procedure



**Figure 5.** Feature importance of DR7 based on backward feature elimination for the first partition, as an example. This process iteratively removes the feature whose removal produces the largest decrease in the validation RMSE (i.e., the greatest improvement in fitting) until no decrease is observed. After the first iteration of removal (the top row), the removal of *skymag-r* decreased the validation RMSE the most and therefore *skymag-r* is removed. In the second iteration (the second row), removing *exptime-g* decreased the validation RMSE the most and therefore it is removed. However, in the ninth iteration, the removal of *mjd-z* did not decrease the validation RMSE and therefore the feature selection stops here, passing the rightmost ten features to the neural network regression. As a result of the process, the importance increases from left to right, and the rightmost ten maps in the figure (*ebv*, *nstar*, *logHI*, *seeing-g*, *skymag-g*, *skymag-z*, *exptime-r*, *exptime-z*, *mjd-g*, *mjd-z*) are the ones that worsen the validation RMSE when being removed from the input layer.

for the first partition of DR7, ranking the input features based on their importance from left to right. This result supports the trends seen in the exploratory analysis in § 2 which indicated strong linear correlations with the stellar density, Galactic extinction, and hydrogen column density in the data (see the correlation matrix in the bottom panel of Fig. 1). The color gradient indicates the relative change in validation Root Mean Squared Error (RMSE) when that particular feature is removed with reference to the baseline. We note that the order of removal is not the same as, for example, the color gradient order of the attributes in the first iteration. We believe it is because, as we remove the redundant features, the relevant importance of the remaining input features changes due to the complex correlations between the removed features and the remaining ones. The remaining features for each of the five partitions are shown in Fig. 6. The attributes *lnHI* and *nstar* are commonly identified as the most important features and then *ebv*, *seeing-g*, *skymag-g*, *skymag-z*, *exptime-r*, *exptime-z*, *mjd-g*, *mjd-z* are commonly identified for all 5 partitions.



**Figure 6.** Important imaging maps identified by the backward feature elimination (*feature selection*) procedure for the five partitions used for DR7. A darker color of a point within each partition represents a more important attribute identified by the feature selection procedure. Note that *nstar* is identified as the most important attribute in all partitions, i.e., across the footprint.

### 3.2.2 Hyper-parameter tuning, training, and testing

We train the hyper-parameters for each partition separately. At each training epoch and for each choice of hyper-parameters, we apply the trained network on the validation fold. We adjust the hyper-parameters accordingly such that the validation MSE is minimized. Our neural network has the following five hyper-parameters: number of hidden layers; number of training epochs  $N_{\text{epoch}}$ ; L2 regularization scale  $\lambda$ ; batch size  $N_{\text{batch}}$ ; Adam’s learning rate. We tune one hyper-parameter at a time. To find the optimal learning rate, we monitor the behavior of the cost function during training. We observe that a learning rate of 0.001 leads to a smoothly decreasing cost function vs. training epochs. We train the network for up to  $N_{\text{epoch}} = 500$  epochs although we implement the early stopping technique with the inertia (or patience) of 10 and the tolerance of  $1.e-4$ : this means the training will be stopped if none of the last 10 epochs achieved a smaller relative error reduction with respect to the minimum validation error, within the tolerance. For the number of hidden layers, we try the following architectures, in which the total number of hidden neurons is fixed at 40 (i.e. roughly twice the number of the features) except for the linear model:

- [0] : no hidden layers
- [40] : one hidden layer of 40 neurons
- [20, 20] : two hidden layers of 20 neurons on each
- [20, 10, 10] : three hidden layers of 20, 10 and 10 neurons
- [10, 10, 10, 10] : four hidden layers of 10 neurons

After finding the best number of layers, we proceed to tune  $\lambda$  by trying powers of 10, e.g., 0.001, 0.01, ..., 1, ..., 1000. Finally, we adjust  $N_{\text{batch}}$  by trying powers of 2, e.g., 128, 256, ..., 4096. The optimal set of the hyper-parameters for each partition is summarized in Tab. 3.

Once the grid search procedure identifies the best performing hyper-parameters out of the predefined ranges introduced in Section 3.2.2, the network is trained with

**Table 3.** The best hyper-parameters for each partition of DR7.

	number of layers	$\lambda$	$N_{\text{batch}}$
Partition 1	[20, 20]	0.001	4096
Partition 2	[20, 10, 10]	0.001	512
Partition 3	[20, 10, 10]	0.001	1024
Partition 4	[20, 10, 10]	0.001	512
Partition 5	[20, 10, 10]	0.001	2048

these hyper-parameters for 10 independent runs, each one with a different initialization of the weights and biases, and then applied on the test set. We compute the median of the predicted test label from the 10 runs and aggregate the results over the 5 different partitions to construct the map of the predicted label ( $\hat{\mathcal{F}}$ ) for the entire footprint. For our default method, the backward feature elimination is conducted for each partition and reduces the number of input features before the hyper-parameter training step, as illustrated in Fig. 7. This process is performed for each partition separately, each partition using a different fold as the test set, until the entire footprint is covered through the 5 test folds. The flow of the feature selection, hyper-parameter tuning, and testing is summarized in Fig. 7.

### 3.3 Mitigation with Multivariate Linear Functions

We use linear and quadratic polynomial functions to model the normalized galaxy density dependence on the imaging attributes (Eq. 9), as the benchmark approaches to be compared with the neural network. Unlike the proposed neural network method, no regularization or dimensionality reduction is performed, and all data are used to train the parameters of the regression models. Despite a lack of any deliberate machinery against over-fitting, we note that over-fitting is less likely to be an issue for this method since the size of the data is much greater than the number of the fitting parameters. Nevertheless, we tried splitting the sample into 60% of the training data to derive the best fit linear coefficients and 20% of the test set (i.e., the same training and test sample size to have a fair comparison with the neural network) to apply the derived coefficients and permuted five times until the test set covers the entire footprint. We find that such a split does not change the results of the linear regression. On the other hand, the limited flexibility of its parameterized form could be a weakness of this method and we believe it is responsible for the differences that the neural network method makes in the comparison presented in Section 4.1.

The contamination model from Eqs. 8 and 9 can be estimated via a multivariate linear function in terms of the standardized imaging attributes ( $\mathbf{s}$ ) as,

$$\hat{\mathcal{F}}(\mathbf{s}|b_0, \alpha) = b_0 + \sum_{m=1}^M \sum_{k=1}^{18} \alpha_{mk} \left( \frac{s_k - \bar{s}_k}{\sigma_k} \right)^m, \quad (12)$$

where  $M$  is the maximum power law index equal to 1 and 2 for the linear and quadratic polynomial model, respectively; the constants  $\bar{s}_k$  and  $\sigma_k$  are respectively the mean and standard deviation of the  $k$ 'th imaging quantity,  $s_k$  (cf.

Tab. 2). The parameters  $b_0$  and  $\alpha_{mk}$  are the intercept and the corresponding coefficients for each term, respectively, which are tuned by minimizing the weighted sum of the squared errors. The output of the regression is applied as the selection mask on the observed galaxy density to eliminate the systematic effects (see Eq., 10).

### 3.4 Angular Clustering Statistics

#### 3.4.1 One point statistics

In the absence of the systematic effects, the galaxy density field should be statistically independent from the imaging attributes and will only depend on the cosmological fluctuations while an individual dataset/mock will be subject to chance correlations within the statistical error. When averaged over many spatial positions, the cosmological fluctuations will be averaged out and therefore the average density should be equal to the mean density once the survey footprint is accounted for. A deviation from the mean density as a function of the imaging attributes, therefore, is indicative of the average dependence of the observed galaxy density on the corresponding imaging attributes. To assess the level of the contamination in the data, we compute the histogram of the spatially averaged galaxy density vs the imaging attributes. For each of the system attribute,  $s_k$ , we prepare 20 bins. For each bin of  $s_k$ , we have:

$$\frac{\bar{n}(s_k)}{\bar{n}_{\text{tot}}} \equiv \frac{1}{\bar{n}} \frac{\sum_{s_k \leq s_{k,i} < s_k + \Delta s_k} n_i^o}{\sum_{s_k \leq s_{k,i} < s_k + \Delta s_k} f_{\text{pix},i}}, \quad (13)$$

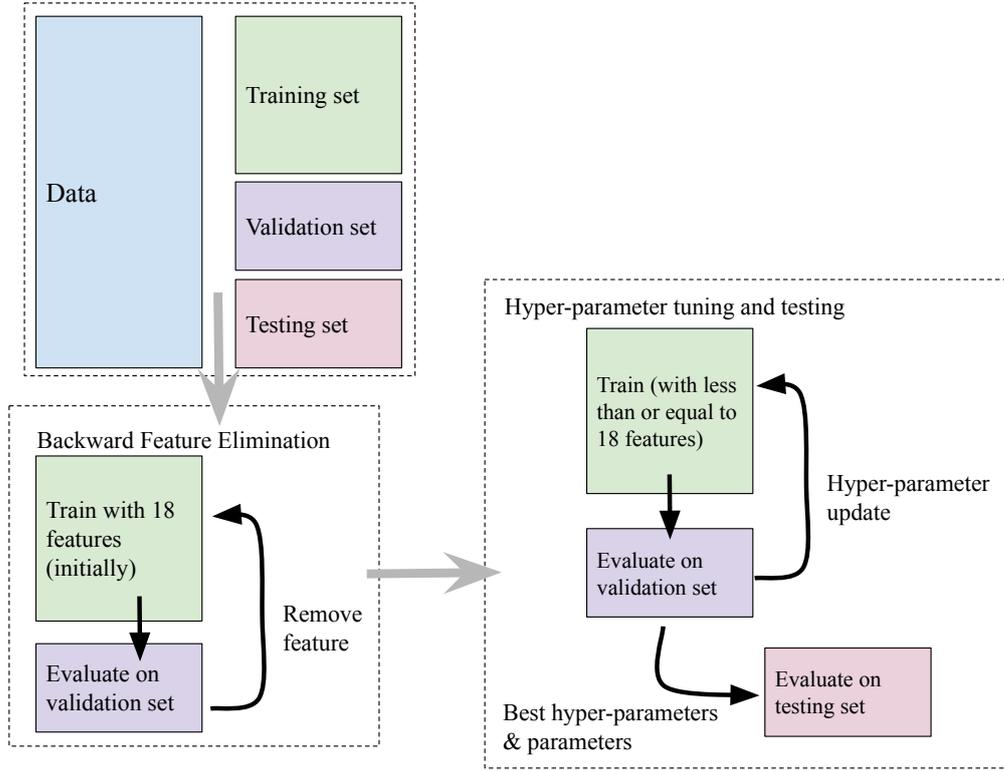
where the indices  $i$  and  $k$  respectively represent the pixel index, and the systematic index;  $\Delta s_k$  is the bin width arranged for different  $s_k$  bins such that each bin contains almost the same amount of effective area, in an attempt to suppress fluctuations in  $\bar{n}(s_k)/\bar{n}_{\text{tot}}$  due to small number statistics. We estimate the error bars using the Jackknife resampling of 20 non-contiguous subsamples of pixels within each imaging attribute bin (see e.g. Ross et al. 2011):

$$\sigma_{\text{Jack}}^2(s_k) = \frac{19}{20} \sum_{j=1}^{20} \left[ \frac{\bar{n}(s_k)}{\bar{n}_{\text{tot}}} - \frac{\bar{n}_j(s_k)}{\bar{n}_{\text{tot}}} \right]^2, \quad (14)$$

where  $\bar{n}_j(s_k)/\bar{n}_{\text{tot}}$  is computed over the entire sample when the  $j$ 'th Jackknife region is excluded. As a result of the adjusted  $\Delta s_k$ , the level of  $\sigma_{\text{Jack}}^2(s_k)$  is almost the same for all  $s_k$ . After mitigation of the systematic effects, one expects that the corrected density field is independent of the imaging attributes, i.e.,  $\bar{n}(s_k)/\bar{n}_{\text{tot}}$  being consistent with unity.

#### 3.4.2 Two-point clustering statistics

The two-point clustering statistic measures the spatial correlation of the galaxy density and has been the main statistic for extracting the cosmological information from galaxy surveys. We use the angular auto and cross two-point clustering statistics of the galaxy density field as well as of the imaging



**Figure 7.** A flow-chart of the backward feature elimination and hyper-parameter tuning for each partition. This entire process is performed five times for each of the five partitions/permutations such that the entire footprint is covered by aggregating the different testing folds.

attributes to estimate the impact of the potential systematics on the cosmological clustering signal and to examine the effectiveness of the mitigation techniques tested in this paper. For pixel  $i$ , we calculate the galaxy overdensity  $\delta_i$  using Eq. 5 and the fluctuation of a given imaging attribute  $\delta_i^s$  as

$$\hat{\delta}_i^s = \frac{S_i}{\hat{s}} - 1, \quad (15)$$

where  $\hat{s}$  is the mean of each imaging attribute weighted with  $f_{\text{pix},i}$ ,

$$\hat{s} = \frac{\sum_i f_{\text{pix},i} S_i}{\sum_i f_{\text{pix},i}} \quad (16)$$

following Ross et al. (2011).

By definition, Eqs. 7 and 15 ensure that the following integral of the observed quantity over the entire footprint vanishes:

$$\sum_i \hat{\delta}_i f_{\text{pix},i} = 0, \quad (17)$$

for both the galaxy as well as imaging attribute fluctuations. We utilize both the angular correlation function and angular power spectrum to extract the cosmological information from the galaxy density field. While our mitigation efficiency is evaluated based on the angular power spectrum, we also inspect the angular correlation function to make sure that the systematics are mitigated in that estimator as well, since both estimators are commonly used in the clustering anal-

ysis and they are complementary to each other given the limited range of data.

- **Angular Correlation Function:** we employ the HEALPix-based estimator to compute the angular correlation function which, for a separation angle  $\theta$ , is defined as (see e.g. Scranton et al. 2002; Ross et al. 2011)

$$\omega^{p,q}(\theta) = \frac{\sum_{ij} \hat{\delta}_i^p \hat{\delta}_j^q \Theta_{ij}(\theta) f_{\text{pix},i} f_{\text{pix},j}}{\sum_{ij} \Theta_{ij}(\theta) f_{\text{pix},i} f_{\text{pix},j}}, \quad (18)$$

where  $p = q$  gives an auto correlation function estimator,  $p \neq q$  gives a cross correlation function estimator, and  $\Theta_{ij}$  is one when two pixels  $i$  and  $j$  are separated from each other within  $\theta$  and  $\theta + \Delta\theta$ , or zero otherwise. Note that our estimator weighs each pixel overdensity with  $f_{\text{pix},i}$  since the pixels with a greater complete area coverage should have a higher signal to noise. Such weight is straightforwardly corrected by the denominator in Eq. 18 unlike its conjugate estimator (i.e., the power spectrum). Since our overdensity map resolution is limited by the pixel size, we set the  $\Delta\theta$  to be the resolution of a pixel ( $\sim 0.23$  deg).

- **Angular Power Spectrum :** one can conveniently expand a coordinate on the surface of a sphere in terms

of spherical harmonics or, if azimuthally symmetric, Legendre polynomials. We define the following estimator for expanding the galaxy overdensity:

$$\hat{\delta}_i = \sum_{\ell=0}^{\ell_{\max}} \sum_{m=-\ell}^{\ell} a_{\ell m} Y_{\ell m}(\theta_i, \phi_i), \quad (19)$$

where  $\theta, \phi$  represent the polar and azimuthal angular coordinates of pixel  $i$ , respectively. The cutoff at  $\ell = \ell_{\max}$  assumes that the signal power is not significant for modes  $\ell > \ell_{\max}$ . We define the following spherical harmonic (SH) transform estimator of overdensity ( $\hat{\delta}$ ) over the total number of non-empty pixels  $N_{\text{pix}}$ :

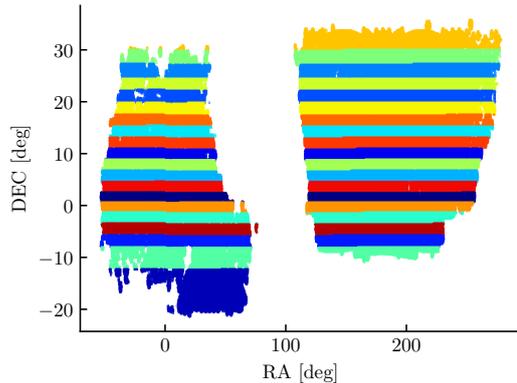
$$\hat{a}_{\ell m} = \frac{4\pi}{N_{\text{pix}}} \sum_{i=1}^{N_{\text{pix}}} \hat{\delta}_i f_{\text{pix},i} Y_{\ell m}^*(\theta_i, \phi_i), \quad (20)$$

where  $*$  represents the complex conjugate, and we again down-weight the overdensity in pixel  $i$  by the completeness ( $f_{\text{pix},i}$ ). Due to the survey window function implicit in the sum over the non-empty pixels and explicit in  $f_{\text{pix},i}$ , our estimator would not return unbiased estimates of the SH coefficients, unless the window function effect is corrected for, and also the expected orthogonalities between different SH modes would not hold. Nevertheless, we define the angular power spectrum estimator as the average of the magnitude of SH coefficients over  $m$ :

$$\hat{C}_{\ell}^{p,q} = \frac{1}{2\ell+1} \sum_{m=-\ell}^{\ell} \hat{a}_{\ell m}^p \hat{a}_{\ell m}^{q*}, \quad (21)$$

where  $p = q$  gives an auto power spectrum,  $p \neq q$  gives a cross power spectrum between the galaxy density and the imaging attributes. In order to compute the angular power spectrum,  $C_{\ell}$ , we make use of the ANAFAST function from HEALPix (Gorski et al. 2005) with the third order iteration of the quadrature to increase the accuracy<sup>3</sup>. Unlike in the angular correlation function, we do not attempt to correct for the survey window function/survey mask effect in the angular power spectrum both for DR7 and the mocks. We rather calculate the window effect on the theoretical models of power spectrum in Appendix A. For the mock test, we use the angular power spectrum observed in the mocks without the contamination model, i.e., the ‘Null’ case, as our baseline to compare with different mitigation methods.

We use the Jackknife resampling technique with 20 equal-area contiguous regions, as shown in Fig. 8, to estimate the error-bars on  $\omega(\theta)$  and  $C_{\ell}$  (see Eq. 14).<sup>4</sup> For both the mock



**Figure 8.** Twenty equal-area contiguous regions used to estimate the Jackknife errorbars for the 2D clustering statistics.

and real datasets, we also utilize the cross power spectra between the galaxy density and various imaging maps to evaluate the performance of the mitigation. In order to estimate the significance of the contamination in  $\hat{C}_{\ell}^{g,g}$  (or  $\omega^{g,g}(\theta)$ ) before and after mitigation, we calculate  $[\hat{C}_{\ell}^{g,s_k}]^2 / \hat{C}_{\ell}^{s_k,s_k}$  (or  $[\omega^{g,s_k}(\theta)]^2 / \omega^{s_k,s_k}(\theta)$ ) as a proxy.<sup>5</sup>

### 3.5 Survey Mocks

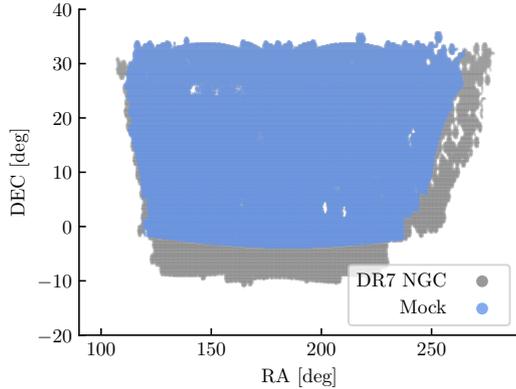
Imaging systematics tend to affect the clustering signal mainly on large scales (Myers et al. 2007; Huterer et al. 2013) and the distribution of galaxies on large scales at moderately low redshift can be well-approximated by a log-normal distribution (Coles & Jones 1991). We therefore believe that log-normal mocks would be sufficient for the purpose of validating our systematic mitigation techniques. We use the *Nbodykit* package (Hand et al. 2017) to generate one hundred log-normal cubic mocks with the box-side length of  $5274h^{-1}\text{Mpc}$  and  $1024^3$  mesh cells, with the input power spectrum matched to the linear power spectrum at  $z = 0.85$  based on the *Planck 2015* cosmology (Ade et al. 2016) (i.e., flat  $\Lambda\text{CDM}$  with  $\Omega_m = 0.3089 \pm 0.0062$ ,  $H_0 = 67.74 \pm 0.46$ ,  $\sigma_8 = 0.8159 \pm 0.0086$ ), with the galaxy bias of 1.5 and the volume density of  $1.947e-4h^3\text{Mpc}^{-3}$  (see e.g. Raichoor et al. 2017). Then, we use the *make\_survey* package (White et al. 2013) to sub-sample the mock galaxies based on the NGC eBOSS ELG redshift distribution in Raichoor et al. (2017) with the redshift cut of  $0.55 < z < 1.5$  and to transform the cubic mocks into survey-like mocks. We do not include redshift-space distortions (RSD) in the mocks as we believe that the systematics mitigation efficiency does not depend on the presence of RSD.

persion in the survey window function of the Jackknife samples. For a real survey the observational condition also varies across the footprint. Therefore, with Jackknife errors, the significance of any improvement on systematics treatment will be conservatively assessed.

<sup>5</sup> These quantities would be the true level of contamination to  $\hat{C}^{g,g}$  if the contamination model is linear and systematics are independent of one another (Ross et al. 2012; Ho et al. 2012).

<sup>3</sup> We refer the reader to <https://healpix.sourceforge.io/pdf/subroutines.pdf>, page 104.

<sup>4</sup> We use the mocks without imaging systematics (null mocks in Section 3.5) to compare the errors from the Jackknife subsamples of one mock with the errors among 100 full DECaLS-like mock footprints, and we find that the former is greater than the latter on  $\ell < \sim 10$  by a factor of 2-3, possibly due to the dis-



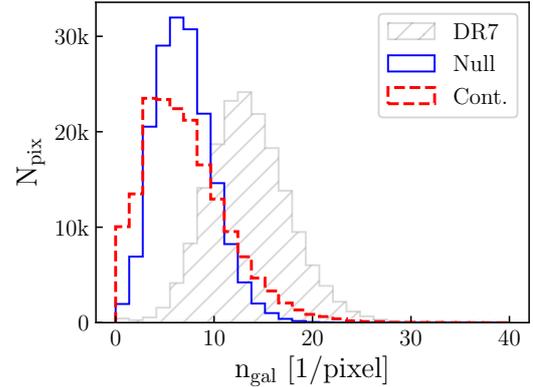
**Figure 9.** The projection of the mock footprint (blue) onto the North Galactic Cap of the DR7 footprint (gray). With this projection, the imaging attributes from the real data are assigned to the mocks.

The survey mocks are then projected onto the two-dimensional sky using HEALPix and overlaid on the NGC footprint of DR7 to be assigned with the DR7 imaging attributes. Fig. 9 illustrates the resulting projection of a simulated survey mock and DR7. Note that the mock footprint (89,672 pixels) is smaller than DR7 (187,257 pixels) almost by a factor of 2. We only use the pixels of the mock that have the DR7 imaging attributes available. The holes (e.g., RA and DEC around 200 and 5 deg) are the pixels that do not have the imaging attributes from the real data. In order to account for the mock survey footprint, we distribute 2,500 random points per  $\text{deg}^2$  within the mock footprint and derive the completeness map for the mocks.

### 3.5.1 Null Mocks

Our goal is to develop a systematics treatment methodology that maximally removes the systematic effects while minimally removes the true cosmological signal. The two aspects may not be simultaneously accomplished, in a way that depends on the signal, noise, and the correlation between the imaging maps and the true galaxy density. In our paper we choose to prioritize losing minimal cosmological information over maximally removing systematics. One way of ensuring this is to check if the mitigation method returns the true clustering in the presence of contaminations, which will be tested using the contaminated mocks. Another way is to check if the mitigation method correctly makes a null operation on the clustering in the absence of contaminations, returning  $\hat{\mathcal{F}} \approx 1$ . To this end, we utilize the 2D projected mocks without introducing any modulation due to imaging attributes in the galaxy density fields. Henceforth, we call this set of simulations, *null* mocks.

Fig. 10 shows the pixel distribution histograms of the number of galaxies per pixel of the mocks in comparison to that of DR7 on the common footprint. The average  $ngal = 7.0/\text{pixel}$  of the null mocks is smaller than  $ngal = 13.3/\text{pixel}$  of the DR7 data (even after accounting for the 5% loss due to tiling completeness and 27% loss due to the redshift range, as stated in Raichoor et al. (2017)). We believe that the difference in  $ngal$  is due to the different *clean photometry* criteria



**Figure 10.** Histogram of the number of galaxies per pixel for DR7 (hatched grey), null (solid blue), and contaminated (dashed red) mocks. All distributions are corrected for the pixel completeness  $f_{\text{pix}}$ . The DR7 distribution is scaled down to account for 5% tiling completeness, and 27% to  $0.55 < \text{reliable redshift} < 1.5$ . The residual difference between the mocks and DR7 shown here might be due to differences between the *clean photometry* criteria applied to eBOSS target selection and those we apply to DR7.

applied to the ELG selection in Raichoor et al. (2017) and to the targets in this paper. The standard deviation of  $ngal$  of the null mocks is 3.0, which is smaller than 4.6 of the DR7 ELGs.

### 3.5.2 Contaminated Mocks

We modulate the mock galaxy density fields using imaging attributes of DR7 and generate the contaminated mocks with additional random noise. The modulation is done based on the best fit coefficients of the imaging attributes and their covariances for  $\mathcal{F}$  (Eq. 8) that we derived from DR7 using our fiducial linear regression model.

In detail, we pick the 10 imaging attributes, i.e.,  $EBV$ ,  $nstar$ ,  $\ln HI$ ,  $seeing - g$ ,  $skymag - g$ ,  $skymag - z$ ,  $exptime - r$ ,  $exptime - z$ ,  $mjd - g$ ,  $mjd - z$  of DR7, which were selected from the feature selection procedure on one of the partitions, and modulate the mock density field  $n$  with  $\mathcal{F}$  that is derived from random deviates of the imaging attribute coefficients while accounting for their covariances. Since the measured covariance of such quantities includes both the cosmological fluctuation and the fluctuations due to the imaging attributes, we rescale the measured covariance matrix of the systematics such that the random fluctuation in  $ngal$  per pixel due to contamination is at a similar level to the cosmological fluctuation from the null mocks. As a result of the random fluctuation in the contamination model we introduced, some of the pixels will be assigned a negative galaxy number. We drop these pixels from our sample. This removes 3.1% of the mock footprint, reducing our mock footprint size from 89,672 to 86,875 pixels. We then introduce the Poisson process, i.e., another random variation step, to ensure the modulated galaxy number per pixel is an integer. These two random variation processes increase the noise in the mock datasets such that the standard deviation of  $ngal$  of the contaminated mocks (= 4.4) is almost the same as that of DR7, despite the

different average  $ngal$  (see Fig. 10).

Therefore our mock contamination is simpler than DR7 in that we adopted a linear model, which is chosen purposely since we do not want to give a priori advantage to our neural network method and also since all methods are capable of reproducing the linear model. Meanwhile, this setup is more challenging than the DR7 data since the mitigation is conducted in the presence of a greater level of noise. Note that, while we included only 10 dominant imaging attributes in the contamination, the remaining attributes in the DR7 data are correlated with these 10 attributes and therefore with the modulated galaxy density. All of the mitigation methods in the following mock test will be challenged to deal with such indirect correlations among the 18 attributes. Note that the effect of the footprint, i.e., the survey window effect, is the same for both the null mocks and contaminated mocks since we chose to apply the selection function on the galaxies while leaving the randoms intact. Therefore the null mocks serve as the baseline for estimating the level of systematics in the contaminated mocks.

## 4 RESULTS

In this section, we present the measurements of the clustering statistics before and after correcting for the systematic effects for the real dataset as well as the simulated ones. We demonstrate that the neural network is capable of learning more structure in the observed galaxy density field due to its greater flexibility beyond a fixed functional form, and therefore it can eliminate more excess clustering which is believed to be due to the imaging systematics. We then show the performance of the neural network and multivariate linear models when applied to the mock datasets.

### 4.1 Mitigating systematics from DR7

In the left panel of Fig. 11, we show the pixel distribution histograms of the selection masks from the three different regression models we consider in this paper. While all three models show fairly consistent selection masks for most of the pixels (note the logarithmic scaling of  $Npix$ ), the neural network method (solid red curve) returns extended tails due to a higher representation flexibility associated with its nonlinear nature. We remove pixels with  $\hat{\mathcal{F}} < 0.5$  or  $> 2.0$  from our data to avoid too aggressive selection correction since we believe none of these methods can be accurate enough for such a long baseline extrapolation. These pixels account for 1.0% of the original data (from 187,257 to 185,781 pixels). In the right panel, we show the spatial distribution of the removed pixels in the case of the neural network selection mask.

Fig. 12 illustrates the spatial distribution of the observed galaxy density before (top left) and after correction (top right) using the neural network selection mask. The bottom panels show the neural network selection mask used for the correction (left) in comparison to the masks derived from the linear (middle) and the quadratic polynomial (right) models. All three masks capture a very similar large scale pattern such as the decrease in the galaxy density

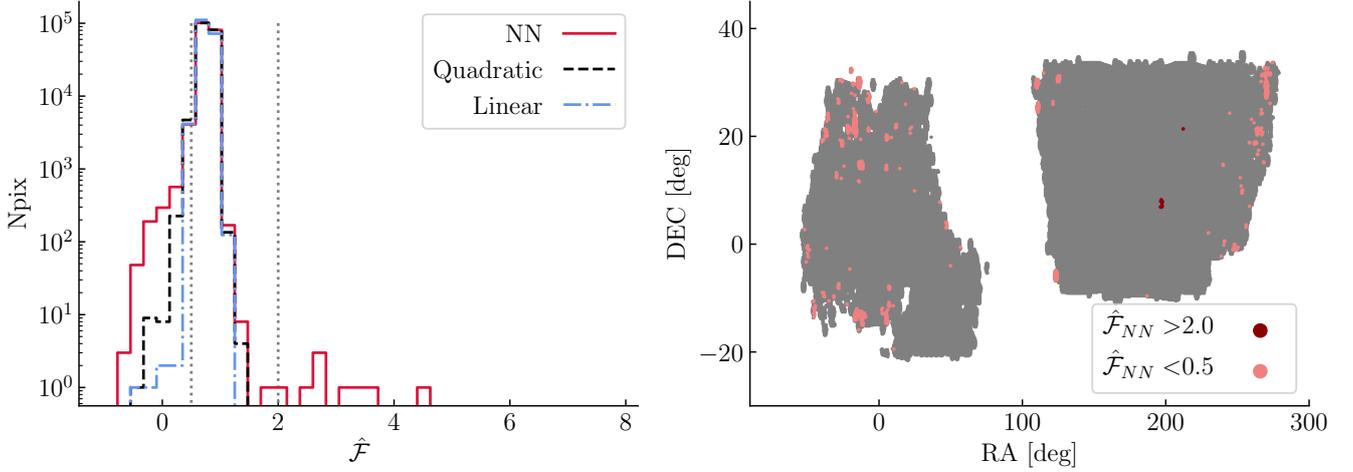
**Table 4.** The  $\chi^2$  values for the measured mean density of the DR7 galaxies vs imaging systematics, presented in Fig. 13. This table presents the cumulative values over all bins and all imaging attributes (i.e.,  $N_{bins} = 20 \text{ bins} \times 18 \text{ attributes}$ ) without accounting for the covariance both between the imaging maps and between different bins<sup>7</sup>.

Correction scheme	$\chi^2$	$N_{bins}$	$\chi^2/N_{bins}$
None	20921.633	360	58.116
Linear	2588.349	360	7.190
Quadratic	2623.006	360	7.286
Neural Network	966.601	360	2.685

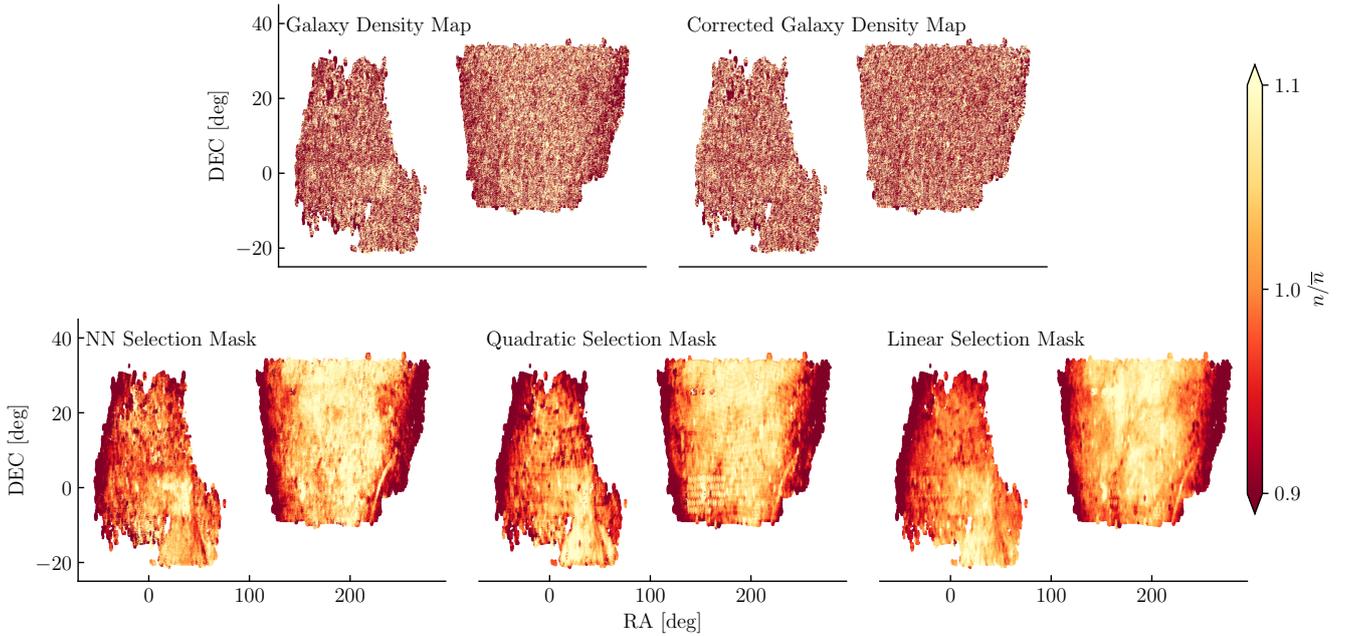
close to the Galactic midplane, which is consistent with the negative correlation coefficients between the galaxy density and the Galactic extinction, hydrogen column density, or stellar density. On smaller angular scales, the three selection masks show different fluctuation details. In the following analyses, we examine which method returns the least contaminated galaxy density distribution.

First, once the systematic effects are corrected for, the mean galaxy density should be independent of the imaging attributes. In Fig. 13 we show the mean number density of galaxies as a function of the imaging attributes. Again, different bins are set to include the same effective pixel area and therefore have the same sampling error. The solid black curve shows the galaxy mean density before correction and the solid red shows the result after correction with the selection mask of the neural network model. The dot-dashed curve represents the correction using the linear polynomial model and the dashed curve is for the quadratic polynomial model. The errorbars are computed using 20 Jackknife sub-samples and shown on only one case for clarity. Similar to what we found from the feature selection procedure, the stellar density, Galactic extinction, and HI density exhibit the strongest dependence before correction. After correction, all three methods return the fractional galaxy density close to unity. To quantify the deviation from unity, we report the  $\chi^2$  statistics in Tab. 4 while ignoring the covariance between the different bins and different imaging attributes. Overall, the neural network achieves the smallest deviation from unity which indicates its highest efficiency in reducing the systematic effects. Ideally we would like to have residual contamination less than the statistical error. Figure 13 and Table 4 implies that we need to further improve the mitigation techniques for future cosmological analyses. In Section 4.3 we provide a more detailed analysis using the same  $\chi^2$  statistics and the mocks to quantify the remaining systematics and assess whether or not the data is clean enough.

We next evaluate the performance of different mitigation techniques using the two-point statistics. We first show the cross power spectra between the DR7 observed galaxy density and various imaging attributes in the form of  $[\hat{C}_\ell^{g,sk}]^2 / \hat{C}_\ell^{sk,sk}$  in Fig. 14. Again, this quantity approximately represents the level of contamination from each attribute to the auto power spectrum of galaxy density and we therefore compare this with the uncertainty in the auto as well as cross power spectrum of galaxies (light and



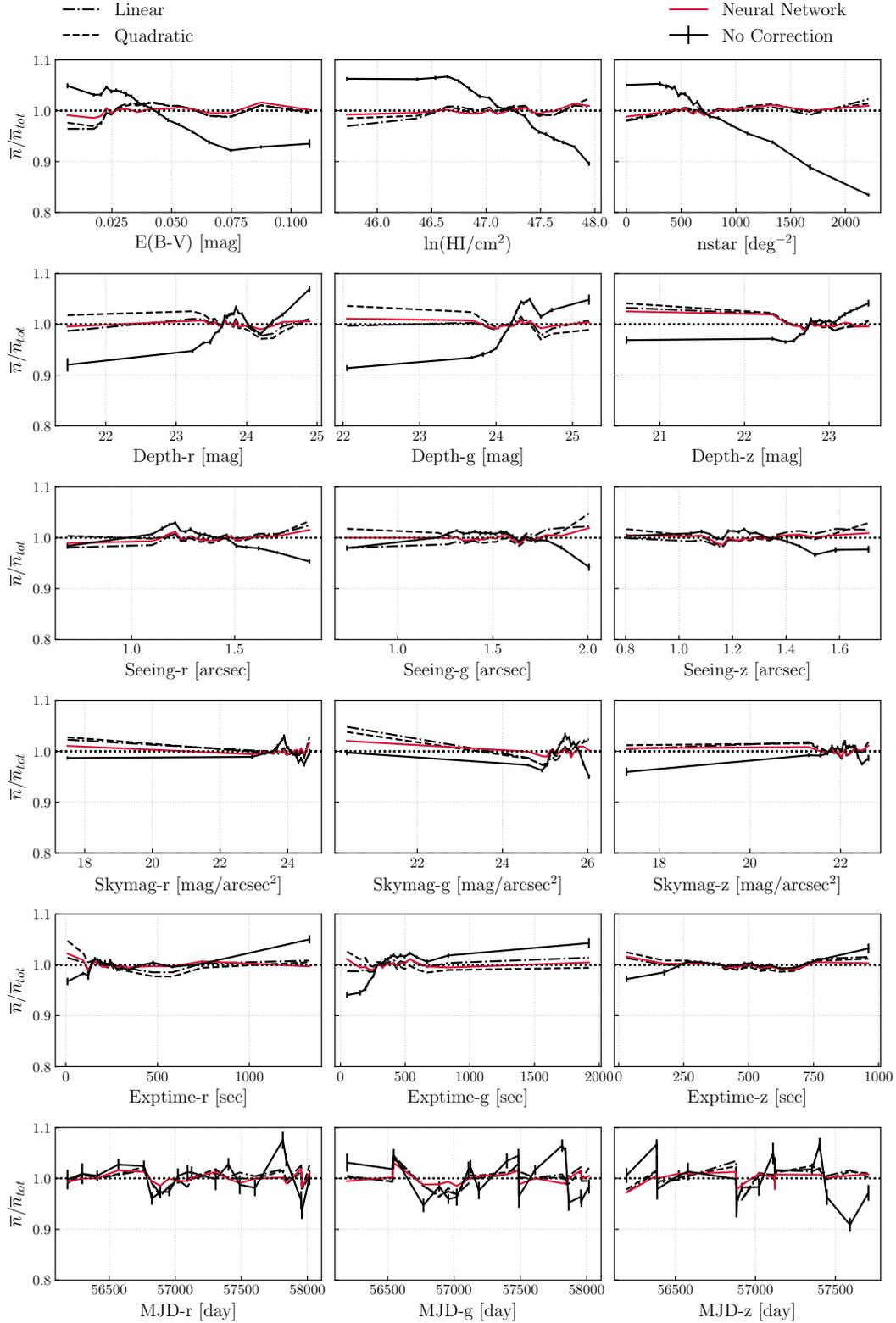
**Figure 11.** The pixel distribution of the selection masks for DR7. *Left:* Distribution of the selection masks (i.e., estimates of the contamination model) derived from different regression models. *Right:* Spatial scatter of the pixels we remove from our data due to the extreme values of the neural network selection mask.



**Figure 12.** *Top:* The normalized observed galaxy density of DR7 before and after systematics treatment using the neural network selection mask, respectively, from left to right. *Bottom:* the selection masks from the Neural Network, quadratic, and linear polynomial models, respectively from left to right. All three selection masks are able to capture the behavior that the galaxy density systematically drops at the footprint boundaries i.e., high extinction regions.

dark gray shades) which are estimated using the Jackknife resampling of 20 equal-area contiguous regions (see Fig. 8). Similarly, we plot  $[\omega^{g,sk}(\theta)]^2/\omega^{sk,sk}(\theta)$  in Fig. 15 to assess the contamination in the auto-correlation function. Fig. 14 and 15 show significant contamination on large scales from *ebv*, *lnHI*, and *nstar* compared to the statistical fluctuation estimated from the Jackknife subsampling of the data. The stellar density map shows the highest cross power spectrum with the galaxy density map, which is in agreement with

the previous results. Qualitatively, all three mitigation techniques perform well and substantially reduce the cross power below  $\ell \sim 30$  and over all separation scales in the cross-correlation function. The neural network method shows a slightly lower cross-power, but this appears to be merely related to the lower amplitude of the corresponding auto galaxy power spectrum compared to the other two cases. We note the spurious peak in the cross-correlation against exptime-z in Fig. 15 near the expected angular



**Figure 13.** The mean number density of the DR7 galaxies as a function of the potential systematics. The solid black curve shows the result before mitigation (*no correction*); the solid red curve is for the result after correcting with the neural network selection mask; the dot-dashed and dashed black curves represent mitigations with the linear and quadratic polynomial selection masks, respectively. The error bars are estimated using the Jackknife resampling of 20 non-contiguous subsamples of pixels within each imaging attribute bin (a total of 20 bins per attribute) and are shown only for one case. This plot again shows that the Galactic foregrounds such as the stellar density introduce a systematic trend in the galaxy density, which indicates a significant contamination by our own galaxy before mitigation. Such systematic trends are mostly removed with any of the three mitigation methods.

location of the BAO feature and such feature necessitates thorough investigations of imaging systematics in analyzing the auto clustering statistics of the spectroscopic data for BAO analysis.

We finally present the effect of the imaging attributes before and after mitigation on the auto galaxy clustering statistics. In Fig. 16, we illustrate the measured two-point clustering statistics for DR7; the measured angular power spectrum without shot-noise subtraction is shown in the left panel, and the HEALPix-based angular correlation function is shown in the right panel. The solid black curve shows the measured clustering before mitigation, while the corrected measurements using the traditional linear, quadratic polynomial, and the default neural network models are shown respectively with the black dot-dashed, black dashed, and solid red curves. In the right panel, the linear and the quadratic polynomial mitigation results are indistinguishable and overlaid.

The comparison between the clustering before correction (solid black curves) and after treatment (solid red, blue, dashed, and dot-dashed curves) suggests that the imaging systematics affect the clustering measurements mostly on large scales, e.g., large separation angles or small multipoles, as expected (see e.g., Myers et al. 2007; Ross et al. 2007; Huterer et al. 2013). We find that all of the mitigation methods are able to reduce such large scale contamination, while there still remains substantial excess clustering on large scales, mostly, with the two traditional linear multivariate methods. The neural network method is much more efficient in reducing such excess. When we investigate the effect of the survey window function on this data, we find that the window effect at  $\ell < 50$  is expected to be less than 5% (more details presented in Appendix A, see Fig. A1).

In comparison to our default neural network model, we also show the measurements mitigated with the neural network model without the feature selection process labeled as ‘plain’ (blue dashed curves), which is very similar to the default case. In the next section, we test the mitigation methods using the mock datasets for which we know the true clustering signals. As will be demonstrated, our default neural network model with the feature selection process is chosen based on this mock test.

## 4.2 Testing the mitigation methods on the mock data

We treat the mocks as if the contamination model was unknown and apply the mitigation pipeline on the mocks as exactly used for the real dataset. After modeling the selection mask for each mock, we remove the pixels whose selection masks values are  $< 0.5$  or  $> 2.0$ . This reduces the mock

footprint size from 86,875 to 86,867 pixels. Again, the mocks do not include the redshift-space distortions.

### 4.2.1 Feature selection of mock galaxies

Fig. 17 shows the distribution of the imaging attributes selected by the feature selection process for all of the five partitions of the 100 null (left) and contaminated (right) mocks. For the null mocks, there is no contamination and the feature selection correctly removes most or all of the imaging attributes, as demonstrated by the sparse distribution of the points in the left panel. The imaging attributes that survived feature selection, probably due to a coincidental correlation with the galaxy density, are randomly distributed. On the other hand, the right panel shows that the feature selection procedure correctly identifies most of the input contamination attributes (marked by ‘\*’ on the y-axis) for the contaminated mocks and almost always selects  $\ln HI$  and  $nstar$ . Indeed, as shown in Fig. 5, these two attributes were the two most significant input contamination.

### 4.2.2 Mean mock galaxy density

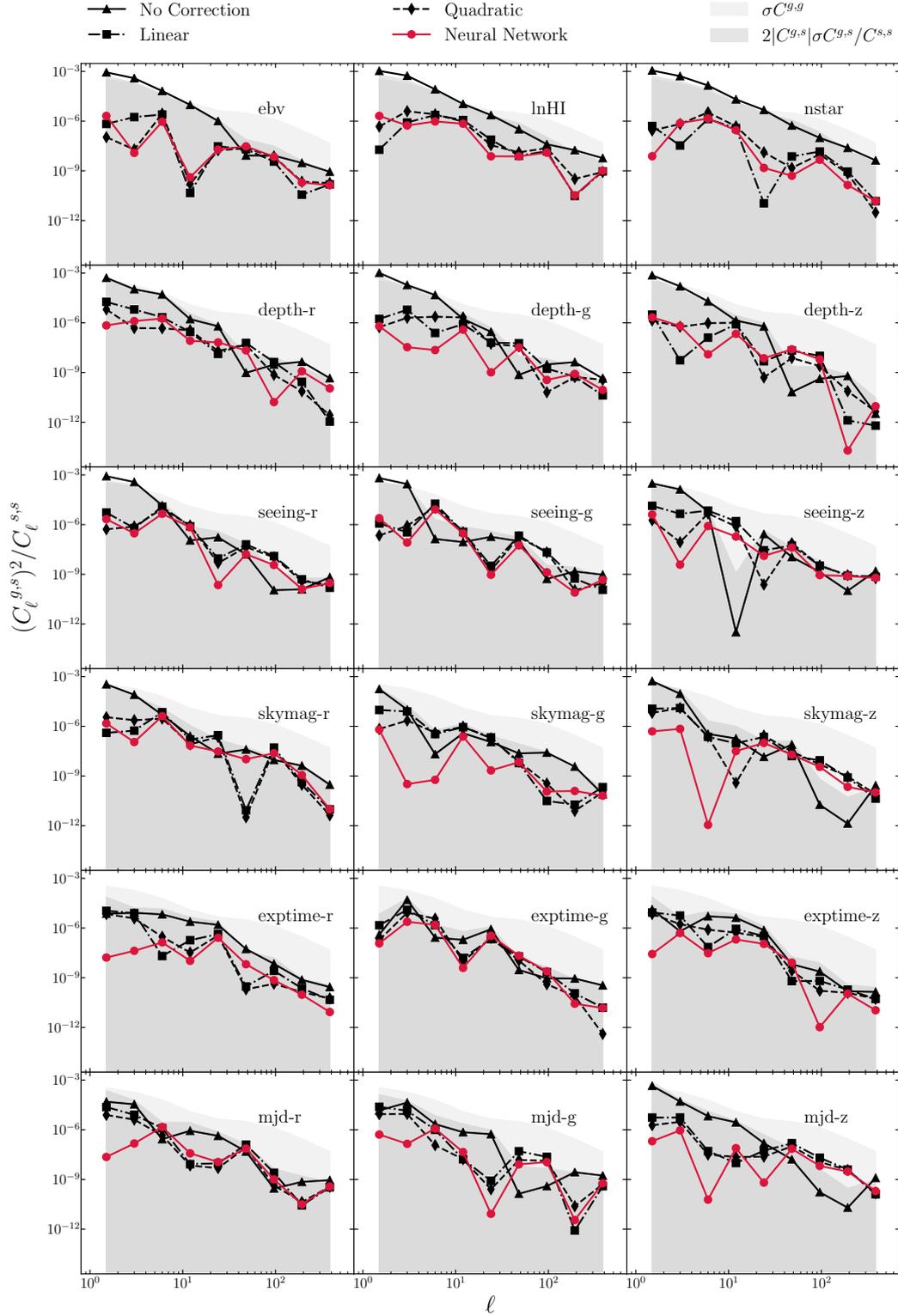
In Fig. 18 we show the number density of mock galaxies, averaged over the 100 mocks, as a function of the imaging attributes. As expected, the galaxy density of the contaminated mocks shows strong or moderate dependence on  $ebv$ ,  $nstar$ ,  $\ln HI$ ,  $seeing - g$ ,  $skymag - g$ ,  $skymag - z$ ,  $exptime - r$ ,  $exptime - z$ ,  $mjd - g$ , and  $mjd - z$  which were indeed the inputs to the contamination model. Meanwhile, the galaxy density also shows strong dependencies on  $mjd - r$ ,  $depth - r$ ,  $depth - g$ , and  $depth - z$  through the correlation between these and the input contamination attributes. Looking at this result alone from a real data perspective, one would not be able to single out the underlying imaging attributes that are directly responsible for the contamination. When the inputs to the mitigation procedure include all of the input contamination maps, Fig. 18 shows that all methods effectively remove the dependence. In subsection 4.3, we discuss further how well the underlying true mean density is recovered after mitigation.

### 4.2.3 Angular power spectrum of mock galaxies

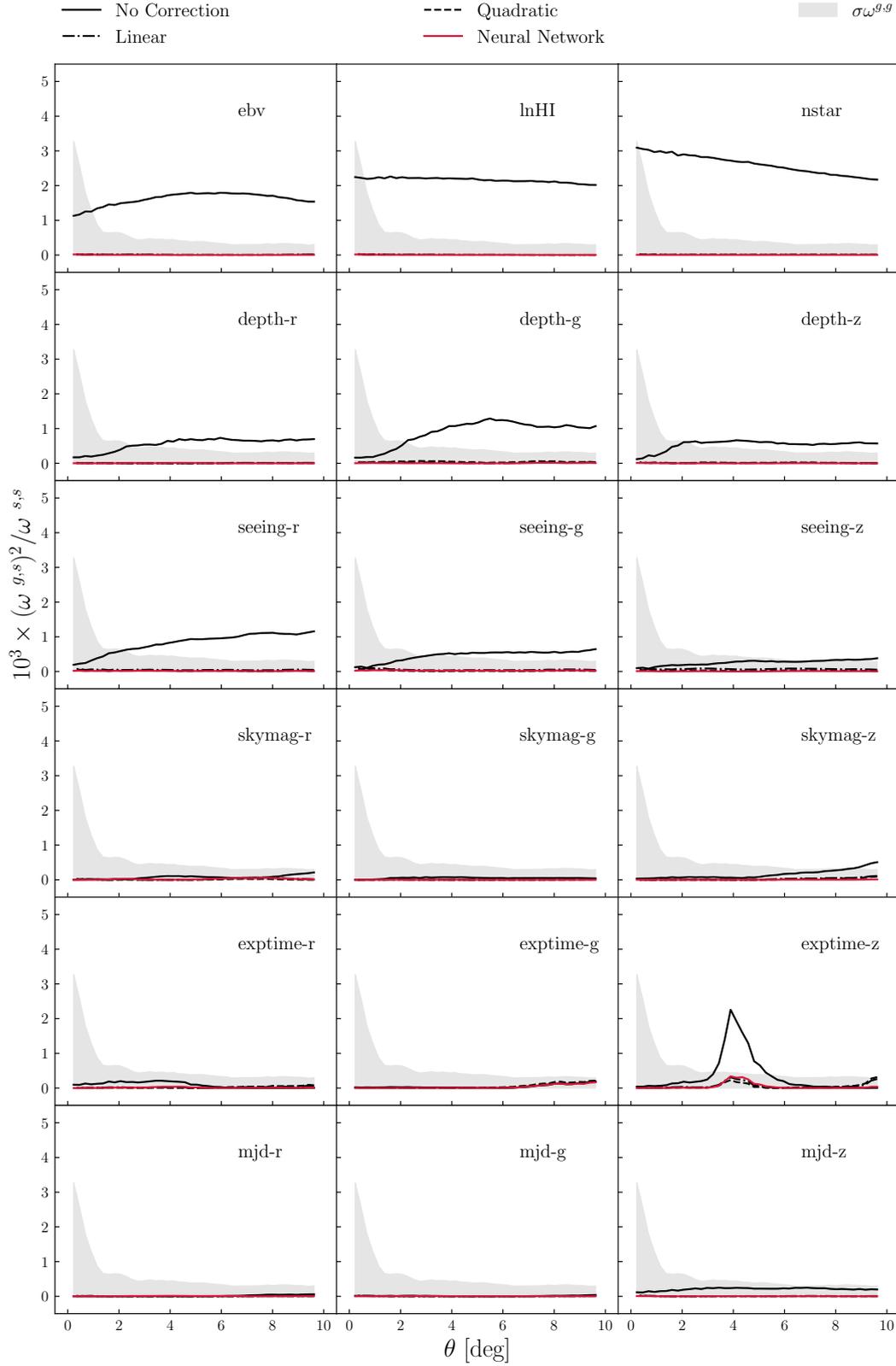
Fig. 19 shows the mean angular power spectrum of the 100 null and 100 contaminated mocks in the left and right panel, respectively, in the top row. In the middle row, we illustrate the remaining bias<sup>8</sup> on clustering after mitigation as an offset from the true power spectrum (i.e., the null power spectrum

<sup>7</sup> Note that the best fit neural network model was applied to the unseen data (i.e., the test set) unlike in the linear and quadratic polynomial models. Nevertheless the neural network method returns the smallest  $\chi^2$ , i.e., the highest efficiency.

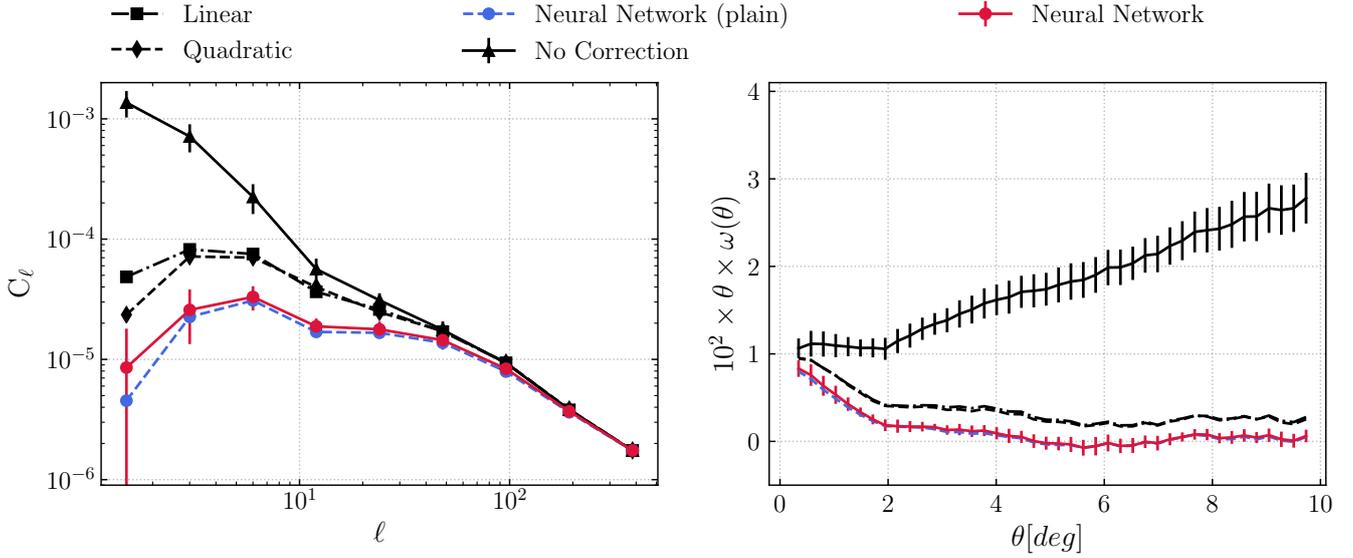
<sup>8</sup> Due to the two-step noise introduced during contamination, the shot noise of the contaminated power spectra is increased by almost a factor of two. We estimate the total offset noise to be around  $3.05 \times 10^{-6}$  from  $\sigma^2(ngal)/\bar{n}^2$  and subtract it from the power spectrum of the contaminated mocks. The additional shot noise is mostly originated from the Poisson process we applied, i.e., precisely  $1/\bar{n}$ , where  $\bar{n}$  is the average number density after contamination, while an extra  $\sim 10\%$  is also due to the noise we added to the contamination model.



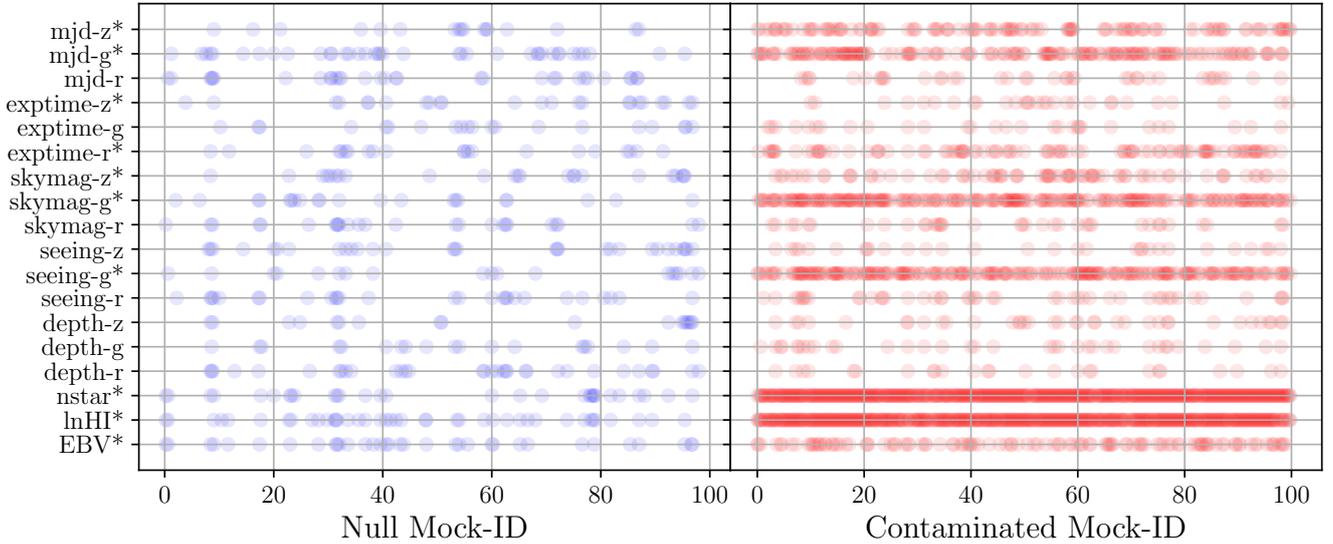
**Figure 14.** The cross power spectrum  $\hat{C}_\ell^{g,s_k}$  between the DR7 observed galaxy density and the imaging attributes  $s_k$  normalized by the auto power spectrum of the imaging attribute  $\hat{C}_\ell^{s_k,s_k}$ . The plotted quantity  $[\hat{C}_\ell^{g,s_k}]^2 / \hat{C}_\ell^{s_k,s_k}$  approximately represents the level of contamination to the auto power spectrum of the galaxy density  $\hat{C}_\ell^{g,g}$ . The light and dark grey shaded regions, respectively, show the Jackknife error estimate of  $\hat{C}_\ell^{g,g}$  and  $[\hat{C}_\ell^{g,s_k}]^2 / \hat{C}_\ell^{s_k,s_k}$  with the galaxy density  $g$  before mitigation. The black solid curve shows the result before mitigation (*no correction*), while the solid red curve shows the result after correcting for the systematics with the neural network selection mask. The dot-dashed and dashed black curves show the corrected results with the linear and quadratic polynomial model selection masks, respectively.



**Figure 15.** The cross correlation function  $\omega^{g,s^k}(\theta)$  between the DR7 observed galaxy density and the imaging attributes  $s_k$  normalized by the auto correlation function of the imaging attribute  $\omega^{s_k,s_k}(\theta)$ . The plotted quantity  $[\omega^{g,s^k}(\theta)]^2 / \omega^{s_k,s_k}(\theta)$  approximately represents the level of contamination to the auto correlation function of the galaxy density  $\omega^{g,g}(\theta)$ . The grey shaded region shows the Jackknife error estimate of  $\omega^{g,g}(\theta)$  before mitigation. All mitigation techniques are able to reduce the excess clustering signal which is due to the imaging systematics.



**Figure 16.** Two-point clustering statistics for DR7. *Left:* the measured angular power spectrum without shot-noise subtraction. *Right:* the HEALPix-based angular correlation function. Solid black curves show the measured statistics without correcting for the systematic effects (*no correction*). The dashed and dot-dashed black curves show the statistics after correcting with linear and quadratic polynomial mitigation methods, respectively. The solid red curves show the results after correcting with our default neural network method. The dashed blue curves show the results mitigated with the neural network method but without the feature selection process. The errors are estimated using the Jackknife resampling with 20 contiguous sub-regions and are shown only for a few cases for clarity (see Fig. 8).

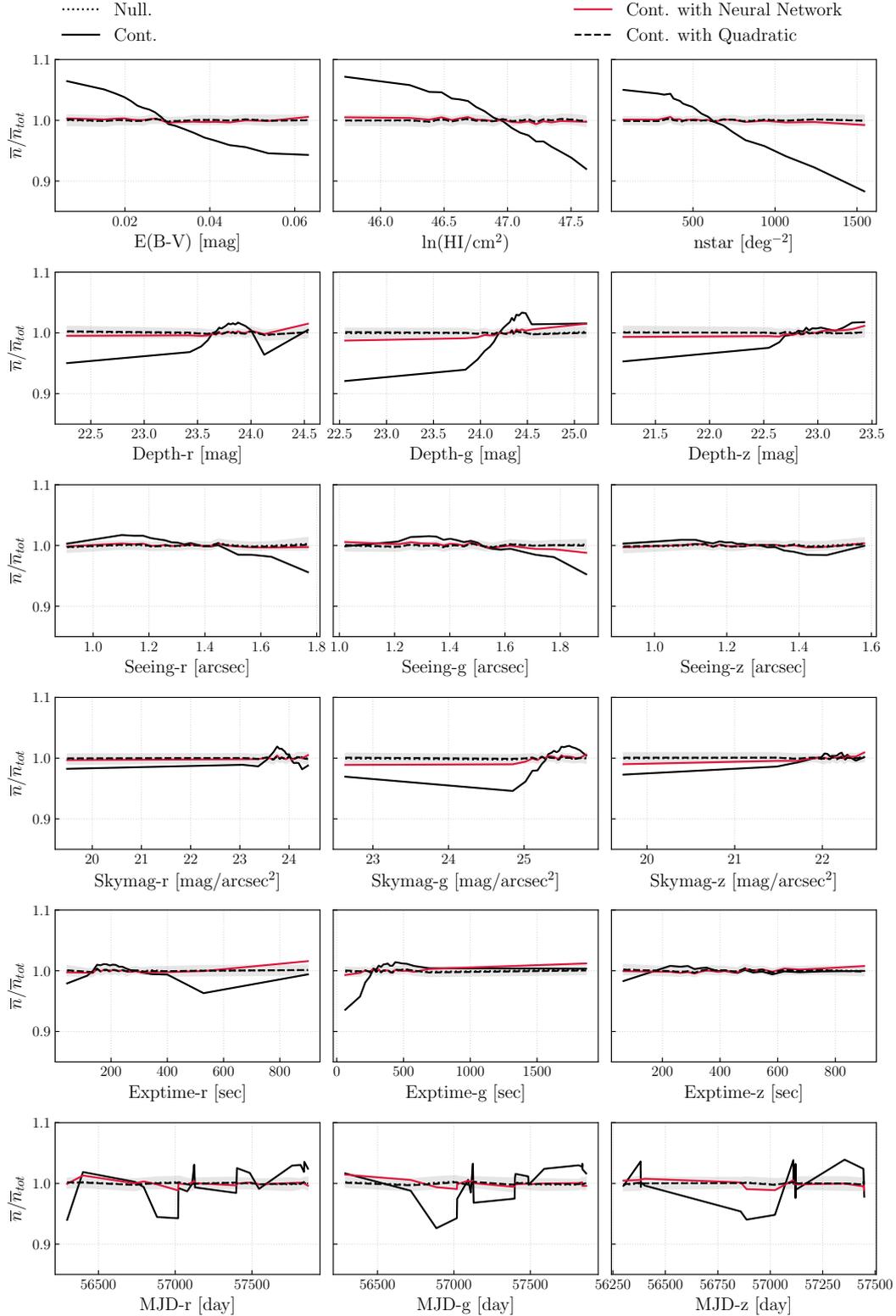


**Figure 17.** Important imaging maps identified in the mocks by the feature selection procedure for the five partitions of the 100 null (left) and contaminated mocks (right). The maps used in the input contamination model were marked by ‘\*’. The right panel shows that the feature selection procedure has identified EBV, Stellar density, skymag-g, seeing-g as important in most of the contaminated realizations whereas in the left panel, no map is consistently selected as important among all of the 100 null mocks.

from the left panel<sup>9</sup>). One can see that the contamination

<sup>9</sup> Note that we use the clustering of the null mocks before mitigation as the ground truth model since the survey window for all of the mocks before and after systematics treatment does not change.

substantially increased power at  $\ell < 50$ . Since the contamination model is based on the linear polynomial model, all three fitting methods, i.e., the linear, quadratic, and neural network, are capable of reproducing the true input contamination, while they perform differently in the presence of the two layers of noise we added and the eight additional imag-



**Figure 18.** The number density of the mock galaxies as a function of the potential systematics averaged over the 100 mock datasets. The grey shaded region illustrates  $1-\sigma$  dispersion in the null mocks. The dotted curve shows the mean density of the null mocks. The black solid curve shows the mean density dependence on the imaging attributes for the contaminated mocks. The solid red curve shows the mean density after correcting for the systematics with the neural network selection mask. The dashed black curve shows the corrected results with the quadratic polynomial model selection mask. The result of the linear polynomial model is almost unity, and therefore is omitted for clarity.

ing attributes that are non-trivially correlated with the ten input contamination attributes.

In the right top and middle panel, we find all three methods effectively remove the contamination over  $\ell < 100$ ; in detail, the linear (black dot-dashed) and the quadratic polynomial (black dashed) methods slightly over-correct power while the neural network method (solid red) slightly under-corrects it. Note that, without the feature selection process (dashed blue), the neural network method would also over-correct the large scale power like the linear and quadratic polynomial models. The left top and middle panel show that, in the absence of contamination, both the linear and quadratic polynomial methods over-correct the large scale power since the fitting methods can always find purely coincidental consistency between the imaging attributes and the cosmic variance. The quadratic method that has a greater freedom appears more prone to such problem. On the other hand, the neural network method without the feature selection process shows a lesser degree of overfitting than the linear methods, probably due to the validation procedure. Our default neural network method, which incorporates feature selection, is the most robust mitigation methodology against overfitting.

The remaining bias can be compared to the typical error expected for such data. The dark and light grey shaded regions in the middle panels indicate the  $1\text{-}\sigma$  confidence regions for the mean and the individual mock of the 100 mocks, respectively. We quantify the significance of such remaining bias by calculating  $\chi^2$ , the sum of the squared offset weighted with the diagonal variance for the mean at each  $\ell$  bin. Note that we use the variance from the 100 null mocks for calculating  $\chi^2$  of the contaminated mocks in order to avoid an advantage of the increased variance after contamination. We find that the default neural network with  $\chi^2 = 0.74$  (reduced  $\chi^2 = 0.12$  with  $dof = 6$ ) recovers the true underlying clustering when applied to the null mocks well within  $1\sigma$  C.L. of the sample variance. We estimate the significance with taking the residual systematics as one extra degree of freedom,

$$\text{max systematics} \sim \sqrt{\chi^2} \quad (22)$$

On the other hand, the linear and quadratic models have systematic biases with more than  $4\sigma$  and  $7.4\sigma$  significance. When applied to the contaminated mocks, we find that the neural network method returns  $\chi^2 = 15.3$  from six  $\ell$  bins ( $\ell \geq 12$ ), while the linear and quadratic methods, respectively, polynomial mitigation return 27.2 and 86.9. The difference in  $\chi^2$  among the different mitigation methods is not significant compared to the  $\chi^2 = 17,466.8$  before mitigation. While the neural network model appears to perform the best,  $\chi^2$  of 15.3 (reduced  $\chi^2$  of 2.6) indicates that the residual is much more significant than the sample variance. However, for the contaminated mocks, we added substantial statistical noises, almost doubling the noise. If we use the covariance of each contaminated/mitigated case, we get  $\chi^2$  of 7.7 (reduced  $\chi^2$  of 1.3) for our default case; that is our residual is at the level of the statistical noise we added in the process of contamination. These  $\chi^2$  values can be translated to  $2.8 - 3.9\sigma$  uncertainties in the residual systematics (see Eq. 22) which implies any particular

cosmological study requires a more thorough analysis of systematics to determine an estimate of the residual systematic uncertainty for the parameters of interest.

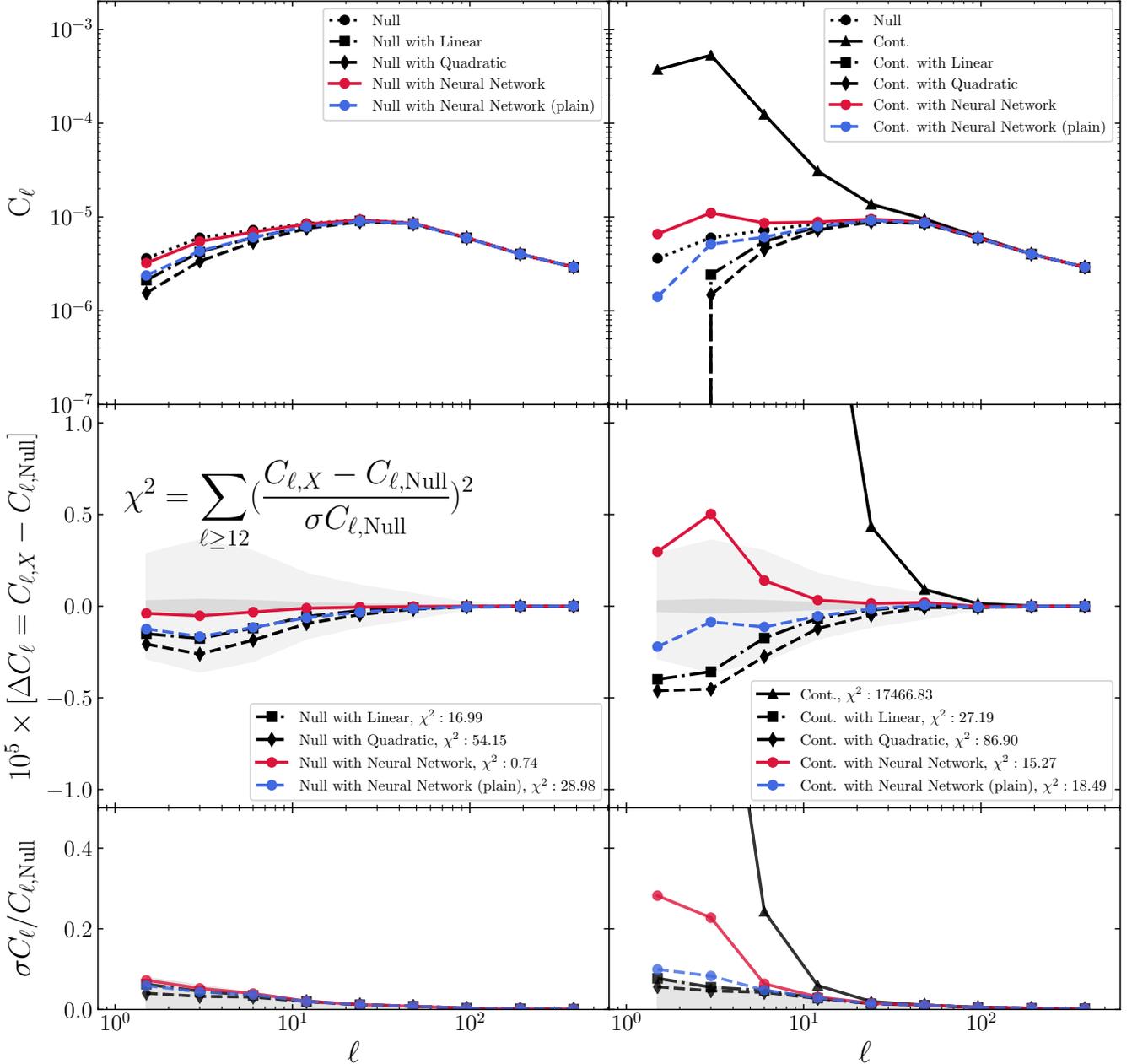
Such  $\chi^2$  can depend on the lower limit of  $\ell$  we consider. In Fig. 20, we investigate the behavior of the remaining bias depending on  $\ell_{min}$ , which shows that the neural network method consistently returns a lower remaining bias for various  $\ell_{min}$  choices. However, for the contaminated mocks, the difference is small and we conclude that all mitigation methods perform similarly for the contaminated mocks.

The bottom panels of Fig. 19 compare the noise introduced by the three different mitigation processes. The right panel shows that the contamination process (solid black) introduces additional noise on large scales relative to the variance of the null mocks (the gray shade). After mitigation, the variance is reduced, which is probably related to the decrease in the large scale power, since the variance of power is proportional to the amplitude of power itself in the Gaussian limit. The quadratic method shows the smallest fractional error on large scales due to the reduced amplitude after correction. In all cases, if we calculate the fractional variance with respect to the measured  $C_\ell$  (e.g.,  $\sigma C_{\ell,NN}/C_{\ell,NN}$  instead of  $\sigma C_{\ell,NN}/C_{\ell,Null}$ ), it agrees with the fractional variance of the null mock (gray shade). Therefore, we do not observe a nontrivial increase in variance by any of the mitigation methods we tested.

#### 4.2.4 Cross power spectrum of mock galaxies and imaging attributes

In Fig. 21, we show the mean cross power spectrum of the 100 mock catalogs and the imaging attributes for the different mitigation techniques. All three methods substantially reduce the cross power with the imaging attributes. The neural network method tends to show a small residual for  $\ell < 10$  that is greater than those of the linear and quadratic polynomial models. The dark shaded region shows the  $1\text{-}\sigma$  confidence region of the mean cross power propagated to  $C_{s,g}^2/C_{s,s}$  and the light shaded region shows the  $1\sigma$  confidence interval of the mean auto-power spectrum of the mocks as shown in Fig. 14. Therefore, we find that these residuals are greater than the statistical noise of  $C_{s,g}^2/C_{s,s}$ , but the effect on the *auto power spectrum* are marginal for  $\ell > 10$ . This excess on small  $\ell$  is partly due to the greater auto galaxy power spectrum amplitude (in Fig. 19) after mitigation than those by the other methods. Meanwhile we still find a residual correlation with skymag-z and mjd-z even after accounting for the effect of the auto power spectrum amplitude. Without the feature selection procedure (dashed blue), the neural network also returns a smaller residual. In essence, we see that once feature selection is applied, the neural network only corrects to a certain level, controlled by the specifics of the feature selection procedure. This protects against over-fitting due to random correlations between the imaging attributes and the galaxy density field.

As a sanity check, if we use the true input linear contamination model to mitigate the systematics (purple dot-dashed line in Fig. 21), the cross-correlation completely



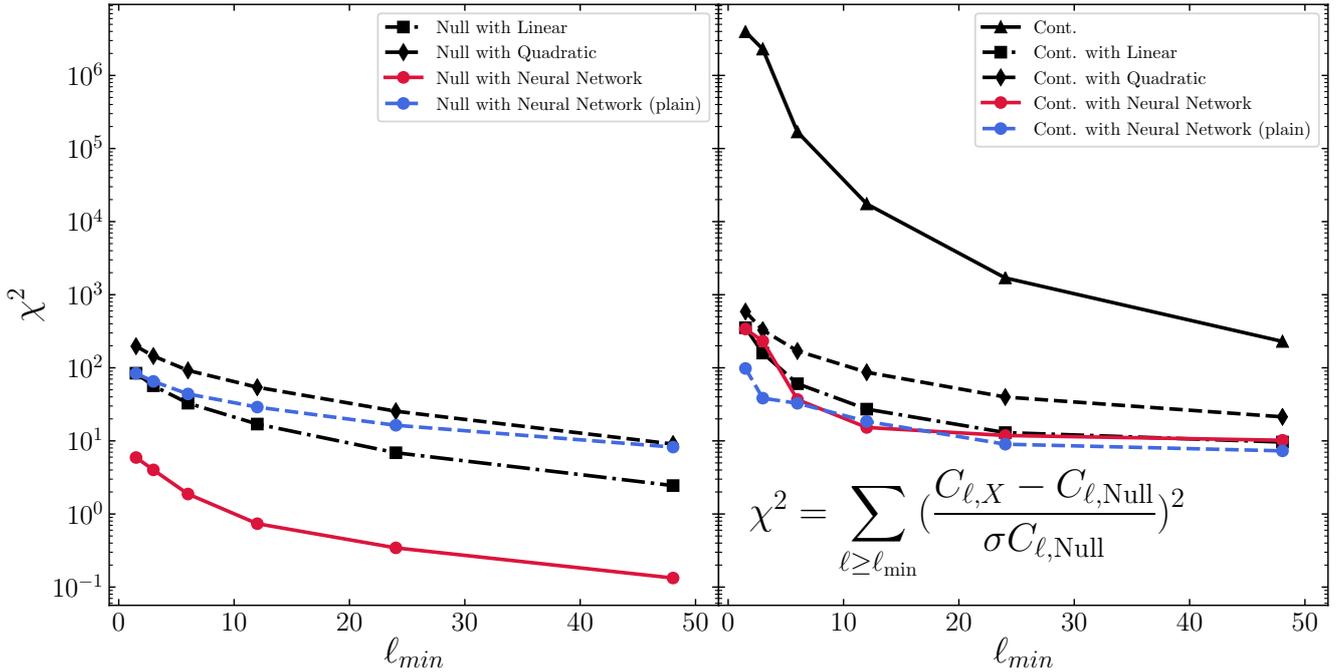
**Figure 19.** *Top row:* The mean angular power spectrum of the 100 contaminated (null) mocks in the right (left) panel. *Middle row:* The mean power spectrum subtracted by the mean of the null mocks to better visualize the remaining bias after each mitigation. The dark grey shaded region shows the  $1-\sigma$  confidence region of the mean of 100 mocks, while the light grey area shows the  $1-\sigma$  confidence region for one mock, calculated from the dispersion of 100 mocks. To account for the increased shot noise during contamination, we remove the same constant power from all contaminated/mitigated power spectra until their small scale power matches that of the null mock power spectrum. We quantify the significance of the remaining bias by calculating  $\chi^2$ , the sum of the squared offset weighted with the diagonal variance of the mean  $C_\ell$  of null mocks over the last six bins ( $\ell \geq 12$ ). The middle panel on the left illustrates the neural network without feature selection ('plain') tends to remove the cosmological clustering signal.

vanishes as expected<sup>10</sup>. For the null mocks, all mitigation

methods return negligible cross-correlation, which is omitted from Fig. 21 for clarity.

<sup>10</sup> The auto-power spectrum of the contaminated mocks mitigated with the true input contamination model (in Fig. 22) returns the smallest residual bias relative to the uncontaminated clustering, as expected.

Overall, while the cross-correlation statistics between the galaxy density and the systematics attributes are a useful indicator for the level of contamination, we find it may be difficult to infer and discriminate the level of contamina-



**Figure 20.** Dependence of  $\chi^2$  on the lowest bin  $\ell_{\min}$  in the null (left) and contaminated (right) mocks. To better quantify the residual bias introduced by each method, we evaluate the dependence of the bias on the lowest bin  $\ell_{\min}$  that is included in the  $\chi^2$  computation. The default neural network method performs significantly better than the conventional methods for the null mocks, mainly because the feature selection procedure successfully prevents the method from regressing out the cosmological clustering signal. For the contaminated mocks, all methods tend to perform similarly, as expected, since all mitigation methods can reproduce the input contamination model.

tion in the density field from such cross-correlation statistics beyond what can be probed by the auto power spectrum.

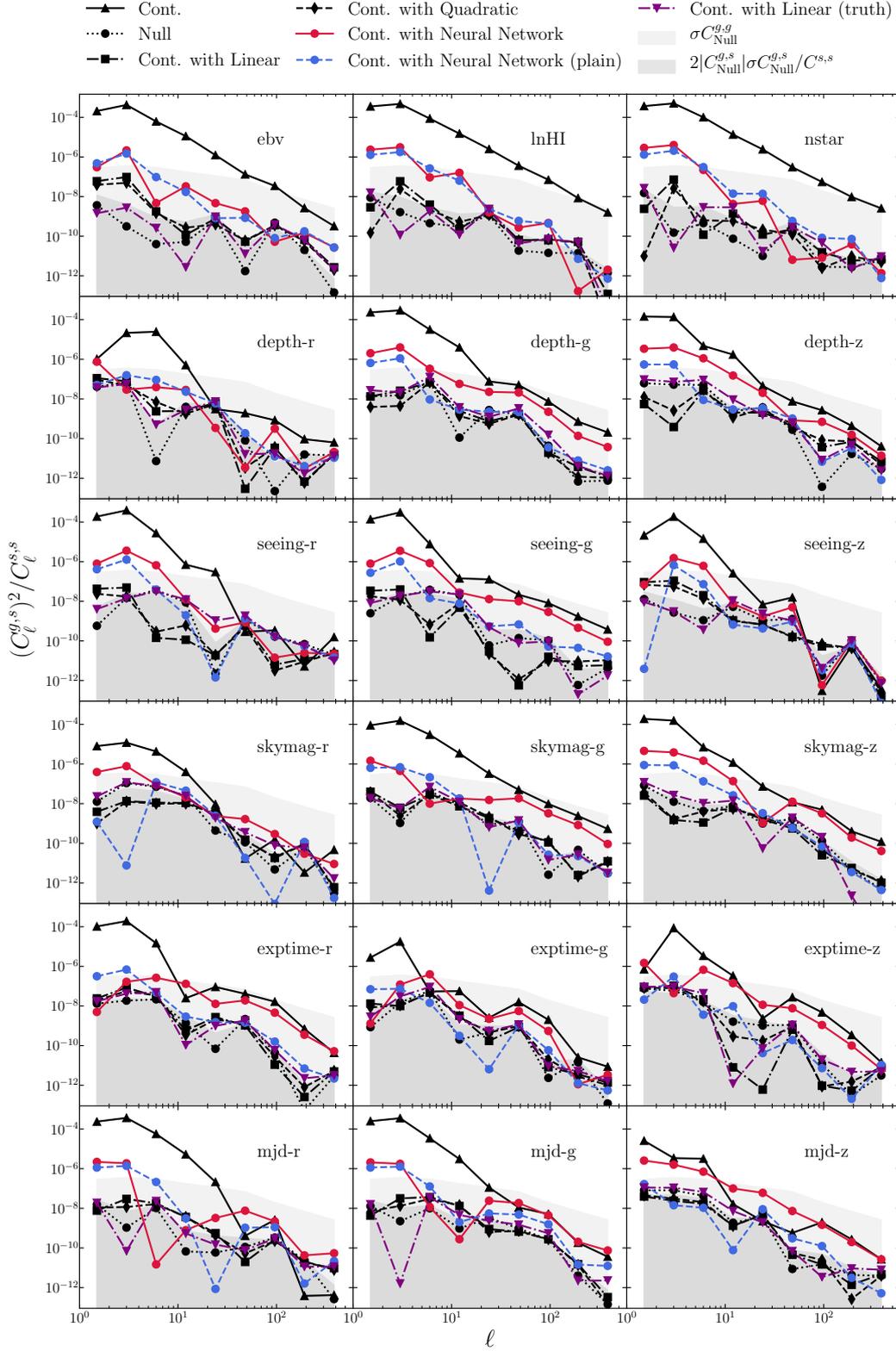
#### 4.2.5 A case with underfitting

It is possible that we may identify only a subset of the contamination attributes for a given data set and attempt to mitigate the contamination based on such limited information. We consider a situation where we input only five imaging attributes to the mitigation procedure: four from the true contamination inputs, i.e., *EBV*, *lnHI*, *nstar*, *skymag-g* and one that is not among the true contamination inputs, but correlated with the contamination inputs, i.e., *depth-r*. The neural network method could be more resilient to such limited information since its nonlinear activation function may allow the mitigation procedure to better utilize the correlation between different input imaging attributes. In Fig. 22, the linear polynomial method with ‘few’ inputs shows a lesser degree of overfitting for the null mocks, compared to the default linear case, while showing under-fitting for the contaminated mocks. This is expected as the freedom of the linear model is now limited. The neural network method with the ‘few’ inputs (without feature selection) returns a very similar pattern as the linear ‘few’ case, which implies that our current default neural network method does not have an advantage over the linear model in such a case despite its greater flexibility. This underfitting case may be worth future investigation, as apparent excess clustering remains in DR7 in Fig. 16 when any method is applied, but

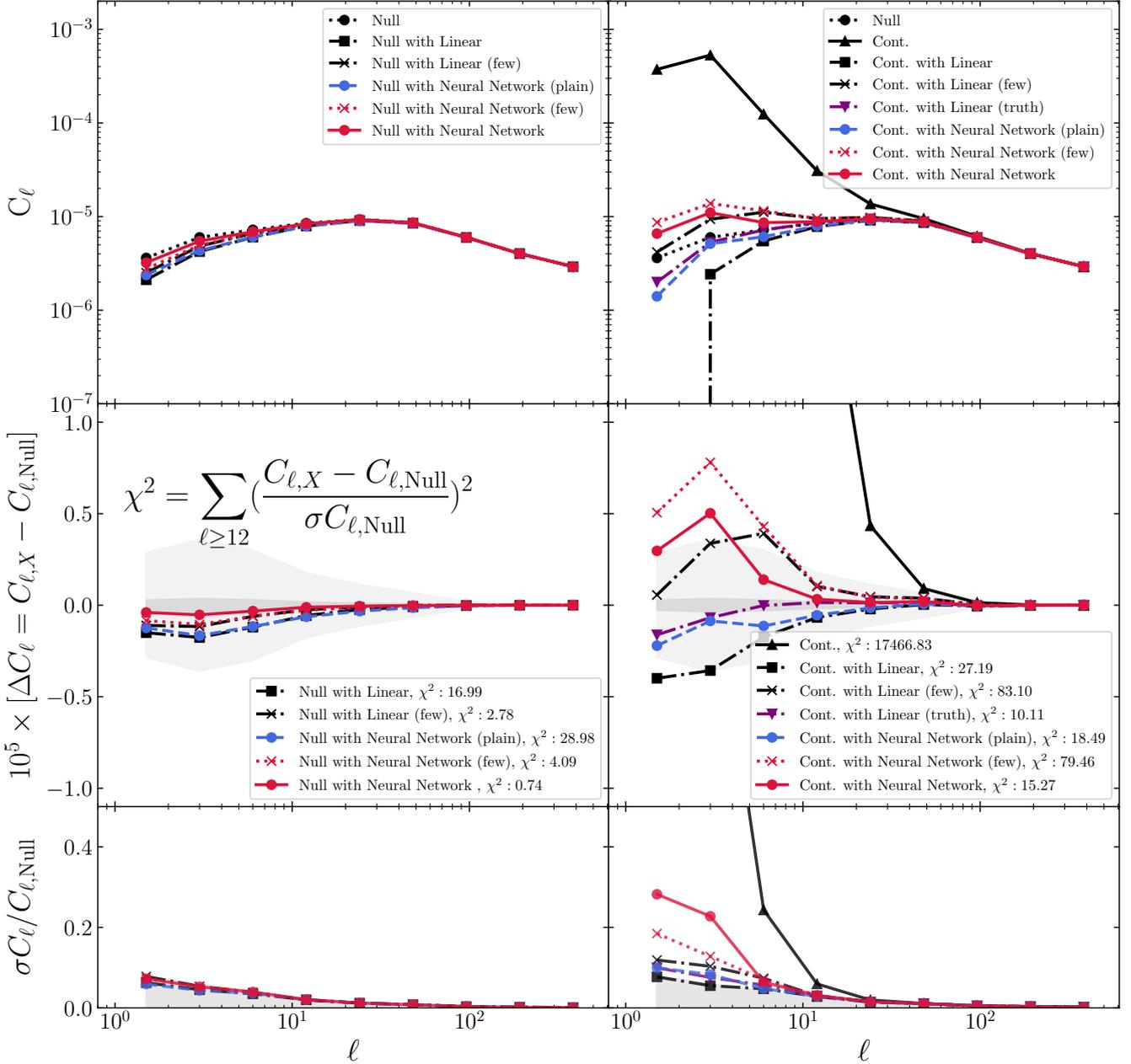
especially in the case of the application of multivariate linear models.

### 4.3 Summary and Discussion

In summary, we find that our default neural network method is more robust against overfitting based on the test with the null mocks. This is due to the feature selection process that appropriately reduces the flexibility of mitigation. Based on the tests with the contaminated mocks, we find that both the linear and neural network methods perform equally well in terms of the residual bias, while the neural network method is more robust against overfitting. The quadratic polynomial method appears to be more prone to the overfitting problem than the other two methods since it has a greater flexibility than the input contamination model, but without a way to suppress the flexibility. All methods do not increase fractional variance during the mitigation process. Note again that we deliberately choose the linear model in contaminating the mocks in this test in order to prevent a disadvantage in using the linear and quadratic polynomial methods. Therefore, the decent performance of the linear method is warranted. In real data, the contamination due to observational effects can take a more complex form as implied by the difference in the mitigation results between the data and our mocks. Therefore our mock test is a conservative estimation of the comparative mitigation capability of our default neural network method.



**Figure 21.** The mean cross power spectrum of the contaminated mock catalogs and the imaging attributes for different mitigation techniques. Our default neural network method with feature selection is shown in solid red while the performance without feature selection (‘Neural Network (Plain)’) is shown in dashed blue. The dark grey shaded region shows the  $1\sigma$  confidence region of the plotted quantity derived from  $2\sigma(C_{g,s}) * C_{g,s}/C_{s,s}$ , and the light grey region shows the typical  $1\sigma$  confidence region of the mean auto power spectrum of the 100 mocks. The mitigation with the ground truth contamination model is shown with a purple dot-dashed curve as ‘cont. with linear (truth)’.



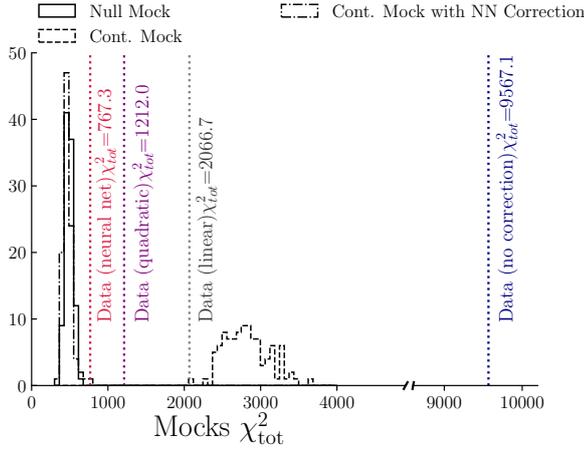
**Figure 22.** Same as Fig. 19 showing the mean auto power spectrum of the 100 null (left) and contaminated (right) mocks mitigated with the fewer imaging maps ‘few’ (to demonstrate under-correction), neural network without the feature selection ‘plain’, and the ground truth contamination model ‘truth’. The mitigation with the ground truth model achieves the lowest residual bias as expected. For the null mocks, using fewer imaging maps prevents over-fitting simply by providing less freedom in the regression model while it leads to underfitting for the contaminated mocks.

While we demonstrated qualitatively and quantitatively that our fiducial neural network method is more robust than the conventional methods, for both DR7 as well as for the mocks, Fig. 13-14 show non-negligible residual contamination compared to the expectations. It is not surprising that the mitigation of imaging systematic effects in the real data is more challenging than that of the linear-model based systematics in our mock tests. In this subsection we attempt to provide a quantitative evaluation

of the residual systematics of DR7.

We quantify the residual systematics in the mean density against the 18 imaging maps using  $\chi^2$  of the mean density diagnostic as in Fig. 13 and Table 4<sup>11</sup>. The null hypothesis is that the total residual squared error of the

<sup>11</sup> To compare these with the mock results, we limit DR7 to the mock footprint, and therefore the  $\chi^2$  results are slightly different



**Figure 23.** Left:  $\chi^2$  distribution for the null mocks (solid black), contaminated mocks (dashed black), and contaminated mocks with NN mitigation (dot-dashed black). The vertical dotted lines overlay the  $\chi^2$  values for the data with the default NN (red), quadratic (purple), and linear (grey) treatments. The  $\chi^2$  statistics before treatment is shown on the right (dark blue).

mean density observed in the data should be consistent with the distribution of the  $\chi^2$  statistics constructed with the null mocks. The vertical lines in Fig. 23 present the  $\chi^2$  values observed in the data for before systematics treatment (9567.1) and after treatment with linear (2066.7), quadratic (1212.0), default Neural Network (767.3), and Neural Network plain (744.4) approaches. As a comparison, we present the distributions of the  $\chi^2$  observed in the null mocks (solid), contaminated mocks (dashed), and contaminated mocks after neural network mitigation (dot-dashed). Fig. 23 illustrates a factor of 12 improvement in terms of the residual  $\chi^2$  when using the neural network. Compared to the conventional quadratic method ( $\chi^2 = 1212$ ), the NN-based method makes a factor of 1.6 improvement.

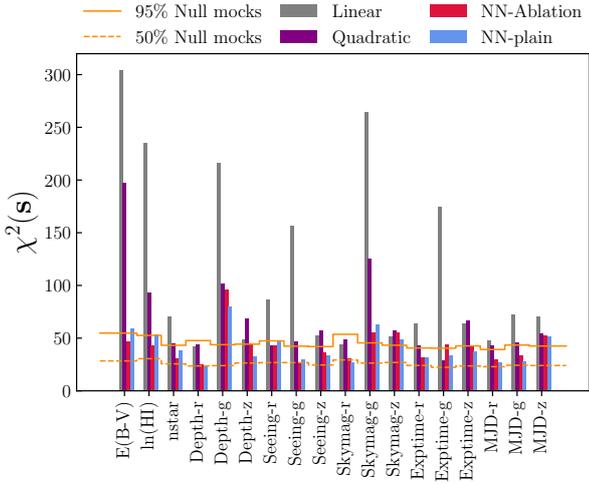
The mean and standard deviation for the distribution of  $\chi^2$  values observed in the null mocks are  $487.1 \pm 4.92$ , respectively. The same quantities observed in the contaminated mocks are  $2819.0 \pm 29.30$ , while Neural Network mitigation decreases these statistics to the mean of  $472.0 \pm 6.59$ . We perform Welch’s t-test (Welch 1947) on the  $\chi^2$  distributions of the null mocks and the mitigated contaminated mocks, and conclude that the two distributions have identical means with  $t$ -statistic = 1.8,  $p$ -value = 0.06 and the neural network recovers the underlying cosmological mean target density. Note that the mean  $\chi^2$  observed in the contaminated mocks, 2819.0, is much smaller than the value observed in the data (the blue vertical line, 9567.1). This indicates that the systematic effects rigorously simulated in the mocks are still not as strong as the systematics in the real data or the effect of the non-diagonal terms, that we ignored in the covariance, is different in the data and in the mocks.

from Fig. 13 and Table 4. The mock footprint is smaller than the data footprint by almost a factor of two.

We perform a hypothesis testing given the  $\chi^2$  value observed in the data after the neural network mitigation and the distribution of  $\chi^2$  values observed in the 100 null mocks. If the imaging maps were independent and normal deviates, the  $\chi^2$  values must follow a Chi-squared distribution for  $dof = 360$  where 360 is the total number of bins. This assumption is not necessarily satisfied given the correlation among imaging systematics, and we indeed find  $\chi^2$  of the null mocks is better fit with  $dof = 487$ . Therefore, we assume that the underlying  $dof$  for evaluating  $\chi^2$  is 487. For DR7, the  $\chi^2$  decreased from 9567 before contamination to 767.3 after our NN-based mitigation. Compared to  $dof = 487$  we expect for the approximate truth, this is a substantial excess, being very unlikely due to random fluctuation: for example, the top 5% of distribution with  $dof = 487$  is at  $\chi^2 = 539.4$ , which is 40% lower than 767.3. Assuming a mere random fluctuation is bounded at this upper 5% and the  $\chi^2$  values of the data are consistent with that of the mocks, this offset could imply that we underestimated the covariance at least by  $\sqrt{(767.3/539.4)} = 19\%$  or that our model, when applied to DR7, missed at least 19% of the systematic effects in the target density.

We further investigate the contribution of each imaging map to the residual systematics in DR7. Fig. 24 presents the  $\chi^2$  vs each imaging map after the linear (grey), quadratic (purple), neural network with feature selection (red), and neural network plain (blue) treatments. The statistics before mitigation are not shown for clarity. We also plot the 50- and 95-th percentiles of the same quantity observed in the null mocks in orange horizontal curves. Depth-g seems to be the main source of the residual systematics which is not mitigated even by the neural network-based methods. All methods also show residual systematics against skymag-g, skymag-z, and MJD-z. The standard treatments show some additional residual systematics against E(B-V), lnHI, and exptime-z. We perform further tests by masking out high extinction and low depth-g regions. We find that applying a more rigorous cut on depth-g (e.g.,  $\text{depth-g} > 24.95$ ) yields a more cleaner mean density at the expense of losing 9% of the data. Although our objective was to avoid applying rigorous cuts on imaging and demonstrate the gain by using non-linear methods, our tests suggest that careful masking based on imaging, particularly depth-g, seems necessary for cosmological studies.

In summary, our results indicate that any cosmological study requires a thorough analysis of systematics to determine an estimate of the residual systematic uncertainty for the parameters of interest. In order to use this sample for a cosmological exploitation, we would need to account for 19% (or more, depending on the assumption on the baseline) additional systematic errors to the statistical errors in the density field level. Our analysis also suggests that additional masking, especially based on depth-g, or improving the method to deal with depth and sky background issues would be the next steps to prepare this data for cosmological analysis. We leave this for future study, as the focus of this study is to compare our non-linear, neural network method to other map-based methods.



**Figure 24.** The breakdown of  $\chi^2$  values observed in the data on the mock footprint after linear (grey), quadratic (purple), neural network with feature selection (red), and neural network plain (blue) treatments. We also plot the 50- and 95-th percentiles of the same quantity observed in the null mocks in orange curves (note that  $\chi^2(s)$  is not a continuous quantity). Given the 5% threshold, we can argue that there exists known residual systematics against E(B-V), depth-g, skymag-g, skymag-z, and MJD-z.

## 5 CONCLUSION

In this paper, we have presented a rigorous application of an artificial neural network methodology to the mitigation of the observational systematics in galaxy clustering measurements of an eBOSS-like ELG sample selected from DR7 (see § 2). We have investigated the galaxy density dependency on 18 imaging attributes of the data (see Fig. 1). We compare the performance of the neural network with that of the traditional, linear and quadratic multivariate regression methods. The key aspects of our neural network methodology are:

- The application of k-fold cross-validation, which implements the training-validation-test split to tune the hyper parameters by evaluating how well the trained network generalizes to the unseen, validation data set and therefore to suppress overfitting when applied to the test set;
- The repeated split process until we cover the entire data footprint as test sets;
- The elimination of redundant imaging maps by the feature selection procedure to further reduce the overfitting problem and therefore protect the cosmological clustering signal.

We apply the output of our pipeline, i.e., the selection mask for the DR7 footprint to the observed galaxy density field. Benchmark selection masks are also produced employing the linear and quadratic polynomial regression. Comparing statistical results before and after applying the selection masks, we find that:

- Galactic foregrounds are the most dominant source

of contamination in this imaging dataset (see Figs. 13, 14, and 15).

- This contamination causes an excess clustering signal in the auto power spectrum and correlation function of the galaxy density field on large scales (see Fig. 16).
- All mitigation techniques e.g., the neural network method as well as the linear multivariate models using the linear and quadratic polynomial functions, are able to reduce the auto and cross clustering signals (see Figs. 15 and 14);
- However, the neural network removes the excess clustering more effectively in the auto power spectrum and correlation function of galaxies (see Fig. 16).

The last result implies that our neural network method has a higher flexibility than both linear multivariate models we tested, and it is therefore capable of capturing the non-linear systematic effects in the observed galaxy density field.

We apply our methodology on two sets of 100 log-normal mock datasets with (‘contaminated mocks’) and without (‘null mocks’) imaging contamination to evaluate how well the ground truth cosmological clustering can be reconstructed in both cases, and therefore to validate the systematic mitigation techniques. All mitigation techniques are applied in the same way we treat the real data. The key results of our mock test are as follows:

- The feature selection procedure is able to identify most of the ten contamination input maps as important for the contaminated mocks while correctly identifying most of the maps as redundant for the null mocks (see Fig. 17).
- All three mitigation methods, i.e., the linear polynomial, quadratic polynomial, and neural network methods, perform similarly in terms of the residual bias in the presence of contamination. This is expected since the contamination model is based on the linear polynomial model which all three methods are capable of reproducing. The default neural network tends to slightly under-correct which is the outcome of the feature selection procedure. On the other hand, the linear and quadratic polynomial methods tend to slightly over-correct (see the right panel of Fig. 19).;
- In the absence of contamination, the neural network is the most robust against regressing out the cosmological clustering. This is mainly due to the feature selection process that appropriately reduces the flexibility of the mitigation (see the left panel of Fig. 19). Based on this result, we implement the feature selection procedure for DR7.
- Using  $\chi^2$  statistics, we quantify the bias and find that for the null mocks, the default neural network recovers the underlying clustering within  $1\sigma$  C.L. (see Eq. 22) while the other methods return more than  $4\sigma$  C.L. bias. For the contaminated mocks, all of the methods return biased clustering with  $2.8 - 3.9\sigma$  C.L. which indicates that it is crucial for cosmological parameter estimation to determine the residual systematic uncertainty in scales sensitive to the parameters of interest (see the middle panels of Figs 19 and 22).

- All methods do not increase fractional variance during the mitigation process (see the bottom row of Fig. 19).

We also employ the mocks to investigate the remaining systematic effects in the data (Figs 13-15). While the neural network methods outperform the conventional methods, we conclude that the data exhibit around 19% residual systematics in the target number density (Figs 23 & 24). Our analysis suggests that a more rigorous masking on  $depth-g$  (e.g.,  $depth-g > 24.95$ ) improves the mean density at the cost of losing 9% of data. To use this sample for cosmology, we therefore suggest a) accounting for 19% (or more, depending on the assumption on the baseline) additional systematic errors to the statistical errors in the density field level; b) performing further analysis of systematics and improvement of the mitigation method to deal with the depth and sky background issues.

To conclude, our analyses illustrate that the neural network method we developed in this paper is a promising tool for the mitigation of the large-scale spurious clustering that is likely raised by the imaging systematics. Our method is more robust against regressing out the cosmological clustering than the traditional, linear multivariate regression methods. Such improvement will be particularly crucial for an accurate measurement of non-Gaussianity from the large-scale clustering of current eBOSS and upcoming DESI and the LSST surveys. Our method is computationally less intensive than other approaches such as the Monte Carlo injection of fake galaxies: analyzing DR7 using our default neural network method requires less than six CPU hours. Application of our methodology on any imaging dataset would be straightforward. Our systematics mitigation methodology pipeline is publicly available at <https://github.com/mehdirezaie/SYSNet>.

## ACKNOWLEDGEMENTS

We would like to thank the anonymous referee for constructive comments. M.R. and H.-J.S. are supported by the U.S. Department of Energy, Office of Science, Office of High Energy Physics under Award Number DE-SC0014329. This research used the Dark Energy Spectroscopic Instrument (DESI) allocation resources of the National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility operated under Contract No. DE-AC02-05CH11231. M.R. would like to thank Stephen Bailey, NERSC, and the DESI collaboration for providing computing resources to carry out the initial phases of this work; Anand Raichoor for providing the stellar mask and discussions on targetting; Nick Hand and Yu Feng for helping with Nbodykit and mock generations; Jeffrey Newman, Arnaud De-Mattia, and Marc Manera for discussions on machine learning and cross-validation; Ted Kisner, Rollin Thomas, and Joel Brownstein for helping with high-performance computing; Dustin Lang, John Moustakas, and Arjun Dey for discussions about the Legacy Surveys and data reduction pipelines. M.R. and H.J. would like to thank Patrick McDonald for discussions on window effects and error analyses. AJR is grateful for support from

the Ohio State University Center for Cosmology and Particle Physics. We would like to appreciate the open-source software and modules that were invaluable to this research: HEALPix, Fitsio, Tensorflow, Scikit-Learn, NumPy, SciPy, Nbodykit, Pandas, IPython, Jupyter, GitHub.

## REFERENCES

- Ade P. A., et al., 2016, *Astronomy & Astrophysics*, 594, A13  
 Aghamousa A., et al., 2016, arXiv preprint arXiv:1611.00036  
 Ahn C. P., et al., 2012, *The Astrophysical Journal Supplement Series*, 203, 21  
 Akrami Y., et al., 2018, arXiv preprint arXiv:1807.06205  
 Alam S., et al., 2017, *MNRAS*, 470, 2617  
 Ata M., et al., 2017, *Monthly Notices of the Royal Astronomical Society*, 473, 4773  
 Bautista J. E., et al., 2018, *ApJ*, 863, 110  
 Bekhti N. B., et al., 2016, *Astronomy & Astrophysics*, 594, A116  
 Bergé J., Gamper L., Réfrégier A., Amara A., 2013, *Astronomy and Computing*, 1, 23  
 Breiman L., 2001, *Machine learning*, 45, 5  
 Brown A., Vallenari A., Prusti T., de Bruijne J., Babusiaux C., Bailer-Jones C., Collaboration G., et al., 2018, arXiv preprint arXiv:1804.09365  
 Chon G., Challinor A., Prunet S., Hivon E., Szapudi I., 2004, *MNRAS*, 350, 914  
 Coles P., Jones B., 1991, *MNRAS*, 248, 1  
 Colless M., et al., 2001, *Monthly Notices of the Royal Astronomical Society*, 328, 1039  
 Crocce M., et al., 2016, *MNRAS*, 455, 4301  
 Cybenko G., 1989, *Mathematics of control, signals and systems*, 2, 303  
 Dahl G. E., Sainath T. N., Hinton G. E., 2013, in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. pp 8609–8613  
 Dalal N., Dore O., Huterer D., Shirokov A., 2008, *Physical Review D*, 77, 123514  
 Dark Energy Survey Collaboration: Fermilab University of Illinois at Urbana-Champaign U. o. C. L. B. N. L. C.-T. I.-A. O., Flaugher B., 2005, *International Journal of Modern Physics A*, 20, 3121  
 Delubac T., et al., 2016, *Monthly Notices of the Royal Astronomical Society*, p. stw2741  
 Devijver P. A., Kittler J., 1982, *Pattern recognition: A statistical approach*. Prentice hall  
 Dey A., et al., 2018, arXiv preprint arXiv:1804.08657  
 Drinkwater M. J., et al., 2010, *Monthly Notices of the Royal Astronomical Society*, 401, 1429  
 Efstathiou G., Ellis R. S., Peterson B. A., 1988, *Monthly Notices of the Royal Astronomical Society*, 232, 431  
 Eisenstein D. J., 2005, *New Astronomy Reviews*, 49, 360  
 Eisenstein D. J., Hu W., Tegmark M., 1998, *The Astrophysical Journal Letters*, 504, L57  
 Elsner F., Leistedt B., Peiris H. V., 2015, *Monthly Notices of the Royal Astronomical Society*, 456, 2095  
 Elvin-Poole J., et al., 2018, *Phys. Rev. D*, 98, 042006  
 Fisher K. B., Davis M., Strauss M. A., Yahil A., Huchra J. P., 1993, *The Astrophysical Journal*, 402, 42  
 Funahashi K.-I., 1989, *Neural networks*, 2, 183  
 Geurts P., Ernst D., Wehenkel L., 2006, *Machine learning*, 63, 3  
 Glorot X., Bordes A., Bengio Y., 2011, in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. pp 315–323  
 Gorski K. M., Hivon E., Banday A., Wandelt B. D., Hansen F. K., Reinecke M., Bartelmann M., 2005, *The Astrophysical Journal*, 622, 759

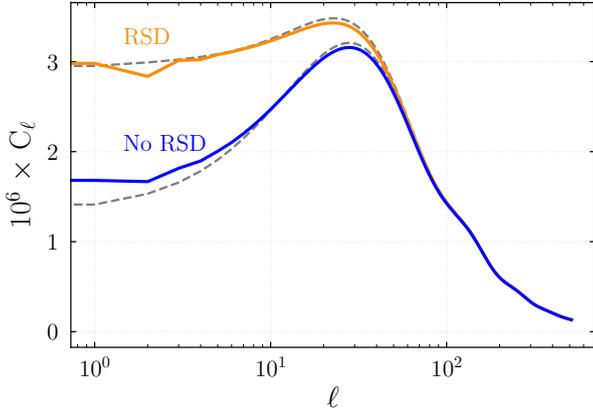
- Guyon I., Elisseeff A., 2003, *Journal of machine learning research*, 3, 1157
- Hamilton A., 1998, in , *The evolving universe*. Springer, pp 185–275
- Hand N., Feng Y., Beutler F., Li Y., Modi C., Seljak U., Slepian Z., 2017, arXiv preprint arXiv:1712.05834
- He K., Zhang X., Ren S., Sun J., 2015, in *Proceedings of the IEEE international conference on computer vision*. pp 1026–1034
- Hivon E., Górski K. M., Netterfield C. B., Crill B. P., Prunet S., Hansen F., 2002, *ApJ*, 567, 2
- Ho S., Hirata C., Padmanabhan N., Seljak U., Bahcall N., 2008, *Physical Review D*, 78, 043519
- Ho S., et al., 2012, *APJ*, 761, 14
- Hoerl A. E., Kennard R. W., 1970, *Technometrics*, 12, 55
- Hornik K., Stinchcombe M., White H., 1989, *Neural networks*, 2, 359
- Huang G.-B., 2003, *IEEE Transactions on Neural Networks*, 14, 274
- Huterer D., Cunha C. E., Fang W., 2013, *Monthly Notices of the Royal Astronomical Society*, 432, 2945
- Ivezic Z., et al., 2008, arXiv preprint arXiv:0805.2366
- John G. H., Kohavi R., Pflieger K., 1994, in , *Machine Learning Proceedings 1994*. Elsevier, pp 121–129
- Jones D., et al., 2018, *The Astrophysical Journal*, 857, 51
- Kaiser N., 1987, *Monthly Notices of the Royal Astronomical Society*, 227, 1
- Kalus B., Percival W. J., Bacon D., Samushia L., 2016, *Monthly Notices of the Royal Astronomical Society*, 463, 467
- Kalus B., Percival W., Bacon D., Mueller E., Samushia L., Verde L., Ross A., Bernal J., 2018, *Monthly Notices of the Royal Astronomical Society*, 482, 453
- Kingma D. P., Ba J., 2014, arXiv preprint arXiv:1412.6980
- Kohavi R., John G. H., 1997, *Artificial intelligence*, 97, 273
- Koller D., Sahami M., 1996, Technical report, *Toward optimal feature selection*. Stanford InfoLab
- Krizhevsky A., Sutskever I., Hinton G. E., 2012, in *Advances in neural information processing systems*. pp 1097–1105
- LSST Science Collaborations et al., 2017, preprint, ([arXiv:1708.04058](https://arxiv.org/abs/1708.04058))
- Lang D., Hogg D. W., Mykytyn D., 2016, *Astrophysics Source Code Library*
- Laurent P., et al., 2017, *Journal of Cosmology and Astroparticle Physics*, 2017, 017
- Leistedt B., Peiris H. V., 2014, *Monthly Notices of the Royal Astronomical Society*, 444, 2
- Leistedt B., Peiris H. V., Mortlock D. J., Benoit-Lévy A., Pontzen A., 2013, *Monthly Notices of the Royal Astronomical Society*, 435, 1857
- Leistedt B., et al., 2016, *ApJS*, 226, 24
- Lin H. W., Tegmark M., Rolnick D., 2017, *Journal of Statistical Physics*, 168, 1223
- Mather J. C., et al., 1994, *The Astrophysical Journal*, 420, 439
- Montufar G. F., Pascanu R., Cho K., Bengio Y., 2014, in *Advances in neural information processing systems*. pp 2924–2932
- Mukhanov V. F., Feldman H. A., Brandenberger R. H., 1992, *Physics Reports*, 215, 203
- Myers A. D., Brunner R. J., Richards G. T., Nichol R. C., Schneider D. P., Bahcall N. A., 2007, *The Astrophysical Journal*, 658, 99
- Nair V., Hinton G. E., 2010, in *Proceedings of the 27th international conference on machine learning (ICML-10)*. pp 807–814
- Peebles P., 1973, *The Astrophysical Journal*, 185, 413
- Perlmutter S., et al., 1999, *The Astrophysical Journal*, 517, 565
- Ponthieu N., Grain J., Lagache G., 2011, *A&A*, 535, A90
- Prakash A., et al., 2016, *The Astrophysical Journal Supplement Series*, 224, 34
- Pullen A. R., Hirata C. M., 2013, *Publications of the Astronomical Society of the Pacific*, 125, 705
- Raichoor A., et al., 2017, *MNRAS*, 471, 3955
- Ramaswamy S., et al., 2001, *Proceedings of the National Academy of Sciences*, 98, 15149
- Riess A. G., et al., 1998, *The Astronomical Journal*, 116, 1009
- Rolnick D., Tegmark M., 2017, arXiv preprint arXiv:1705.05502
- Ross A. J., Brunner R. J., Myers A. D., 2007, *The Astrophysical Journal*, 665, 67
- Ross A. J., et al., 2011, *Monthly Notices of the Royal Astronomical Society*, 417, 1350
- Ross A. J., et al., 2012, *MNRAS*, 424, 564
- Ross A. J., et al., 2013, *Monthly Notices of the Royal Astronomical Society*, 428, 1116
- Ross A. J., et al., 2017, *MNRAS*, 464, 1168
- Ruder S., 2016, arXiv preprint arXiv:1609.04747
- Rybicki G. B., Press W. H., 1992, *The Astrophysical Journal*, 398, 169
- Sánchez A. G., Baugh C. M., Angulo R., 2008, *Monthly Notices of the Royal Astronomical Society*, 390, 1470
- Schlafly E. F., Finkbeiner D. P., 2011, *The Astrophysical Journal*, 737, 103
- Schlegel D. J., Finkbeiner D. P., Davis M., 1998, *The Astrophysical Journal*, 500, 525
- Scranton R., et al., 2002, *The Astrophysical Journal*, 579, 48
- Seo H.-J., Eisenstein D. J., 2003, *The Astrophysical Journal*, 598, 720
- Slosar A., Seljak U., Makarov A., 2004, *Physical Review D*, 69, 123003
- Smoot G. F., et al., 1992, *The Astrophysical Journal*, 396, L1
- Suchyta E., et al., 2016, *MNRAS*, 457, 786
- Szapudi I., Prunet S., Pogosyan D., Szalay A. S., Bond J. R., 2001, *ApJ*, 548, L115
- Tamura S., Tateishi M., 1997, *IEEE Transactions on Neural Networks*, 8, 251
- Tegmark M., 1997, *Physical Review D*, 55, 5895
- Tegmark M., Hamilton A. J., Strauss M. A., Vogeley M. S., Szalay A. S., 1998, *The Astrophysical Journal*, 499, 555
- Thomas S. A., Abdalla F. B., Lahav O., 2011a, *Physical review letters*, 106, 241301
- Thomas S. A., Abdalla F. B., Lahav O., 2011b, *Monthly Notices of the Royal Astronomical Society*, 412, 1669
- Welch B. L., 1947, *Biometrika*, 34, 28
- White M., Tinker J. L., McBride C. K., 2013, *Monthly Notices of the Royal Astronomical Society*, 437, 2594
- Wright E. L., et al., 2010, *The Astronomical Journal*, 140, 1868
- York D. G., et al., 2000, *The Astronomical Journal*, 120, 1579
- Zou H., et al., 2017, *Publications of the Astronomical Society of the Pacific*, 129, 064101

## APPENDIX A: WINDOW FUNCTION

The observed density field of targets does not cover the full sky, due to the galactic plane obscuration. This means that the pseudo-power spectrum  $\hat{C}_\ell$  obtained by the direct Spherical Harmonic Transforms of a partial sky map (eq. 21), differs from the full-sky angular spectrum  $C_\ell$ . However, their ensemble average is related by (Hivon et al. 2002; Ponthieu et al. 2011)

$$\langle \hat{C}_\ell \rangle = \sum_{\ell'} M_{\ell\ell'} \langle C_{\ell'} \rangle, \quad (\text{A1})$$

where  $M_{\ell\ell'}$  represents the mode-mode coupling from the partial sky coverage. This is known as the Window Function effect and a proper assessment of this effect is crucial for a



**Figure A1.** Window corrected theory  $C_\ell$  for two different models with and without Redshift Space Distortions respectively in orange and blue. The dashed curves show the theoretical models before window convolution. The effect of the window is around 5% in redshift and 20% in real space. The theory with redshift space distortions uses  $galaxy\ bias = 2$  and the surface density  $n(z)$  of NGC eBOSS ELG (Tab. 4 of Raichoor et al. 2017) and assuming the fiducial cosmology of Ross et al. (2012); Ho et al. (2012).

robust measurement of the large-scale clustering of galaxies. We follow a similar approach to that of (Szapudi et al. 2001; Chon et al. 2004) to model the window function effect on the theoretical power spectrum  $C_\ell$  rather than correcting the measured pseudo-power spectrum from data. First, we compute the two-point correlation function of the window,

$$RR(\theta) = \sum_{i,j>1} f_{pix,i} f_{pix,j} \Theta_{ij}(\theta), \quad (\text{A2})$$

where  $\Theta_{ij}(\theta)$  is one when the pixels  $i$  and  $j$  are separated by an angle between  $\theta$  and  $\theta + \Delta\theta$ , and zero otherwise. Next, we normalize the  $RR$  by  $\sin(\theta)\Delta\theta$  to account for the area and total number of pairs such that  $RR(\theta = 0) = 1$ . We fit a polynomial on  $RR$  to smooth out the wiggles raised by noise. Then, we multiply the theoretical correlation function  $\omega(\theta)$  by the window paircount,

$$\omega^{WC}(\theta) = \omega(\theta) RR(\theta), \quad (\text{A3})$$

and finally use the Gaussian-Quadrature algorithm to transform the window convolved theory correlation function  $\omega^{WC}$  to  $C_\ell^{WC}$ ,

$$C_\ell^{WC} = 2\pi \int \omega^{WC}(\theta) P_\ell(\theta) d\theta. \quad (\text{A4})$$

Fig. A1 shows the DECaLS window effect on two theoretical models of  $C_\ell$  with and without redshift space distortions. The window effect for the model without Redshift Space Distortions (RSD) is around 20% but that for the model with RSD is less than 5% due to the flat power spectrum at the low  $\ell$  limit. We find a consistent pattern in the mocks; the window effect on the clustering of the mocks is between 5-15%.

This paper has been typeset from a  $\text{T}_\text{E}\text{X}/\text{L}^\text{A}\text{T}_\text{E}\text{X}$  file prepared by the author.