

# Photometric redshifts for the SDSS Data Release 12

Róbert Beck<sup>1,2,\*</sup>, László Dobos<sup>1</sup>, Tamás Budavári<sup>3,4,2</sup>,  
Alexander S. Szalay<sup>2</sup> and István Csabai<sup>1</sup>

<sup>1</sup>*Department of Physics of Complex Systems, Eötvös Loránd University, 1117 Budapest, Hungary*

<sup>2</sup>*Department of Physics and Astronomy, The Johns Hopkins University, Baltimore, MD 21218, USA*

<sup>3</sup>*Department of Applied Mathematics and Statistics, The Johns Hopkins University, Baltimore, MD 21218, USA*

<sup>4</sup>*Department of Computer Science, The Johns Hopkins University, Baltimore, MD 21218, USA*

Accepted 2016 Month Day. Received 2016 Month Day; in original form 2016 Month Day

## ABSTRACT

We present the methodology and data behind the photometric redshift database of the Sloan Digital Sky Survey Data Release 12 (SDSS DR12). We adopt a hybrid technique, empirically estimating the redshift via local regression on a spectroscopic training set, then fitting a spectrum template to obtain K-corrections and absolute magnitudes. The SDSS spectroscopic catalog was augmented with data from other, publicly available spectroscopic surveys to mitigate target selection effects. The training set is comprised of 1,976,978 galaxies, and extends up to redshift  $z \approx 0.8$ , with a useful coverage of up to  $z \approx 0.6$ . We provide photometric redshifts and realistic error estimates for the 208,474,076 galaxies of the SDSS primary photometric catalog. We achieve an average bias of  $\Delta z_{\text{norm}} = 5.84 \times 10^{-5}$ , a standard deviation of  $\sigma(\Delta z_{\text{norm}}) = 0.0205$ , and a  $3\sigma$  outlier rate of  $P_o = 4.11\%$  when cross-validating on our training set. The published redshift error estimates and photometric error classes enable the selection of galaxies with high quality photometric redshifts. We also provide a supplementary error map that allows additional, sophisticated filtering of the data.

**Key words:** galaxies: emission lines – galaxies: stellar content – galaxies: starburst – galaxies: active – methods: data analysis.

## 1 INTRODUCTION

Photometric redshift estimation has become a vital technique in the field of astronomy, as it enables measuring the distance of a much larger number of objects than what would be achievable through a spectroscopic survey. The Sloan Digital Sky Survey is one of the largest public collections of both photometric and spectroscopic measurements, with 208,478,448 galaxies in its photometric catalog (York et al. 2000; Gunn et al. 1998; Doi et al. 2010), and, as of Data Release 12 (Alam et al. 2015), 2,274,081 galaxy spectra in the continually expanded spectroscopic catalog (Eisenstein et al. 2011; Smee et al. 2013).

The purpose of this paper is to give a detailed description of the methods and data we used in creating the photometric redshift database of SDSS DR12, released to the public in January 2015. We chose an empirical technique, local linear regression, to estimate the redshift and its error, utilising a training set of 1,976,978 elements, assembled

from DR12 spectroscopy and data from other spectroscopic surveys (listed in Sec. 3.1). Additionally, we computed the maximum likelihood spectral template fit to the photometry, using the composite spectrum atlas of Dobos et al. (2012), to obtain additional information such as K-corrections, spectral type, and rest-frame absolute magnitudes.

The main goal of our photometric redshift catalog is to complement the estimated redshift with a reasonable assessment of the estimation error for the wide variety of galaxies in the SDSS photometric survey. The inclusion of spectroscopic data from other surveys means that we have more reference points for distant and faint bluer objects, up to  $z \approx 0.8$  and  $r \approx 21.5$  mag, which would be less well-represented in SDSS spectroscopy due to target selection (Eisenstein et al. 2001; Dawson et al. 2013). We also published an error map in support of this goal, as it highlights problematic regions in the space of galaxy colours where there are overlapping galaxies at different redshifts, leading to reduced accuracy.

The structure of the paper is as follows. In Sec. 2, we describe our empirical method of redshift estimation, and outline the template fitting procedure. Sec. 3 explains how

\* E-mail: beckrob23@caesar.elte.hu, dobos@complex.elte.hu, csabai@complex.elte.hu

the training set was compiled. Our results are presented and discussed in Sec. 4. We give pointers on using the database in Sec. 5. We provide a summary in Sec. 6.

Throughout the paper, broad-band magnitudes are quoted in the SDSS `asinh` magnitude system (Lupton, Gunn & Szalay 1999), and are dereddened according to Schlegel, Finkbeiner & Davis (1998). Following the recommendations of Scranton et al. (2005), for galaxy magnitudes we use the SDSS `cModelMag` magnitudes, and scale magnitude errors according to Eq. 15 in Scranton et al. (2005), while for galaxy colours we use SDSS `modelMag` magnitudes. Similarly to other SDSS applications, we adopt WMAP 5-year + SNe + BAO best-fitting cosmological parameters:  $\Omega_{\Lambda} = 0.726$ ,  $\Omega_m = 0.2739$ ,  $\Omega_r = 0.0001$  and  $H_0 = 70.5 \text{ km/s/Mpc}$  (Hinshaw et al. 2009).

The photometric database can be accessed via *Sky-Server*<sup>1</sup>. More information on the data used for this study, and program source code are available on the web site of the paper<sup>2</sup>. Colour versions of the figures are available in the online version of the paper.

### 1.1 Photometric redshift estimation

In the literature, there are two main approaches to estimating redshifts from broad-band photometry: the empirical, and the template-based approach.

Empirical methods generally utilise a supervised machine learning algorithm to find patterns in a training set – with both broad-band magnitude and redshift values – that allow prediction for cases when the redshift is not known. The ‘similarity’ of galaxies is usually defined in a metric space, the dimensions of which are some combination of the broad-band magnitudes and colours, perhaps with some scaling applied – we will refer to this as the colour and magnitude space, or, even more concisely, the colour space. The metric is generally chosen to be Euclidean distance within the colour space. Galaxies with a small distance between them – i.e. local galaxies – are considered to be similar, therefore their redshifts are also assumed to be similar. This assumption is then used in an algorithm for estimating galaxies with unknown redshifts. Examples of machine learning tools used for this purpose include artificial neural networks (Collister et al. 2007; Reis et al. 2012; Brescia et al. 2014), local polynomial fits (Csabai et al. 2007), random forests (Carliles et al. 2010), and boosted decision trees (Gerdes et al. 2010).

The template-based approach generally starts with a set of spectral templates and filter transmission curves, computes synthetic photometric magnitudes from them at various redshifts, and records the redshifts of templates that best reproduce the observed photometry. The choice of spectral templates is a crucial element of these methods. Galaxy templates can be computed theoretically using stellar models, an assumed initial mass function and stellar evolutionary tracks, which is a process known as stellar population synthesis (Fioc & Rocca-Volmerange 1997; Bruzual & Charlot 2003; Maraston & Strömbäck 2011; Vazdekis et al. 2012). The modelling of emission lines in such models – which

can contribute significantly to broad-band magnitudes (Atek et al. 2011) – is a difficulty because of the number of extra parameters needed to model the interstellar medium, but additional theoretical assumptions (Stasińska 1984; Fioc & Rocca-Volmerange 1997; Ferland et al. 2013) or empirical line estimation (Gyóry et al. 2011; Beck et al. 2016) can still be used. Alternatively, sets of measured galaxy spectra can be used to compile a library of empirical spectral templates, where the inclusion of measured lines is relatively straightforward (Yip et al. 2004; Dobos et al. 2012; Marchetti et al. 2013).

Template-based photometric redshift estimation methods in the literature include simple  $\chi^2$ -minimisation with a well-calibrated and wide set of templates (Arnouts et al. 2002; Ilbert et al. 2006), full Bayesian analysis using an empirical prior but relatively fewer templates (Benítez 2000; Coe et al. 2006), and Bayesian analysis using a linear combination of templates (Brammer, van Dokkum & Coppi 2008). Additional refinements include template corrections based on objects with known redshifts (Budavári et al. 2000; Csabai et al. 2000; Feldmann et al. 2006), and wavelength-dependent weighting of template errors (Brammer, van Dokkum & Coppi 2008).

The template-based approach has notable advantages over the empirical one: a training set with known redshifts is not required, and additional physical properties are implicitly estimated, since the entire template spectral energy distribution (SED) is known. However, unknown systematics in the photometric measurements are not accounted for, as opposed to empirical methods, where these are contained within the training set. Additionally, empirical techniques generally perform considerably better than template-based ones within the object type and redshift coverage of the training set (Csabai et al. 2003). However, the extrapolating capabilities of empirical methods are typically poor.

As in previous releases, to utilise the extensive spectroscopic sample of the SDSS, we elected to use an empirical method for estimating the redshift and its error, local linear regression. To get the best of both worlds, we combined this with a template fitting step that uses the photometric redshift, yielding additional physical information. We detail our methods in Sec. 2.

### 1.2 Difficulties in photometric redshift estimation

There are two main factors that are detrimental to the accuracy of photometric redshift estimation, regardless of the specific approach taken: the overlap in photometric colour space between different galaxy types, and the measurement errors in the photometry. While these are of different origin, their effect is very much intertwined.

The first factor, overlap in broad-band colour space, is a purely physical phenomenon. When the available colours cannot differentiate between morphological types, i.e. when different galaxy types have the same colours at different redshifts, there simply is not enough data to give an unequivocal answer to the question of what the redshift is. In such cases, the assumption that the broad-band magnitudes and colours uniquely determine the redshift does not hold, there are degeneracies in the colour–redshift relation (Benítez 2000).

The second factor, photometric measurement errors, is a major issue. The measurement errors can greatly exacerbate

<sup>1</sup> <http://skyserver.sdss.org/CasJobs/>

<sup>2</sup> <http://www.vo.elte.hu/papers/2016/photoz/>

the effects of overlap, blurring the divisions in colour space between different galaxy types, and also between galaxies of the same type but with differing redshifts (Benítez 2000). Additionally, when the measurement errors are not estimated accurately, or when photometric errors in different bands are correlated (Scranton et al. 2005), the assumption of uncorrelated Gaussian errors, used in many methodologies, simply does not hold (Budavári 2009).

These issues can be mitigated by improvements in the instrumentation. A better camera and telescope can reduce photometric errors (Ivezic et al. 2008; Tonry et al. 2012), while a large selection of filters (Wolf et al. 2003), or filters designed specifically for photometric redshift estimation (Budavári et al. 2001) can remove degeneracies.

In Sec. 4.2, we discuss how these factors affect our results.

## 2 METHODS

### 2.1 Local linear regression

Following Csabai et al. (2007) and earlier SDSS releases, we adopted a local (or piecewise) linear model to describe how the redshifts of galaxies depend on broad-band colours and magnitudes. The locality allows the model to follow the complex relationship between these properties, while using a polynomial of just the first order means that a relatively small number of galaxies is enough to fit the parameters. Taking just a few neighbouring galaxies into account helps preserve the local aspect of the model, and, as opposed to a simple average of the neighbours, the linear fit can follow subtle colour-dependent trends in the redshift.

Let  $i$  be the index of a galaxy in the set  $Q$  of galaxies to be estimated (query set), and let us denote the redshift of the  $i$ -th galaxy with  $z_i$  and its coordinates in the  $D$ -dimensional colour and magnitude space with the vector  $\mathbf{d}_i$ . Let us use  $j$  to index galaxies in the training set  $T$ , which is a collection of galaxies with both coordinate and redshift measurements –  $\mathbf{d}_j$  and  $z_j$ , respectively. Thus, our local linear model can be formulated in the following way:

$$z_i \approx c_i + \mathbf{a}_i \mathbf{d}_i = z_{\text{phot},i} \quad (1)$$

$z_{\text{phot},i}$  denotes the photometric redshift estimate. The parameter  $c_i$  is a constant offset, while components of the vector  $\mathbf{a}_i$  are linear coefficients. These parameters describe our model in the local neighbourhood of galaxy  $i$  – to determine them, we need to extract the local empirical relationship present in the training set,  $T$ . We do this by first finding the  $k$ -nearest neighbours of galaxy  $i$  within  $T$ , i.e. the  $k$  galaxies whose  $\mathbf{d}_j$  coordinates are the closest to  $\mathbf{d}_i$  in terms of Euclidean distance. Let us denote the set of nearest neighbours by  $NN$ . The parameters can then be determined using standard linear regression, by minimising the expression

$$\chi_i^2 = \sum_{j \in NN} \frac{(z_j - c_i - \mathbf{a}_i \mathbf{d}_j)^2}{w_j} \quad (2)$$

where  $w_j$  is a weight that could e.g. represent uncertainties in  $z_j$  and  $\mathbf{d}_j$ , or it could be a function of the distance

between  $\mathbf{d}_i$  and  $\mathbf{d}_j$ . The summation runs over the nearest neighbours, and the  $\chi_i^2$ -minimisation has to be done for every galaxy  $i$  within  $Q$ . The error of the photometric redshift  $z_{\text{phot},i}$  can be estimated by how well the thus fitted hyperplane reproduces the  $z_j$  redshifts of the nearest neighbours – we compute the RMS of the deviations from the fit:

$$\delta z_{\text{phot},i} \approx \sqrt{\frac{\sum_{j \in NN} (z_j - c_i - \mathbf{a}_i \mathbf{d}_j)^2}{k}} \quad (3)$$

In our current implementation, we have  $D = 5$  dimensions, and the components of the vectors  $\mathbf{d}_i$  and  $\mathbf{d}_j$  are the  $r$ -band magnitude, and the  $u-g$ ,  $g-r$ ,  $r-i$ ,  $i-z$  colours. All five dimensions are scaled to have zero mean and unit standard deviation. The nearest neighbours are weighted equally,  $w_j = 1$  for every  $j$ . These choices were made to optimise the accuracy of the photo- $z$  estimation. We use  $k = 100$  to have enough data points to determine the parameters and the error, but still preserve the locality of the model. The exact choice of  $k$  does not significantly impact the results, however.

We assume that the error of the spectroscopic redshift is negligible, i.e.  $z_j = z_{\text{spec},j}$ . Generally, this is a reasonable approximation because spectroscopic redshifts are much more accurate than photometric redshift estimates. However, it is important to note that there is a non-negligible percentage of spectroscopic redshift failures corresponding to a given quality cut in a survey (see Sec. 3.1 for a discussion of failure rates). If the failures are correlated with spectral type and colour, this systematic error in  $z_{\text{spec}}$  will be included in our training set, and thus propagate through to our  $z_{\text{phot}}$  estimates. Still, our best reference points for estimation are the redshifts published by spectroscopic surveys.

As an additional refinement of our method, when there are neighbours with outlying redshifts, we perform the computations twice to eliminate them. Neighbours that satisfy  $3\delta z_{\text{phot},i} < |z_j - c_i - \mathbf{a}_i \mathbf{d}_j|$  are discarded from the set  $NN$ , and the fit is redone for the limited set of  $l < k$  nearest neighbours, as needed. Also, we flag galaxies that lie outside the bounding box of the nearest neighbours in the  $D$ -dimensional colour and magnitude space. In such cases, we perform an extrapolation using the fitted hyperplane as opposed to an interpolation, therefore we can expect less reliable results (see Sec. 4.2 for more details).

Once the photometric redshift of the query point has been determined using this empirical method, we follow up with a spectral template fitting step, as described in the following section.

### 2.2 Spectral template fit

Let us denote redshift with  $z$ , galaxy type with  $t$ , measured magnitudes and magnitude errors with  $m$  and  $\Delta m$ , and synthetic magnitudes with  $s$ . Let us index the  $D$ -dimensional magnitude space with  $p$ , and, again, index galaxies in the query set with  $i$ . Under this notation, traditional maximum likelihood template-based photometric redshift estimation methods (Bolzonella, Miralles & Pelló 2000; Csabai et al. 2000; Arnouts et al. 2002; Ilbert et al. 2006) solve the following problem:

$$(z_{\text{phot},i}; t_i; m_{0,i}) = \arg \min_{(z;t;m_0)} \sum_{p=1}^D \left( \frac{m_{p,i} - (s_p(z, t) - m_0)}{\Delta m_{p,i}} \right)^2 \quad (4)$$

Here the constant offset in magnitude,  $m_0$ , is a scaling factor for the total flux with respect to the synthetic total flux. Generally,  $z$  and  $t$  iterate over a pre-determined list of redshifts and galaxy templates, while the best-fitting  $m_0$  can be calculated analytically for a given  $z$  and  $t$ .

In our hybrid approach, instead of iterating over  $z$ , we use the empirically determined photometric redshift, as described in Sec. 2.1. This way, we enjoy the benefit of higher redshift accuracy due to the extensive training set, while also fitting a galaxy template with a known SED. Thus, the expression we solve becomes:

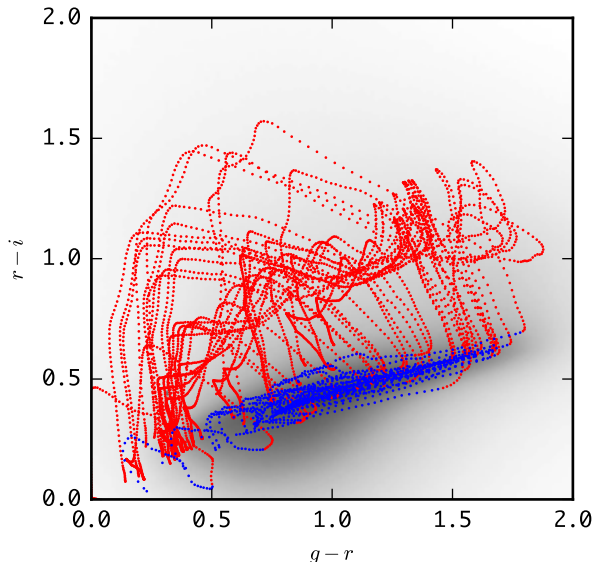
$$(t_i; m_{0,i}) = \arg \min_{(t;m_0)} \sum_{p=1}^D \left( \frac{m_{p,i} - (s_p(z = z_i, t) - m_0)}{\Delta m_{p,i}} \right)^2 \quad (5)$$

where  $z_i$  is computed using Eq. 1. As for the list of templates, we use the composite spectrum atlas of Dobos et al. (2012), which has been assembled from SDSS spectra, takes emission lines into account, and contains extreme red and blue galaxy types in addition to the more frequently occurring ones. Dobos et al. (2012) also published synthetic photometric magnitudes in the SDSS filter set for a grid of redshift values. Fig. 1 shows the coverage of the templates in  $g-r$ ,  $r-i$  colours – the dense galaxy regions are well-covered by the composite spectrum atlas. For the set of all photometric galaxies, the fitted synthetic magnitudes are within  $3\Delta m$  of the measured  $m$  for 82.0%, 89.7%, 96.2%, 97.2% and 97.3% of cases, respectively, for the  $u$ ,  $g$ ,  $r$ ,  $i$  and  $z$  broad-band magnitudes. The normalized error distributions are roughly Gaussian, with the exception of the  $u$ -band, where it is asymmetric. Considering the redshift estimation errors and outlier rate in the unfiltered galaxy set (see Sec. 4 and Tab. 3 for more details), the within- $3\Delta m$  ratios are relatively high, which shows that the templates adequately describe the fitted galaxies.

Once we have found the best-fitting spectral template  $t_i$ , we determine other values of physical interest. The  $DM$  distance modulus and  $D_L$  luminosity distance are computed using the redshift and our assumed cosmology. Knowing the SED of the template, we also calculate observed-frame synthetic magnitudes,  $K$ -corrections to redshifts 0 and 0.1, and rest-frame absolute magnitudes (see Sec. A1 for exact definitions of these).

### 3 TRAINING SET

Our training set initially consisted of the entire spectroscopic galaxy catalog of the Sloan Digital Sky Survey Data Release 12. This includes the earlier main galaxy and LRG samples, and also the more recent BOSS sample. The main galaxy sample consists of a wide variety of galaxies, with no cuts on colour, although it is rather limited in terms of redshift (Strauss et al. 2002). The LRG sample provides an expanded redshift coverage, however, it has specifically targeted luminous red galaxies (Eisenstein et al. 2001). The BOSS sample



**Figure 1.** The colour space coverage of the spectral templates from Dobos et al. (2012) in  $g-r$ ,  $r-i$  dimensions. Blue dots show templates with  $z \leq 0.35$ , while red dots correspond to redshifts  $0.35 < z < 0.7$ . The template colours are superimposed on a grayscale density map of SDSS photometric measurements.

extends much deeper than the former two, and has somewhat relaxed the sharp colour cuts of the LRG sample, but it is still targeted towards massive galaxies (Dawson et al. 2013), likely resulting in non-negligible selection effects.

On the other hand, the photometric galaxy catalog of the SDSS has no such selection effects, and our goal is to provide photometric redshifts for the entire catalog, not just a subset of morphological types or colour. Since it would be advantageous to have a wider selection of galaxy types and colours even at higher redshifts, we decided to extend our training set by cross-matching galaxies in the Sloan photometric catalog with spectroscopic measurements of other, publicly available surveys.

#### 3.1 Data from other surveys

The spectroscopic surveys with which we extended our training set are listed in Tab. 1, with references. For each survey, we used the published redshift quality flag to select only the reasonably confident redshift measurements, with confidences  $\geq 95\%$  (with the exception of PRIMUS, where  $\geq 92\%$ ).

We cross-matched the galaxies from other surveys with SDSS primary photometric galaxy measurements by using J2000 right ascension and declination coordinates and published astrometric errors. We followed the probabilistic methodology of Budavári & Szalay (2008), assumed Gaussian errors, and calculated the Bayes factor of Eq. 16 in Budavári & Szalay (2008), which is the ratio of the likelihood that the two measurements are of the same source, and the likelihood that they are of separate sources:

$$B = \frac{L(\text{same source})}{L(\text{separate sources})} = \frac{2}{\sigma_1^2 + \sigma_2^2} \exp \left\{ -\frac{\psi^2}{2(\sigma_1^2 + \sigma_2^2)} \right\} \quad (6)$$

Survey Name	References	Quality flag
2dF	Colless et al. (2001, 2003)	$Quality = 4, 5$
6dF	Jones et al. (2004, 2009)	$Q = 3, 4$
DEEP2	Davis et al. (2003); Newman et al. (2013)	$ZQUALITY = 3, 4$
GAMA	Driver et al. (2011); Baldry et al. (2014)	$NQ = 4$
PRIMUS	Coil et al. (2011); Cool et al. (2013)	$Q = 4$
VIPERS	Garilli et al. (2014); Guzzo et al. (2014)	$[zflg] \bmod 10 = 3, 4$
VVDS	Le Fèvre et al. (2004); Garilli et al. (2008)	$ZFLAGS \bmod 10 = 3, 4$
WiggleZ	Drinkwater et al. (2010); Parkinson et al. (2012)	$Qop = 4, 5$
zCOSMOS	Lilly et al. (2007, 2009)	$[Class] \bmod 10 = 3, 4$

**Table 1.** Information about the external spectroscopic surveys we used to expand our training set.

Here  $\sigma_1$  and  $\sigma_2$  are the astrometric errors of two given galaxies, and  $\psi$  is the angular separation between them. We accepted matches with  $B > 10,000$ , thus ensuring that we only used rather certain matches.

Galaxies with existing SDSS spectrometry were excluded from the cross-match, and where we found multiple matches for the same Sloan galaxy, we selected the one with the smallest redshift error.

In total, we found 168,834 matches with reasonable redshift confidence. However, later filtering steps greatly limited the number we could utilise, to 76,193.

Multiple matches provide an opportunity to test whether the published spectroscopic redshift failure rates are correct and whether the cross-match itself works reliably. Of the total of 1,012 multiple matches in the filtered sample, 171 were between two PRIMUS measurements, 769 had one PRIMUS object, and 72 did not include PRIMUS. (We are handling PRIMUS separately because of its lower confidence level and higher redshift error compared to other surveys.)

The only-PRIMUS set had a standard deviation of  $\sigma(\Delta z_{\text{spec}}) = 0.00473$  and a  $3\sigma$  outlier rate of  $P_o = 9.36\%$  (outliers were removed iteratively). Individual PRIMUS measurements typically have an accuracy of  $\sim 0.005$ , therefore the deviation is even below what one would expect, and the outlier rate is also well below the theoretically expected  $1 - 0.92 \times 0.92 = 15.4\%$ .

The PRIMUS–other survey set is described by the numbers  $\sigma(\Delta z_{\text{spec}}) = 0.00472$  and  $P_o = 8.97\%$ . The deviation is roughly the accuracy of PRIMUS, just as expected, while the outlier rate is again below the expected  $1 - 0.95 \times 0.92 = 12.6\%$ .

The non-PRIMUS set had a standard deviation of  $\sigma(\Delta z_{\text{spec}}) = 0.00071$  and an outlier rate of  $P_o = 29.2\%$ . The deviation is negligible for our purposes, but the outlier rate is significantly larger than the expected  $1 - 0.95 \times 0.95 = 9.75\%$ . The observed discrepancies might be due to a number of reasons, below we list a few.

- The spectroscopic redshift failures could be correlated, which would reduce the combined outlier rate. Especially the only-PRIMUS set could be affected, where the same survey measured the same object twice.

- Galaxies could be erroneously cross-matched due to an underestimation of astrometric accuracy – DEEP2, the survey responsible for 71.4% of outliers in the non-PRIMUS set, uses the Canada France Hawaii Telescope, which quotes the USNOA 2.0 astrometric error of  $0.5''$  (Coil et al. 2004), the highest value of all the external surveys.
- Overlapping galaxies in the field of view could compro-

mise spectroscopic redshifts, and could also lead to incorrect cross-matches.

Furthermore, it is important to note that the non-PRIMUS set only had 72 matches, of which 21 were outliers. This is a rather small sample size, and might not be representative. On the whole, the confidence levels derived from multiple matches are in line with – or better than – the expectations based on the published numbers of the surveys.

### 3.2 Filtering the training set

While our goal was to assemble a training set with as wide a coverage in redshift and colour space as possible, the inclusion of objects with too large photometric errors would diminish our ability to find the most similar reference galaxies. The ‘true’ nearest neighbours may be scattered away due to errors, with less similar galaxies taking their place. To alleviate this problem, we introduced photometric error cuts to the training set. Additionally, we filtered out galaxies with outlying colours, which both eliminates erroneous measurements, and also more clearly defines the boundaries of our training set in colour space.

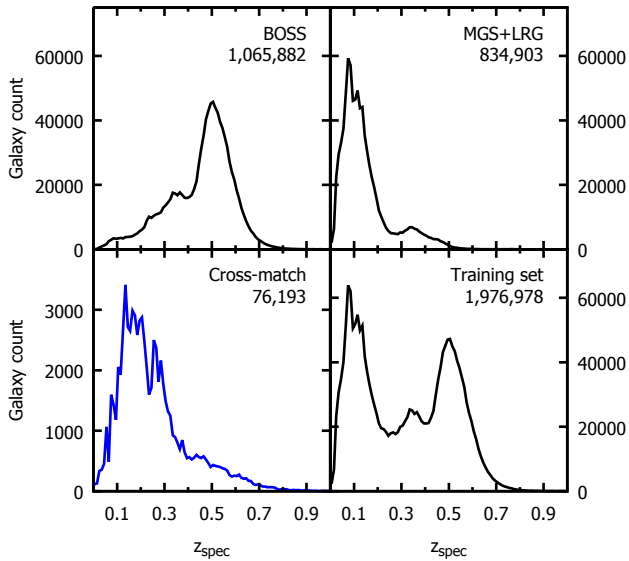
The exact parameters of the cuts were determined empirically, with the following criteria in mind:

- optimise the photometric redshift estimation results,
- leave no region empty in the space of broad-band colours, if otherwise within the coverage of the training set,
- keep fainter and higher-redshift measurements of sufficient accuracy.

The final values of the photometric error and colour cuts are as follows:

$$\begin{aligned}
 \Delta r &< 0.15 \\
 \Delta(g - r) &< 0.225 \\
 \Delta(r - i) &< 0.15 \\
 \Delta(i - z) &< 0.25 \\
 -0.911 &< (u - g) < 5.597 \\
 0.167 &< (g - r) < 2.483 \\
 0.029 &< (r - i) < 1.369 \\
 -0.452 &< (i - z) < 0.790
 \end{aligned} \tag{7}$$

Magnitudes are in the SDSS *ugriz* filter system, with errors scaled following Scranton et al. (2005) (see also Sec. 1). The colour cuts correspond to filtering out the highest and lowest 0.5% of data for the  $(u - g)$  colour, and 1%



**Figure 2.** The redshift distribution of our entire training set, and subsets of it: the BOSS spectroscopic sample, the pre-BOSS spectroscopic sample that includes the main galaxy sample and the LRG sample (MGS+LRG), and the additional galaxies cross-matched from other surveys. In the top right corner of each panel, we indicate the corresponding subset, and the total number of galaxies within that subset. Note the different scale of the cross-match subset.

for the other three colours. The reason for having no limit on  $\Delta u$ , and for having relaxed criteria on  $(u - g)$  compared to other colours is that the errors of the SDSS  $u$ -band are generally much larger than that of other bands, and even galaxies with fairly secure photometric redshifts can have very large  $u$ -band errors.

Only those galaxies were included in the training set that fulfilled all of Eq. 7. Additionally, SDSS galaxies with unsecure spectroscopic redshifts were also cut: the spectroscopic error flag `SpeczWarning` had to either take the value `OK` or `MANY_OUTLIERS` (the latter rarely signifying a real error according to the documentation). This spectroscopic error flag cut filters out a higher and higher fraction of galaxies as the redshift increases – with more distant galaxies typically having lower signal-to-noise spectra – but there is no indication of a specific redshift being preferentially eliminated, which otherwise could have pointed to a systematic incompleteness in our training set. The redshift distribution of the finalised training set is shown in Fig. 2.

## 4 RESULTS

### 4.1 Accuracy of photometric redshifts

To evaluate the performance of our methods, we randomly divided the training set into two equal-sized subsets, and performed cross-validation, estimating the photometric redshifts of one half using the other half as the training set (and vice versa). The resulting photometric redshifts could then be contrasted with the spectroscopic redshifts. Fig. 3 shows the photometric redshift  $z_{\text{phot}}$ , the estimation error  $z_{\text{phot}} - z_{\text{spec}}$ , and the estimation error divided by the re-

ported photometric redshift error  $\delta z_{\text{phot}}$ , as functions of the spectroscopic redshift.

Using the normalized redshift estimation error  $\Delta z_{\text{norm}} = \frac{z_{\text{phot}} - z_{\text{spec}}}{1 + z_{\text{spec}}}$ , we achieve an average bias of  $\overline{\Delta z_{\text{norm}}} = 5.84 \times 10^{-5}$ , a standard deviation of  $\sigma(\Delta z_{\text{norm}}) = 0.0205$ , and an outlier rate of  $P_o = 4.11\%$ . Outliers are defined as  $|\Delta z_{\text{norm}}| > 3\sigma(\Delta z_{\text{norm}})$ , and are removed iteratively. While most galaxies in the training set have fairly small estimation errors, on Fig. 3 it is apparent that there are redshift ranges where there is a non-negligible bias, up to  $\Delta z = 0.01$  or  $0.5 \delta z_{\text{phot}}$ .

Still, in biased regions between 58% and 76% of galaxies are within  $\pm 1 \delta z_{\text{phot}}$ , and between 86% and 98% of galaxies are within  $\pm 2 \delta z_{\text{phot}}$ . Thus, the confidence interval  $z_{\text{phot}} \pm \delta z_{\text{phot}}$  can be reasonably used in applications, as it will contain a fairly high fraction of galaxies even when there is bias in the estimation, and the distribution is not centered on  $z_{\text{phot}}$ .

Additionally, on Fig. 4, we plotted the probability density function of  $(z_{\text{phot}} - z_{\text{spec}})/\delta z_{\text{phot}}$  alongside a standard normal distribution. There is a small overall bias, but otherwise the two distributions match rather well, highlighting that our method for estimating the error of the photometric redshift gives a fair assessment of the estimation accuracy.

Another issue visible on Fig. 3 is that the estimation accuracy declines dramatically from around  $z = 0.6$ , where the number count of the training set falls off. These high-redshift galaxies occupy sparsely sampled regions in colour space, as evidenced by the fact that 94% of them are above the 50th, and 68% are above the 75th percentile of nearest neighbour bounding box volume. Sparse regions are more likely to include a non-negligible amount of galaxies scattered there due to high photometric errors. In the case of high-redshift galaxy regions, the scattered galaxies are also more likely to have lower redshifts, hence the negative estimation bias.

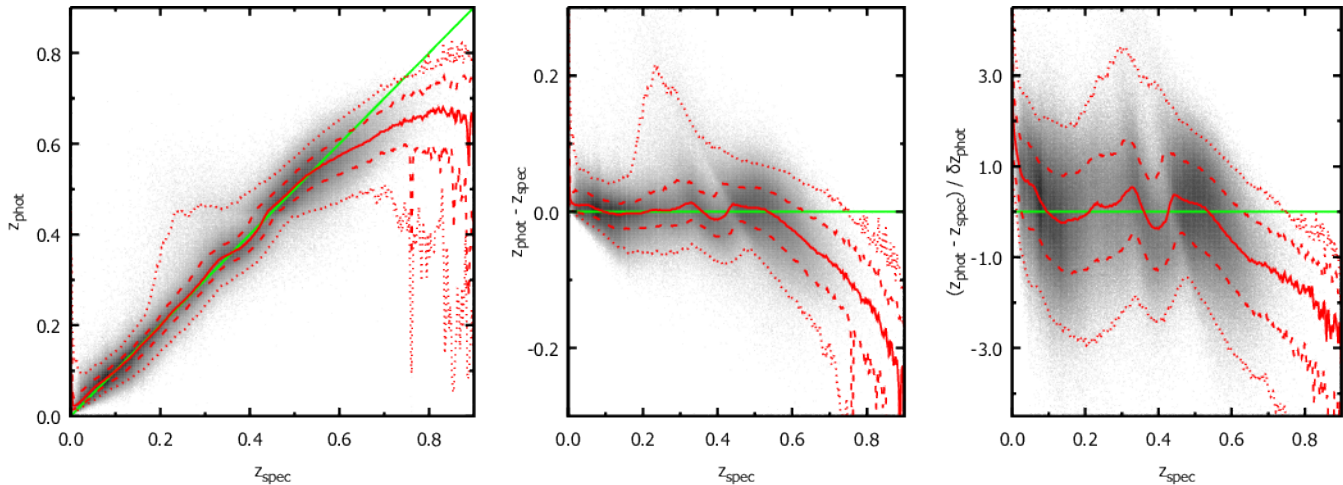
In the following section, we will go into more detail concerning the biases and errors.

### 4.2 Discussion of biases and errors

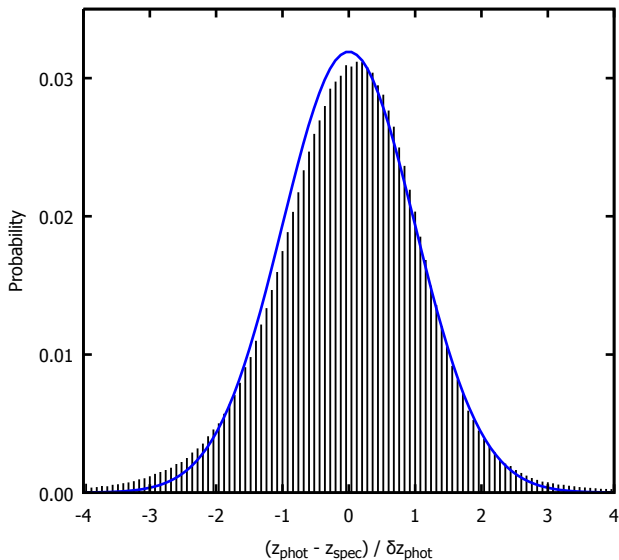
Here we discuss how the issues outlined in Sec. 1.2, namely overlap and errors in the photometry, affect our results.

When galaxies of different redshifts overlap in colour space, the nearest neighbours that we find with our algorithm will have a bimodal, or even multimodal redshift distribution. In this case, the estimated redshift will lie between the different peaks in the distribution, and the estimated error will increase accordingly. Additionally, galaxies belonging to a peak at a smaller redshift will have a positive estimation bias, while galaxies that correspond to a peak at a higher  $z$  will be estimated with a negative bias. When the overlap is 'perfect', it is impossible to decide which is the appropriate galaxy group, but when the loci of groups in colour space have a slight offset between them, the closer galaxy group will more strongly constrain the fitted hyperplane locally. Assuming a mixture of two Gaussian distributions of equal weight but with different means, our method will estimate the redshift closer to the correct peak, as opposed to a simple  $k$ -nearest neighbour average, which would give the centerpoint, the average of the means. This effect becomes





**Figure 3.** The photometric redshift ( $z_{\text{phot}}$ ) as a function of spectroscopic redshift ( $z_{\text{spec}}$ ),  $z_{\text{phot}} - z_{\text{spec}}$  as a function of  $z_{\text{spec}}$ , and  $(z_{\text{phot}} - z_{\text{spec}})/\delta z_{\text{phot}}$  as a function of  $z_{\text{spec}}$ , where  $\delta z_{\text{phot}}$  is the photometric redshift error estimate. The galaxy density of our training set is shown in grayscale – we took the logarithm of galaxy counts so that even individual galaxies can be seen. The red solid, dashed and dotted lines represent the median, 68% and 95% confidence regions of the training set, respectively. The green line shows  $z_{\text{spec}} = z_{\text{phot}}$ , i.e. what would be the perfect estimation. See the text for a discussion.



**Figure 4.** The normalized histogram of the scaled redshift error,  $(z_{\text{phot}} - z_{\text{spec}})/\delta z_{\text{phot}}$ , is plotted in black for our training set. The blue line shows the standard normal distribution. See the text for a discussion.

less noticeable when one of the groups is underrepresented in the training set, or when photometric errors are large enough to sufficiently mix the groups in colour space. To remove some of the degeneracy, in addition to the colours, we also used the  $r$ -band magnitude in our local linear regression. In Sec. 4.3, we describe an error map that helps quantify the effects of overlap.

Photometric measurement errors strongly limit the accuracy of photometric redshifts in the SDSS catalog – both the training set and the query set are affected, especially the fainter galaxies. We introduced several photometric error classes for the galaxies to quantify the dependence of redshift estimation errors on errors in the photometry: class

Class	$\Delta r_{\text{max}}$	$\Delta(g-r)_{\text{max}}$	$\Delta(r-i)_{\text{max}}$	$\Delta(i-z)_{\text{max}}$
1	0.15	0.225	0.15	0.25
2	0.18	0.25	0.18	0.28
3	0.21	0.30	0.21	0.31
4	0.24	0.35	0.24	0.34
5	0.27	0.40	0.27	0.37
6	0.30	0.45	0.30	0.40

**Table 2.** The maximum photometric error values that a galaxy belonging to a photometric error class was allowed to have. The errors were scaled following Scranton et al. (2005). Each galaxy is placed in the lowest possible class. Class 7 contains galaxies that could not be placed in any other class.

1 matches the error limits of the training set, and the subsequent classes (2 – 7) contain galaxies with progressively higher errors. The exact limits were determined empirically based on the photometric error and redshift error distributions, with the aim of giving a sequence from useful photometric redshifts to highly inaccurate ones. The class identifiers have a negative sign if the local linear regression was an extrapolation, i.e. the estimated galaxy lay outside the bounding box of the nearest neighbours. For example, class  $-1$  denotes galaxies that match the error limits of the training set, but were estimated with an extrapolation. This way, class  $-1$  also includes galaxies that did not satisfy the colour cuts of Eq. 7, and therefore are not within the training set. Classes 2 – 7 and  $(-2) - (-7)$  do contain galaxies with spectroscopic redshifts, specifically those that did not fulfill the error limits of Eq. 7 – these galaxies can be used to test the redshift estimation accuracy in a given class.

Tab. 2 gives the photometric error limits used for each class, while Tab. 3 lists the redshift estimation bias, standard deviation, outlier rate, and the spectroscopic and photometric galaxy count in the classes. It is clear that higher photometric errors correspond to sharply increasing biases and deviations.

We emphasise here that, since the training set only contains galaxies of class 1 or  $-1$ , the redshift error estimate

Class	$\overline{\Delta z_{\text{norm}}}$	$\sigma(\Delta z_{\text{norm}})$	$P_o$	$N_{\text{spec}}$	$N_{\text{phot}}$
1	$6.11 \times 10^{-5}$	0.0204	4.07%	1,957,234	42,410,836
2	-0.0033	0.0333	4.03%	77,281	5,657,368
3	-0.0057	0.0331	3.93%	68,610	5,240,766
4	-0.0082	0.0369	4.45%	36,218	3,955,814
5	-0.0107	0.0412	5.00%	19,110	2,970,417
6	-0.0127	0.0486	5.33%	10,674	2,232,881
7	-0.0222	0.0823	3.80%	16,563	6,950,249
-1	$4.22 \times 10^{-4}$	0.0289	5.71%	19,744	2,001,544
-2	-0.0051	0.0549	11.2%	5,940	1,421,618
-3	-0.0081	0.0514	8.97%	10,262	2,848,424
-4	-0.0104	0.0567	7.65%	11,200	4,098,896
-5	-0.0150	0.0643	6.68%	10,917	5,118,595
-6	-0.0165	0.0728	6.06%	10,350	5,862,776
-7	-0.0488	0.1410	2.30%	86,574	117,703,892

**Table 3.** The average redshift estimation bias  $\overline{\Delta z_{\text{norm}}}$ , standard deviation  $\sigma(\Delta z_{\text{norm}})$ , outlier rate  $P_o$ , number of spectroscopic galaxies  $N_{\text{spec}}$  and number of photometric galaxies  $N_{\text{phot}}$  in each photometric error class, with  $\Delta z_{\text{norm}} = \frac{z_{\text{phot}} - z_{\text{spec}}}{1 + z_{\text{spec}}}$ . Outliers are defined to have  $|\Delta z_{\text{norm}}| > 3\sigma(\Delta z_{\text{norm}})$ , and are removed iteratively. We indicate extrapolation in the local linear regression with a negative sign in front of the class identifier.

( $\delta z_{\text{phot}}$ ) is expected to be an accurate representation of the estimation error only when the query galaxy also belongs to class 1 or -1 (and satisfies Eq. 7, if class -1). As we show in Sec. 4.3, the redshift estimation errors are dependent on the position in colour space. A higher photometric error class leads to additional variance in  $z_{\text{phot}}$ , which therefore should ideally be characterized as a function of the position in colour space. However, as shown in Tab. 3, there are relatively few spectroscopic galaxies in the higher error classes, and we do not have a good enough coverage to allow a detailed treatment of this phenomenon. As a crude first approximation for other classes, the class-wide extra variance with respect to class 1 can be added according to Tab. 3.

### 4.3 The redshift error map

As described in the previous section, the presence of biases and higher errors in the redshift estimation is strongly dependent on the position of a given galaxy in the space of broad-band colours. To provide a tool for filtering out regions in the colour space where these issues are the most prominent, we compiled and published an error map.

The error map gives the redshift estimation results – as computed on the training set – for a 3D grid in  $r$ -band magnitude, and  $g-r$ ,  $r-i$  colours. For each bin in the grid, we report the galaxy count, the average  $z_{\text{spec}}$ , the average  $z_{\text{phot}}$ , the rms of  $z_{\text{phot}} - z_{\text{spec}}$ , the average  $\delta z_{\text{phot}}$ , and the average standard deviation of the redshifts of the neighbours,  $\sigma(z_{\text{NN}})$ . With the help of this map, it is possible to flag galaxies in sparsely populated regions, or in regions with high estimation errors (which also indicate possible biases).

To illustrate, we computed these measures for a 2D projection of the 3D map, where the  $r$ -band magnitude has been summed over, and the grid remains in  $g-r$ ,  $r-i$  colours. In Fig. 5, the galaxy count distribution is shown as a function of the two colours. The pronounced discontinuity – a diagonal line – is a target selection effect produced by the colour cut of the CMASS subsample in the BOSS survey (Dawson et al. 2013), leading to a sparsely populated region below the cut. In Fig. 6, we plotted three measures of the redshift error, all of which show similar behaviour. The esti-

mation error is highest where the redshifts of the neighbours have a larger deviation, i.e. where there is overlap of galaxies with differing redshifts. Since the photometric error limit of even error class 1 is as high as  $\Delta(g-r)_{\text{max}} = 0.225$  and  $\Delta(r-i)_{\text{max}} = 0.15$  in these two dimensions, such mixing between different redshifts is to be expected. The reported error follows the actual error closely, drawing the same overall picture of the dependence of estimation errors on the location in colour space, which also supports our result in Sec. 4.1 that our redshift error estimates are accurate. Additionally, sparsely populated regions in Fig. 5 correspond to higher errors in Fig. 6. Since sparse regions could be occupied either by exotic galaxy types or galaxies that were scattered there due to high photometric errors, it is not surprising that their redshift estimation is inaccurate.

## 5 USING THE DATABASE

The photometric redshift database has been made public along with the SDSS DR12. It can be accessed via *Sky-Server*<sup>3</sup>, it is the Photoz table within the DR12 context. The redshift error map is contained in the table PhotozErrorMap, also in the DR12 context. Refer to Sec. A for a description of each column in these tables.

### 5.1 Best practices

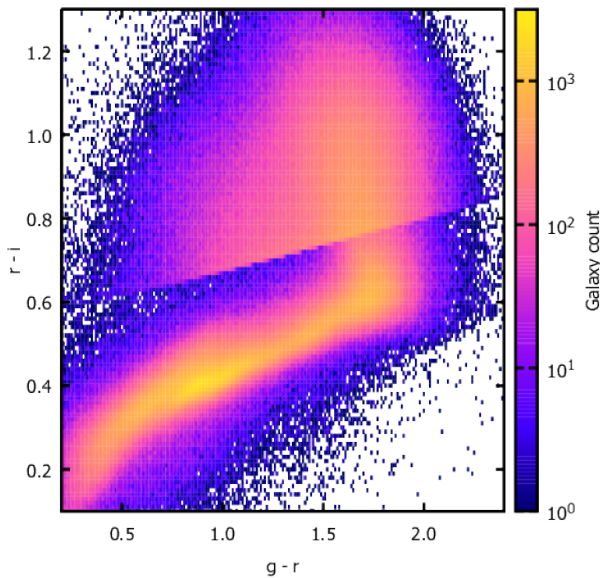
Here we intend to give a few pointers on how best to utilise our database.

As a first step, we recommend only using galaxies with a photometric error class of 1, because that is when the redshift error estimate is expected to be accurate. If more galaxies are desired, follow the instructions at the end of Sec. 4.2 and use Tab. 3 for including additional error classes.

When nearest neighbours in the local linear regression are outliers, they are excluded from the fit. However, having

<sup>3</sup> <http://skyserver.sdss.org/CasJobs/>





**Figure 5.** The galaxy count distribution of the training set, on a 2D grid of  $g-r$ ,  $r-i$  broad-band colours. The published 3D map also includes the  $r$ -band magnitude, this 2D version is used here for illustrative purposes. The discontinuity is caused by the colour cut of the CMASS sample (Dawson et al. 2013). See the text for a discussion.

too many outliers may indicate that the given galaxy is difficult to estimate, therefore a limit should be put on the minimum number of nearest neighbours used out of the total of  $k = 100$ , e.g. a minimum of  $l = 97$ . Additionally, the desired accuracy may be achieved with a cut based on the redshift error estimate, e.g. only using galaxies with  $\delta z_{\text{phot}} < 0.03$ .

The linear fitting algorithm can fail when it encounters a singular or near-singular matrix – such cases are indicated by  $z_{\text{phot}} = -9999$ , therefore those galaxies should be excluded. If required, the redshift of the first nearest neighbour, and the average redshift of the 100 nearest neighbours are still available, however, in this case, there is no redshift error estimate (also flagged with  $\delta z_{\text{phot}} = -9999$ ).

When small biases are a critical issue, the redshift error map of Sec. 4.3 should be used for leaving out galaxies that are located in a high-error region in colour space, but that otherwise have a low reported  $\delta z_{\text{phot}}$ .

Additionally, the volume of the bounding box of the nearest neighbours in the colour space could also be used for filtering – a volume that is very large means that galaxies of very different colours are used in the local linear regression, therefore the estimated redshift could be compromised. To give an idea of potential limits, a bounding box volume cut of  $\text{nnVol} < 2$  filters out  $\approx 2.5\%$  of training set galaxies that are in the sparsest colour space regions,  $\text{nnVol} < 1$  eliminates  $\approx 5\%$ , while  $\text{nnVol} < 0.45$  cuts  $\approx 10\%$ . The galaxies thus filtered out have estimation accuracies of  $\sigma(\Delta z_{\text{norm}}) = 0.0464$ , 0.0400 and 0.0342, respectively, with the  $3\sigma$  outlier rate around 7% for all three cases.

In Fig. 7, we illustrate the hazards of using non-filtered photometric redshift data. Without implementing any cuts, the errors and the number of catastrophic failures are visibly much larger. Also, it is important to remember that we are only analysing spectroscopic measurements, which on

average have much more accurate photometry than the rest of the SDSS photometric catalog – the fraction of catastrophic outliers is expected to be much larger for the unfiltered photometric sample. On the other hand, using too stringent cuts might unnecessarily limit the redshift coverage, or the colour space coverage of the sample. For this reason, we recommend experimenting with different filtering choices to find the one most appropriate for the task at hand.

## 6 SUMMARY

We described in detail how we created the photometric redshift database of SDSS DR12.

After a brief overview of the photometric redshift estimation literature, we defined the local linear regression method that we use for the redshift and redshift error estimation, and described the spectral template fitting step that followed it. We gave an account of the data and methods that went into assembling the training set. We evaluated the accuracy of our estimation via cross-validation on the training set, then we discussed the errors and biases that we encountered. We introduced photometric error classes, and a 3D redshift error map to help quantify the errors and filter out inaccurately estimated galaxies. We also provided recommendations for using the database, and choosing appropriate filtering criteria.

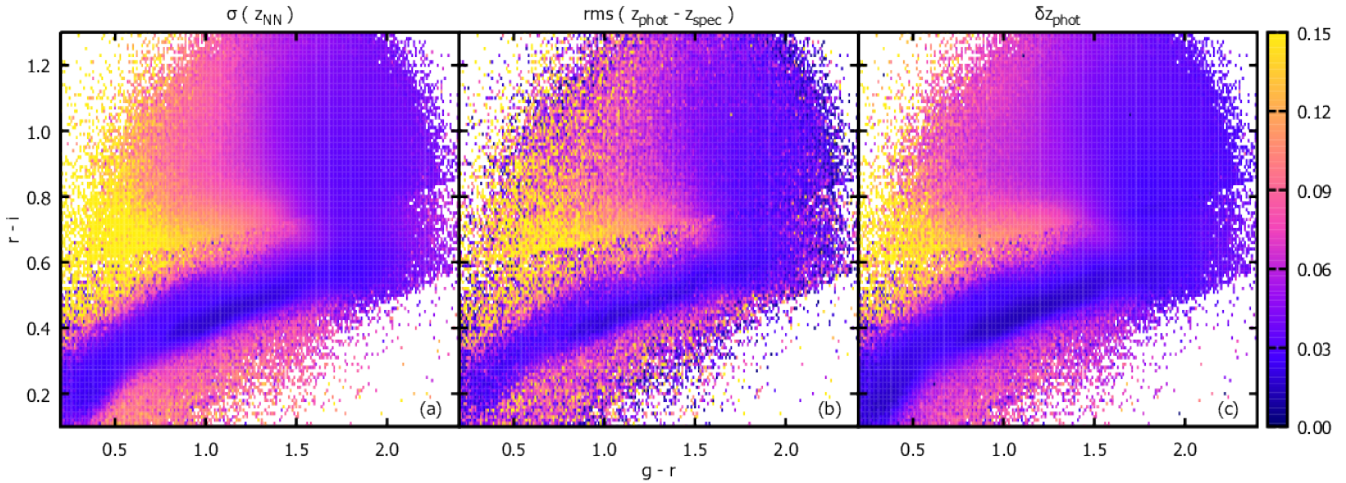
Our photometric redshift estimates are relatively accurate, with a standard deviation of  $\sigma(\Delta z_{\text{norm}}) = 0.0205$ , and an acceptable  $3\sigma$  outlier rate of  $P_o = 4.11\%$ . The reported redshift error is a realistic estimate of the actual redshift estimation error (see Fig. 4). While we observed redshift-dependent biases of up to  $\Delta z = 0.01$ , the  $z_{\text{phot}} \pm \delta z_{\text{phot}}$  confidence intervals provide a reasonably good approximation of the spectroscopic redshift (see Sec. 4.1). However, from  $z \approx 0.6$ , the coverage of our training set drops sharply, therefore so does the accuracy of our photometric redshifts.

In addition to the redshift error estimate, we provide further tools that allow users to select measurements of the desired accuracy. These include the photometric error class, the 3D redshift error map, and the bounding box volume of the nearest neighbours (see Sec. 4.2, Sec. 4.3 and Sec. 5.1, respectively).

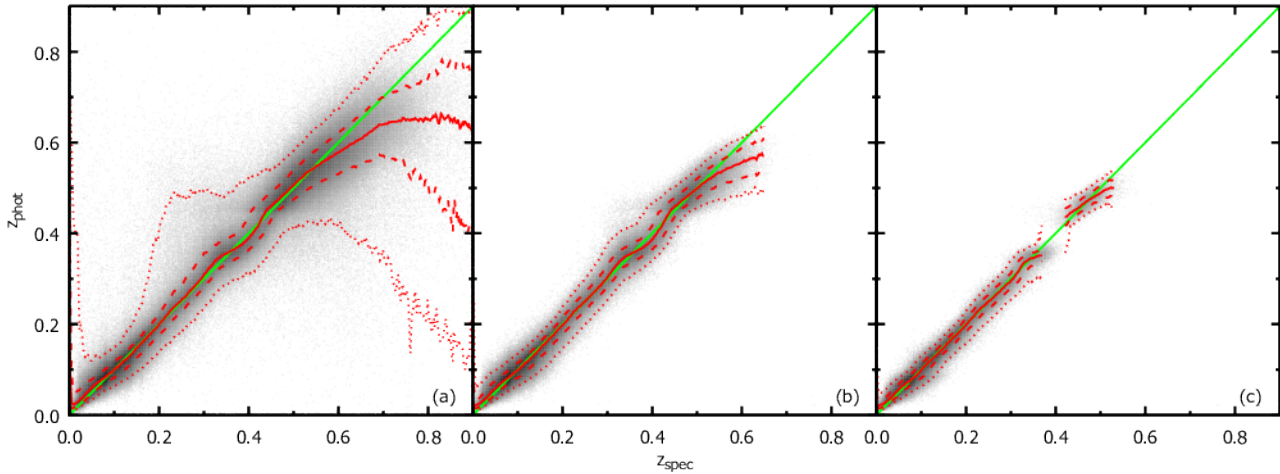
As opposed to purely empirical methods, our hybrid method fits a spectral template, which allows us to provide K-corrections and absolute magnitudes (see Sec. 2.2 and Sec. A1).

In later releases, we intend to expand our training set as new data comes available, and also review our training set and methods with the intention of reducing biases and extending the useful coverage to higher redshifts. Having a less sparse sampling of high-redshift galaxies in photometric colour space would help reduce the pronounced negative bias in their redshift estimation.

The photometric redshift database, corresponding documentation and tools are available online on the appropriate SDSS DR12 webpages.



**Figure 6.** Photometric redshift estimation results for the training set, on a 2D grid of  $g-r$ ,  $r-i$  broad-band colours. The published 3D map also includes the  $r$ -band magnitude, this 2D version is used here for illustrative purposes. Panel (a) shows the average standard deviation of the redshifts of the nearest neighbours ( $\sigma(z_{NN})$ ), panel (b) displays the rms of  $z_{\text{phot}} - z_{\text{spec}}$ , the actual estimation error, while panel (c) shows the average reported estimation error ( $\delta z_{\text{phot}}$ ). We note that the outliers have not been removed from the rms computation, therefore panel (b) is noisier. See the text for a discussion.



**Figure 7.** The photometric redshift ( $z_{\text{phot}}$ ) as a function of spectroscopic redshift ( $z_{\text{spec}}$ ), for three different subsets of all available spectroscopic measurements. The galaxy density is shown in grayscale – we took the logarithm of galaxy counts so that even individual galaxies can be seen. The red solid, dashed and dotted lines represent the median, 68% and 95% confidence regions of the data, respectively. On panel (a), there is no selection, all galaxies are shown. On panel (b), we only included galaxies of photometric error class 1 or -1, and with a reported redshift error of  $\delta z_{\text{phot}} < 0.03$ . On panel (c), galaxies of photometric error class 1 and with  $\delta z_{\text{phot}} < 0.02$  are shown. We note that on panel (c), the more biased region around  $z_{\text{spec}} = 0.4$  has been almost completely filtered out. See the text for a discussion.

## ACKNOWLEDGMENTS

The realisation of this work was supported by the Hungarian OTKA NN grants 103244 and 114560.

## REFERENCES

- Alam S. et al., 2015, *ApJS*, 219, 12  
 Arnouts S. et al., 2002, *MNRAS*, 329, 355  
 Atek H. et al., 2011, *ApJ*, 743, 121  
 Baldry I. K. et al., 2014, *MNRAS*, 441, 2440  
 Beck R., Dobos L., Yip C.-W., Szalay A. S., Csabai I., 2016, *MNRAS*, 457, 362  
 Benítez N., 2000, *ApJ*, 536, 571  
 Bolzonella M., Miralles J.-M., Pelló R., 2000, *A&A*, 363, 476  
 Brammer G. B., van Dokkum P. G., Coppi P., 2008, *ApJ*, 686, 1503  
 Brescia M., Cavaoti S., Longo G., De Stefano V., 2014, *A&A*, 568, A126  
 Bruzual G., Charlot S., 2003, *MNRAS*, 344, 1000  
 Budavári T., 2009, *ApJ*, 695, 747  
 Budavári T., Szalay A. S., 2008, *ApJ*, 679, 301  
 Budavári T., Szalay A. S., Connolly A. J., Csabai I., Dickinson M., 2000, *AJ*, 120, 1588  
 Budavári T., Szalay A. S., Csabai I., Connolly A. J., Tsvetanov Z., 2001, *AJ*, 121, 3266  
 Carliles S., Budavári T., Heinis S., Priebe C., Szalay A. S.,

2010, ApJ, 712, 511  
Coe D., Benítez N., Sánchez S. F., Jee M., Bouwens R., Ford H., 2006, AJ, 132, 926  
Coil A. L. et al., 2011, ApJ, 741, 8  
Coil A. L., Newman J. A., Kaiser N., Davis M., Ma C.-P., Kocevski D. D., Koo D. C., 2004, ApJ, 617, 765  
Colless M. et al., 2001, MNRAS, 328, 1039  
Colless M. et al., 2003, ArXiv Astrophysics e-prints  
Collister A. et al., 2007, MNRAS, 375, 68  
Cool R. J. et al., 2013, ApJ, 767, 118  
Csabai I. et al., 2003, AJ, 125, 580  
Csabai I., Connolly A. J., Szalay A. S., Budavári T., 2000, AJ, 119, 69  
Csabai I., Dobos L., Trencsényi M., Herczegh G., Józsa P., Purger N., Budavári T., Szalay A. S., 2007, Astronomische Nachrichten, 328, 852  
Davis M. et al., 2003, in Society of Photo-Optical Instrumentation Engineers (SPIE) Conference Series, Vol. 4834, Discoveries and Research Prospects from 6- to 10-Meter-Class Telescopes II, Guhathakurta P., ed., pp. 161–172  
Dawson K. S. et al., 2013, AJ, 145, 10  
Dobos L., Csabai I., Yip C.-W., Budavári T., Wild V., Szalay A. S., 2012, MNRAS, 420, 1217  
Doi M. et al., 2010, AJ, 139, 1628  
Drinkwater M. J. et al., 2010, MNRAS, 401, 1429  
Driver S. P. et al., 2011, MNRAS, 413, 971  
Eisenstein D. J. et al., 2001, AJ, 122, 2267  
Eisenstein D. J. et al., 2011, AJ, 142, 72  
Feldmann R. et al., 2006, MNRAS, 372, 565  
Ferland G. J. et al., 2013, Revista Mexicana de Astronomía y Astrofísica, 49, 137  
Fioc M., Rocca-Volmerange B., 1997, A&A, 326, 950  
Garilli B. et al., 2014, A&A, 562, A23  
Garilli B. et al., 2008, A&A, 486, 683  
Gerdes D. W., Sypniewski A. J., McKay T. A., Hao J., Weis M. R., Wechsler R. H., Busha M. T., 2010, ApJ, 715, 823  
Gunn J. E. et al., 1998, AJ, 116, 3040  
Guzzo L. et al., 2014, A&A, 566, A108  
Gyórfy Z., Szalay A. S., Budavári T., Csabai I., Charlot S., 2011, AJ, 141, 133  
Hinshaw G. et al., 2009, ApJS, 180, 225  
Ilbert O. et al., 2006, A&A, 457, 841  
Ivezic Z. et al., 2008, ArXiv e-prints  
Jones D. H. et al., 2009, MNRAS, 399, 683  
Jones D. H. et al., 2004, MNRAS, 355, 747  
Le Fèvre O. et al., 2004, A&A, 417, 839  
Lilly S. J. et al., 2009, ApJS, 184, 218  
Lilly S. J. et al., 2007, ApJS, 172, 70  
Lupton R. H., Gunn J. E., Szalay A. S., 1999, AJ, 118, 1406  
Maraston C., Strömbäck G., 2011, MNRAS, 418, 2785  
Marchetti A. et al., 2013, MNRAS, 428, 1424  
Newman J. A. et al., 2013, ApJS, 208, 5  
Parkinson D. et al., 2012, Phys. Rev. D, 86, 103518  
Reis R. R. R. et al., 2012, ApJ, 747, 59  
Schlegel D. J., Finkbeiner D. P., Davis M., 1998, ApJ, 500, 525  
Scranton R., Connolly A. J., Szalay A. S., Lupton R. H., Johnston D., Budavári T., Brinkman J., Fukugita M., 2005, ArXiv Astrophysics e-prints  
Smees S. A. et al., 2013, AJ, 146, 32

Stasińska G., 1984, A&AS, 55, 15  
Strauss M. A. et al., 2002, AJ, 124, 1810  
Tonry J. L. et al., 2012, ApJ, 750, 99  
Vazdekis A., Ricciardelli E., Cenarro A. J., Rivero-González J. G., Díaz-García L. A., Falcón-Barroso J., 2012, MNRAS, 424, 157  
Wolf C., Meisenheimer K., Rix H.-W., Borch A., Dye S., Kleinheinrich M., 2003, A&A, 401, 73  
Yip C. W. et al., 2004, AJ, 128, 585  
York D. G. et al., 2000, AJ, 120, 1579

## APPENDIX A: THE PHOTOMETRIC REDSHIFT TABLES IN SDSS DR12

Here we give a description of each column in the published tables, either referencing a concept used in the article, or detailing it here. With  $\{\text{ugriz}\}$  we denote that there is a column for each of the five SDSS  $\text{ugriz}$  broad-band magnitudes, with the single corresponding (capitalized) letter present in the column name.

### A1 The Photoz table

- **objID** – the SDSS objID of the query galaxy.
- **z** –  $z_{\text{phot},i}$  in Eq.1, i.e. the photometric redshift. It takes the value  $-9999$  when there was an error in the fitting algorithm.
- **zErr** –  $\delta z_{\text{phot},i}$  in Eq.3, i.e. the photometric redshift error estimate. It takes the value  $-9999$  when there was an error in the fitting algorithm.
- **nnCount** –  $l$ , the number of nearest neighbours used in the local linear regression, with outliers excluded from the total of  $k = 100$ , as described in Sec. 2.1. It takes the value  $-9999$  when there was an error in the fitting algorithm.
- **nnVol** – the volume of the bounding box of the  $k = 100$  nearest neighbours.
- **photoErrorClass** – the photometric error class described in Sec. 4.2, Tab. 2 and Tab. 3.
- **nnObjID** – the SDSS objID of the first nearest neighbour.
- **nnSpecz** – the spectroscopic redshift ( $z_{\text{spec}}$ ) of the first nearest neighbour.
- **nnFarObjID** – the SDSS objID of the farthest, 100th nearest neighbour.
- **nnAvgZ** – the average redshift of the  $k = 100$  nearest neighbours.
- **distMod** – the distance modulus ( $DM_i$ ) corresponding to **z**, if available, or **nnAvgZ**. See the end of Sec. 1 for the adopted cosmology.
- **lumDist** – the luminosity distance in  $Mpc$  corresponding to **z**, if available, or **nnAvgZ**. See the end of Sec. 1 for the adopted cosmology.
- **chisq** – the  $\chi^2$  value of the spectral template fit, i.e.  $\chi^2 = \sum_{p=1}^D \left( \frac{m_{p,i} - (s_p(z=z_i, t=t_i) - m_{0,i})}{\Delta m_{p,i}} \right)^2$ , using the notation of Eq. 5.
- **rnorm** – the residual Euclidean norm of the spectral template fit, i.e.  $\left( \sum_{p=1}^D (m_{p,i} - (s_p(z=z_i, t=t_i) - m_{0,i}))^2 \right)^{0.5}$ , using the notation of Eq. 5.

$t_i$	Name	$t_i$	Name	$t_i$	Name
1	Red P	15	Blue P	29	RED 1
2	Red H $\alpha$	16	Blue H $\alpha$	30	RED 2
3	Red SF	17	Blue SF	31	RED 3
4	Red A+HII	18	Blue A+HII	32	RED 4
5	Red L	19	Blue L	33	RED 5
6	Red S	20	Blue S	34	SF 1
7	Red all	21	Blue all	35	SF 2
8	Green P	22	All P	36	SF 3
9	Green H $\alpha$	23	All H $\alpha$	37	SF 4
10	Green SF	24	All SF	38	SF 5
11	Green A+HII	25	All A+HII		
12	Green L	26	All L		
13	Green S	27	All S		
14	Green all	28	All all		

**Table A1.** The name in Dobos et al. (2012) that corresponds to the  $t_i$  (or `bestFitTemplateID`) template identifier used in this article.

- `bestFitTemplateID` –  $t_i$ , the identifier of the best-fitting spectral template. See Tab. A1 for the corresponding names in Dobos et al. (2012).

- `synth{ugriz}` – the synthetic magnitude of the best-fitting spectral template, i.e.  $s_p(z = z_i, t = t_i) - m_{0,i}$ , using the notation of Eq. 5.

- `kcorr{ugriz}` – the  $K$ -correction to  $z = 0$ , i.e.  $K_{p,i}(z = 0) = s_p(z = z_i, t = t_i) - s_p(z = 0, t = t_i)$ , using the notation of Eq. 5.

- `kcorr{ugriz}01` – the  $K$ -correction to  $z = 0.1$ , i.e.  $K_{p,i}(z = 0.1) = s_p(z = z_i, t = t_i) - s_p(z = 0.1, t = t_i)$ , using the notation of Eq. 5.

- `absMag{ugriz}` – the rest-frame absolute magnitude of the galaxy, i.e.  $m_{p,i} - K_{p,i}(z = 0) - DM_i$ , using the notation of Eq. 5.

## A2 The PhotozErrorMap table

- `CellID` – The unique identifier of the cell in the grid. The grid spans the  $r$ -band magnitude, and the  $g - r$ ,  $r - i$  colours.

- `rMag` – The centerpoint of the cell in  $r$ -band magnitude. Linear size of a cell: 0.5.

- `gMag_Minus_rMag` – The centerpoint of the cell in  $g - r$  colour. Linear size of a cell: 0.01.

- `rMag_Minus_iMag` – The centerpoint of the cell in  $r - i$  colour. Linear size of a cell: 0.01.

- `countInCell` – The number of training set galaxies within the cell (denoted below with  $N$ ).

- `avgPhotoZ` – The average photometric redshift of training set galaxies in the cell, i.e.  $\frac{\sum_{i=1}^N z_{\text{phot},i}}{N}$ , using the notation of Sec. 2.1.

- `avgSpectroZ` – The average spectroscopic redshift of training set galaxies in the cell, i.e.  $\frac{\sum_{i=1}^N z_{\text{spec},i}}{N}$ , using the notation of Sec. 2.1.

- `avgRMS` – The rms of the redshift estimation for training set galaxies in the cell, i.e.  $\left(\frac{\sum_{i=1}^N (z_{\text{phot},i} - z_{\text{spec},i})^2}{N}\right)^{0.5}$ , using the notation of Sec. 2.1.

- `avgEstimatedError` – The average redshift error estimate for training set galaxies in the cell, i.e.  $\frac{\sum_{i=1}^N \delta z_{\text{phot},i}}{N}$ , using the notation of Sec. 2.1.

- `avgNeighborZStDev` – The average standard deviation of the redshifts of the  $k = 100$  nearest neighbours, for every training set galaxy in the cell. Denoting the standard deviation of the  $z_{\text{spec}}$  of the neighbours with  $\sigma_i(z_{NN})$ , it is  $\frac{\sum_{i=1}^N \sigma_i(z_{NN})}{N}$ .

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