

PHOTOMETRIC REDSHIFTS FOR QUASARS IN MULTI-BAND SURVEYS

M. BRESCIA^{1,2}, S. CAVUOTI², R. D'ABRUSCO³, G. LONGO^{2,4}, AND A. MERCURIO¹

¹ INAF-Astronomical Observatory of Capodimonte, via Moiariello 16, I-80131 Napoli, Italy; brescia@oacn.inaf.it

² Department of Physics, University Federico II, via Cinthia 6, I-80126 Napoli, Italy

³ Harvard-Smithsonian Center for Astrophysics, 60 Garden Street, Cambridge, MA 02138, USA

Received 2013 February 28; accepted 2013 May 23; published 2013 July 17

ABSTRACT

The Multi Layer Perceptron with Quasi Newton Algorithm (MLPQNA) is a machine learning method that can be used to cope with regression and classification problems on complex and massive data sets. In this paper, we give a formal description of the method and present the results of its application to the evaluation of photometric redshifts for quasars. The data set used for the experiment was obtained by merging four different surveys (Sloan Digital Sky Survey, *GALEX*, UKIDSS, and *WISE*), thus covering a wide range of wavelengths from the UV to the mid-infrared. The method is able (1) to achieve a very high accuracy, (2) to drastically reduce the number of outliers and catastrophic objects, and (3) to discriminate among parameters (or features) on the basis of their significance, so that the number of features used for training and analysis can be optimized in order to reduce both the computational demands and the effects of degeneracy. The best experiment, which makes use of a selected combination of parameters drawn from the four surveys, leads, in terms of Δz_{norm} (i.e., $(z_{\text{spec}} - z_{\text{phot}})/(1 + z_{\text{spec}})$), to an average of $\Delta z_{\text{norm}} = 0.004$, a standard deviation of $\sigma = 0.069$, and a median absolute deviation, $\text{MAD} = 0.02$, over the whole redshift range (i.e., $z_{\text{spec}} \leq 3.6$), defined by the four-survey cross-matched spectroscopic sample. The fraction of catastrophic outliers, i.e., of objects with photo- z deviating more than 2σ from the spectroscopic value, is $<3\%$, leading to $\sigma = 0.035$ after their removal, over the same redshift range. The method is made available to the community through the DAMEWARE Web application.

Key words: catalogs – galaxies: distances and redshifts – methods: data analysis – quasars: general – surveys

Online-only material: color figures

1. INTRODUCTION

Photometric redshifts (hereafter photo- z) provide an estimate of the redshift of sources obtained using photometry instead of spectroscopy. They are in fact driven (1) by the shape of the broadband continuum of the object's spectroscopic emission and (2) by a limited number of strong spectral features (i.e., the one at 4000 \AA , the Ly α forest, and the Lyman limit), which are still recognizable after the integration of the spectral energy distribution (SED) sampled by the filter's transmission function.

At the cost of lower accuracy, photo- z offer several advantages with respect to their spectroscopic counterparts: (1) being derived from intermediate/broadband imaging, photo- z are much more effective in terms of observing time; (2) they may allow one to probe objects much fainter than the spectroscopic flux limit; and (3) under specific conditions, they allow one to correct some biases, such as those encountered at high redshift where, as it has been noted (Fernandez-Soto et al. 2001), spectroscopy is pushed to its limits both by the low signal-to-noise ratio (S/N) in the spectra and by the fact that in many cases, even when a good S/N is achieved, the lack of features in the observed spectral range may undermine the estimation of a trustworthy redshift (Lanzetta et al. 1998).

The latter aspect becomes crucial when photometric redshift methods are applied to quasars (QSO) and, in particular, to the construction and characterization of the large, complete samples which are required by modern cosmology. In fact, quasar samples have always been, and still are, constructed either by compiling lists of more or less serendipitous discoveries

obtained with different techniques and selection criteria (Veron-Cetty & Veron 2000), or via a two-step process where the first one consists in the identification of QSO candidates from multi-wavelength surveys, and the second one requires the spectroscopic validation of the candidates. In practice, due to the large amount of observing time required by spectroscopy, the latter step is usually optimized by applying the spectroscopic validation procedure just to a more or less significant subsample of the candidates and then by extrapolating the resulting statistics to the whole sample. Modern surveys are usually so deep and extensive that the number of candidates rapidly becomes too large to be handled with the latter approach. On the other hand, modern multi-wavelength digital surveys also provide such a wealth of information (multi-band high-accuracy photometry) that it becomes feasible to approximate the SED of objects over a quite large range of redshifts (Richards et al. 2001a, 2001b; Budavari et al. 2001; Wolf et al. 2004), thus minimizing the need for spectroscopic follow-up.

In the last few years, it has in fact been demonstrated that after having provided an accurate enough photometry and significant wavelength coverage, it is possible to obtain samples of photometrically selected quasars matching the low contamination and high completeness (D'Abrusco et al. 2009; Bovy et al. 2012) required by many fields of modern cosmology. The relevance of these photometric samples will increase more and more in the near future, when the new generation of deeper and more accurate surveys will allow us to access larger and more complete samples of QSOs. These *photometric* samples are in fact already being used for a variety of applications such as the measurement of the integrated Sachs–Wolfe effect (Giannantonio et al. 2008), the cosmic magnification bias (Scranton et al. 2005), and the clustering of quasars on large (Myers et al. 2006) and

⁴ Visiting Associate, California Institute of Technology, Pasadena, CA 91125, USA.

small (Hennawi et al. 2006) scales, to note just a few. Since both candidate selection and photometry redshift estimates are performed on the same data (colors in many bands), it is also apparent that for the same samples, photometric data alone should carry enough information to characterize in an almost univocal way the SED, and therefore also to derive accurate estimates of photometric redshifts (D’Abrusco et al. 2009; Laurino et al. 2011; Bovy et al. 2012).

It goes without saying that the utility of the photometric samples goes hand in hand with the development of photo- z methods capable of providing accurate enough estimates of the redshifts.

In this paper, we use a new empirical method, named Multi Layer Perceptron with Quasi Newton Algorithm (MLPQNA), and apply it to the evaluation of photometric redshift of quasars. In Section 2, we discuss the data sets used for the experiments, and in Section 3 we present both a detailed description of the MLPQNA method and the statistical indicators used throughout the paper. We wish to stress that the lack of a common agreement on such indicators is among the main obstacles in comparing the performances of different methods. In Section 4, we describe the experiments performed in order to select the best combination of input parameters, bands, and network topology. The results of these experiments are summarized and discussed in Section 5, where we also present the final performances of the best experiments. Finally, we compare our results with those available in the literature and draw some general conclusions.

A short Appendix provides the reader with the math behind the Quasi Newton Algorithm (QNA).

2. THE DATA SET

The sample of quasars used in the experiments described in this paper is based on the spectroscopically selected quasars from the Sloan Digital Sky Survey (SDSS) DR7 database (table *Star* of the SDSS database). According to the spectroscopic classification index (*index SP* or *specClass*) provided in the SDSS-DR7 release (Schneider et al. 2010), we selected quasars for which a reliable measure of the spectroscopic redshifts (with $zConf > 0.90$) is available.

We then cross-matched the SDSS sample quasars identified as point sources with clean measured photometry in all filters (*ugriz*), with the latest versions of the data sets from *Galaxy Evolution Explorer* (GALEX; Martin et al. 2005), UKIDSS (Lawrence et al. 2007), and *WISE* (Wright et al. 2010). These three surveys observed large fractions of the sky in the ultraviolet, near-infrared, and middle-infrared spectral intervals, respectively. After the cross-matching, we obtained a series of multi-band catalogs, defined as follows.

SDSS (DR7) (Aihara et al. 2011) has observed $\sim 1.4 \times 10^4$ deg² of the sky in five bands (*ugriz*) covering the [3551, 8931] Å wavelength range. Photometric SDSS observations reach the limiting magnitude of 22.2 in the *r* band (95% completeness for point sources; Abazajian et al. 2009).

GALEX (DR6/7) (Martin et al. 2005) is a two-band survey (*nuv, fuv* for near and far ultraviolet, respectively) covering the [1300,3000] Å wavelength interval. The GALEX photometric survey has observed the whole sky to the near ultraviolet limiting magnitude $nuv = 20.5$.

UKIDSS (DR9) (Lawrence et al. 2007) has been designed to be the SDSS infrared counterpart and covers ~ 7000 deg² of the sky in the *YJHK* near-infrared bands covering the ~ 0.9 – 2.4 μ m spectral range down to the limiting magnitude $K = 18.3$. The

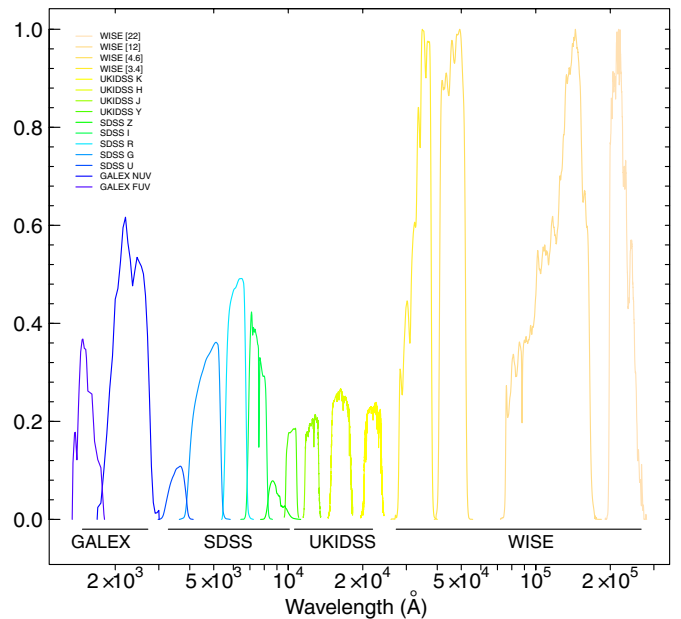


Figure 1. Transmission curves for all filters in the four surveys considered. (A color version of this figure is available in the online journal.)

Large Area Survey has imaged ~ 4000 deg² (overlapping with the SDSS), with the additional *Y* band down to the limiting magnitude of 20.5.

The *WISE* mission (Wright et al. 2010) has observed the entire sky in the mid-infrared spectral interval at 3.4, 4.6, 12, and 22 μ m with an angular resolution of 6'.1, 6'.4, 6'.5, and 12'.0 in the four bands, achieving 5σ point source sensitivities of 0.08, 0.11, 1, and 6 mJy in unconfused regions on the ecliptic, respectively. The astrometric accuracy of *WISE* is $\sim 0''.50$, $0''.26$, $0''.26$, and $1''.4$ for the four *WISE* bands, respectively.

The transmission curves of all filters related to the four surveys are shown in Figure 1. All these surveys present a large common overlap region and overall good astrometry with comparable astrometric accuracy. In order to cross-match the catalogs, we used a maximum radius $r = 1''.5$ to associate the optical quasars with counterparts in each of the three catalogs. Afterward, we rejected all sources containing one or more missing data in any of their photometric parameters. In this case, the term *missing data* means undefined numerical values underlying either not detected or contaminated magnitude measurements. This last step is crucial in empirical methods since the presence of missing data might affect their generalization capabilities (Marlin 2008).

The resulting number of objects in the data sets used for the experiments are the following:

1. SDSS: $\sim 1.1 \times 10^5$;
2. SDSS \cap GALEX: $\sim 4.5 \times 10^4$;
3. SDSS \cap UKIDSS: $\sim 3.1 \times 10^4$;
4. SDSS \cap GALEX \cap UKIDSS: $\sim 1.5 \times 10^4$;
5. SDSS \cap GALEX \cap UKIDSS \cap WISE: $\sim 1.4 \times 10^4$.

An additional data set was produced by decimating the final *four-survey* cross-matched catalog. This data set was used to perform the preliminary feature selection or *pruning* phase (see Section 4.2) and consisted of $\sim 3.8 \times 10^3$ objects, each observed in 15 bands (4 UKIDSS, 2 GALEX, 5 SDSS, and 4 WISE) and with accurate spectroscopic redshift estimates. The decimation was needed to reduce the computational time needed to perform

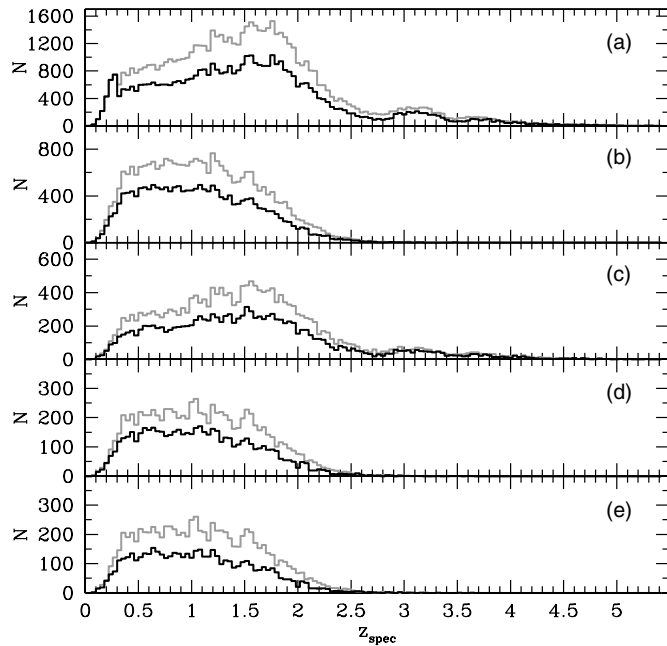


Figure 2. Histograms of the spectroscopic redshift distribution in the five survey cross-matched samples as derived from the SDSS spectroscopic data. (a) SDSS, (b) SDSS+GALEX, (c) SDSS+UKIDSS, (d) SDSS+GALEX+UKIDSS, and (e) SDSS+GALEX+UKIDSS+WISE. The gray dotted line is the training sample. The black line is the test sample.

the large number of experiments described in what follows. For some bands there were multiple measurements (i.e., magnitude measured accordingly to different definitions), and therefore we are left with a total of 43 different features.

Finally, when producing training and test sets, we made sure that they had compatible spectroscopic redshift distributions (see Figure 2).

3. THE METHOD

This section is dedicated to the description of the machine learning method used for the experiments. All mathematical details are reported in the [Appendix](#).

3.1. Multi Layer Perceptron

From a technical point of view, the MLPQNA method is a multi layer perceptron (MLP; Bishop 2006) neural network (NN) trained by a learning rule based on the QNA; in other words and as it is synthesized in the acronym, MLPQNA differs from more traditional MLPs' implementations in the way the optimal solution of the regression problem is found. In previous papers, most of the characteristics of the method have been described in the contexts of both classification (Brescia et al. 2012) and regression (Cavuoti et al. 2012b).

According to Bishop (2006), feed-forward NNs (in their various implementations) provide a general framework for representing nonlinear functional mappings between a set of input variables (also called features) and a set of output variables (the targets). The training of an NN can be in fact seen as the search for the function which minimizes the errors of the predicted values with respect to the true values available for a small but significant subsample of objects in the same parameter space (PS). This subset is also called *training set* or *knowledge base* (KB). The final performances of a specific NN depend on many factors, such as topology, the way the minimum of

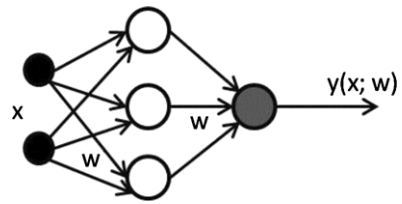


Figure 3. Scheme of a Multi Layer Perceptron general architecture for two input variables, one hidden layer with three neurons and one output value.

the error function is searched and found, the way errors are computed, as well as the intrinsic quality of training data.

The formal description of a feed-forward NN with two computational layers is given in Equation (1):

$$y_k = \sum_{j=0}^M w_{kj}^{(2)} g \left(\sum_{i=0}^d w_{ji}^{(1)} x_i \right). \quad (1)$$

Equation (1) can be better understood by using a graph like the one shown in Figure 3. The input layer (x_i) is made of a number of neurons (also known as perceptrons) equal to the number of input variables (d); the output layer, on the other hand, will have as many neurons as the output variables (k).

In the general case, the network may have an arbitrary number of hidden layers (in the depicted case there is just one hidden layer with three neurons), each of which can be formed by an arbitrary number of neurons (M). In a fully connected feed-forward network, each node of a layer is connected to all the nodes in the adjacent layers. Each connection is represented by an adaptive weight ($w_{kj}^{(l)}$) that can be regarded as the strength of the synaptic connection between neurons k and j , while the response of each perceptron to the inputs is represented by a nonlinear function g , referred to as the *activation function*. All the above characteristics, the topology of the network and the weight matrix of its connections, define a specific implementation and are usually called a *model*.

The model, however, is only part of the story. In fact, in order to find the model that best fits the data in a specific problem, one has to provide the network with a set of examples, id est, of objects for which the final output is known by independent means. These data, already defined as *training set* or *KB*, are used by the network to find the optimal model.

In our implementation, we choose as learning rule the QNA, which differs from the Newton Algorithm in how the Hessian of the error function is computed. Newtonian models are variable metric methods used to find local maxima and minima of functions (Davidon 1968) and, in the case of MLPs, they can be used to find the stationary (i.e., the zero gradient) point of the learning function. A complete mathematical description of the MLP with the QNA model is reported in the [Appendix](#).

The model has been made available to the community through the DATA Mining & Exploration Web Application RESOURCE (DAMEWARE⁵; Cavuoti et al. 2012a).

3.2. The Implementation of MLPQNA

In this work, we use our implementation of the QNA based on the limited-memory BFGS (L-BFGS; Byrd et al. 1994), where BFGS is the acronym composed of the names of the four inventors (Broyden 1970; Fletcher 1970; Goldfarb 1970; Shanno 1970).

⁵ http://dame.dsf.unina.it/beta_info.html

To summarize, the algorithm for MLP with QNA is the following. Let us consider a generic MLP with $w^{(t)}$ being the weight vector at time (t).

1. Initialize all weights $w^{(0)}$ with small random values (typically normalized in $[-1, 1]$), and set the constant error tolerance ε and $t = 0$;
2. present to the network the whole training set and calculate $E(w^{(t)})$ as the error function for the current weight configuration;
3. if $t = 0$, then $d^{(t)} = -\nabla E^{(t)}$;
4. otherwise $d^{(t)} = -\nabla E^{(t-1)} + Ap + Bv$, where $p = w^{(t+1)} - w^{(t)}$ and $v = g^{(t+1)} - g^{(t)}$;
5. calculate $w^{(t+1)} = w^{(t)} - \alpha d^{(t)}$, where α is obtained by a line search equation (see Equation (A6) in the appendix);
6. calculate A and B for the next iteration, as reported in Equation (A19);
7. if $E(w^{(t+1)}) > \varepsilon$, then $t = t + 1$ and go to (2), otherwise STOP.

As is known, all *line search* methods, which are based on techniques that search for the minimum error by exploring the error function surface, are likely to get stuck in a local minimum and many solutions to this problem have been proposed (Floudas & Jongen 2005). In order to optimize the convergence of the gradient descent analysis (GDA; see the Appendix), Newton's method uses the information on the second-order derivatives. By having information on the second derivatives, QNA is more effective in avoiding local minima of the error function and more accurate in the error function trend follow-up, thus revealing a *natural* capability to find the absolute minimum error of the optimization problem (Shanno 1990).

In the L-BFGS version of the algorithm, in the case of high dimensionality (i.e., input data with many parameters), the amount of memory required to store the Hessian is too large, along with the machine time required to process it. Therefore, instead of using a complete number of gradient values to generate the Hessian, we can use a smaller number of values to approximate it.

By the way, if the convergence slows down, the performance may even increase. This statement might seem paradoxical but, although the convergence is measured by the number of iterations, the performance depends on the number of the processor's time units spent calculating the result.

Related to the computational cost, there is also the strategy adopted in terms of stopping criteria for the method. As is known, the process of adjusting the weights based on the gradients is repeated until a minimum is reached. In practice, one has to decide the stopping condition of the algorithm. Among the possible criteria, the algorithm could be terminated when (1) the Hessian approximation error becomes sufficiently small (by definition, the gradient will be zero at a minimum), (2) the maximum number of iterations is reached, in terms of a fixed threshold, and (3) is based on the cross validation.

The cross validation can be used to monitor generalization performance during training and to terminate the algorithm when there is no more improvement. Statistically significant results come out by trying multiple independent data partitions and then averaging the performances. There are several variants of cross validation methods (Sylvain & Celisse 2010). In particular, we have chosen the k -fold cross validation, which is particularly suited in the presence of a scarcity of known data samples (Geisser 1975). The mechanism, also known as *leave-one-out*, is quite simple, since it consists of dividing the

training set of N samples into k subsets ($k > 1$). The model is then trained on $k - 1$ subsets and validated by testing it on the left-out subset. This procedure is then iterated, leaving out each time a different subset for validation, and its mean-squared error is averaged on all cycles.

Regarding the MLP topology, a significant contribution came from a seminal paper by Bengio & LeCun (2007). In fact, they re-analyzed the implications of the Haykin pseudoteorem (Haykin 1998), proving that complex problems in which the mapping function is highly nonlinear and the local density of data in the PS is very variable are better matched by *deep* networks with more than one hidden computational layer.

3.3. Statistical Indicators

In order to evaluate and reciprocally compare the experiments described in the next section, we adopted the following definitions:

$$\text{bias}(x) = \frac{\sum_{i=1}^N x_i}{N} \quad (2)$$

$$\sigma(x) = \sqrt{\frac{\sum_{i=1}^N \left[x_i - \left(\frac{\sum_{i=1}^N x_i}{N} \right) \right]^2}{N}} \quad (3)$$

$$\text{MAD}(x) = \text{median}(|x|) \quad (4)$$

$$\text{NMAD}(x) = 1.48 \times \text{median}(|x|) \quad (5)$$

$$\text{rms}(x) = \sqrt{\frac{\sum_{i=1}^N x_i^2}{N}}, \quad (6)$$

where σ is the standard deviation, MAD is the median absolute deviation, NMAD is the normalized MAD, and rms is the root mean square. The term x in all of the above expressions may be either Δz defined as

$$\Delta z = (z_{\text{spec}} - z_{\text{phot}}) \quad (7)$$

or the normalized residuals Δz_{norm} defined as

$$\Delta z_{\text{norm}} = (z_{\text{spec}} - z_{\text{phot}})/(1 + z_{\text{spec}}). \quad (8)$$

4. THE EXPERIMENTS

Our approach is based on machine learning methods, and therefore it needs to be as automatic as possible in order to optimize the decisional support to the user (in this case the astronomer). Therefore, the results of the individual experiments as well as their comparison with others have to be evaluated in a consistent and objective manner through a homogeneous set of statistical indicators.

In what follows, we shall discuss the general experiment workflow and the outcome of the experiment phases.

4.1. The Knowledge Base and Model Setup

For machine-learning-supervised methods, it is a common practice to use the available KB to obtain at least three disjoint subsets for every experiment: one (training set) for training purposes, i.e., to train the method in order to acquire the hidden correlation among the input features, which is needed to perform the regression; the second one (validation set) to check the

training, in particular against loss of generalization capabilities (a phenomenon also known as overfitting); and the third one (test set) to evaluate the overall performances of the model. As a rule of thumb in the case of machine learning methods, these sets should be populated with, respectively, 60%, 20%, and 20% of the objects in the KB (Kearns 1996). In our case, however, we reduced the training+validation data amount (from 80% to 60%), driven by the past experience with the very accurate regression capabilities of the model also in the case of smaller KBs (Brescia et al. 2012; Cavuoti et al. 2012b), obtaining implicitly the possibility to verify its prediction performance on a larger test set, as well as a faster execution of the training phase. Furthermore, in order to ensure a proper coverage of the KB in the PS, the data objects were indeed divided up among the three data sets by random extraction, and this process is usually iterated several times to minimize the possible biases induced by fluctuations in the coverage of the PS, namely, small differences in the redshift distribution of training and test samples used in the experiments.

The first two criteria used to decide the stopping condition of the algorithm, as mentioned at the end of Section 3.2, are mainly sensitive to the choice of specific parameters and may lead to poor results if the parameters are improperly set. The cross validation does not suffer from such a drawback; it can avoid overfitting the data and is able to improve the generalization performance of the model. However, when compared to the traditional training procedures, the cross validation is much more computationally expensive. Therefore, by exploiting the cross validation technique (see Section 3.2), training and validation were indeed performed together using $\sim 60\%$ of the objects as a training + validation set, and the remaining $\sim 40\%$ as a test set.

The automatized process of the cross validation was done by performing 10 different training runs with the following procedure: (1) splitting of the training/validation set into 10 random subsets, each one composed of 10% of the data set; (2) at each training run, we applied 90% of the data set and then excluded 10% for validation.

As remarked in Section 3.2, the k -fold cross validation is able to avoid overfitting on the training set (Bishop 2006) with an increase of the execution time estimable around $k - 1$ times the total number of runs (Cavuoti et al. 2012b).

In terms of the internal parameter setup of the MLPQNA, we need to consider the following topological parameters.

1. *Input layer.* A variable number of neurons, corresponding to the pruned number of survey parameters used in all experiments, up to a maximum number of 43 nodes (all available features).
2. *Neurons on the first hidden layer.* A variable number of hidden neurons, depending on the number N of input neurons (features in the data set), equal to $2N + 1$ as a rule of thumb.
3. *Neurons on the second hidden layer.* A variable number of hidden neurons, ranging from 0 (to ignore the second layer) to $N - 1$.
4. *Output layer.* One neuron, returning the reconstructed photo- z value.

For the QNA learning rule, we heuristically fixed the following values as best parameters for the final experiments:

1. step: 0.0001 (one of the two stopping criteria. The algorithm stops if the approximation error step size is less than this value. A step value equal to zero means using the parameter MaxIt as a unique stopping criterion);

2. res: 40 (number of restarts of the Hessian approximation from random positions, performed at each iteration);
3. dec: 0.1 (regularization factor for weight decay. The term $\text{dec} * ||\text{network weights}||^2$ is added to the error function, where the network weights is the total number of weights in the network, directly depending on the total number of neurons inside. When properly chosen, the generalization performances of the network are highly improved);
4. MaxIt: 8000 (max number of iterations of the Hessian approximation. If zero, the step parameter is used as a stopping criterion);
5. CV (k): 10 (k -fold cross validation, with $k = 10$);
6. error evaluation: mean square error (between the target and network output).

Using these parameters, we obtained the statistical results reported in Section 4.4.

4.2. Selection of Features

As is known, supervised machine learning models are powerful methods for learning the hidden correlation between input and output features from training data. Of course, their generalization and prediction capabilities strongly depend on the intrinsic quality of data (S/N), the level of correlation among different features, and the amount of missing data present in the data set (Ripley 1996). It is obvious that some, possibly many, of the 43 parameters listed in Table 1 may not be independent and that their number needs to be reduced in order to speed up the computation (which scales badly with the number of features). This is a common problem in data mining and there is a wide body of literature on how to optimize the selection of features that are most relevant for a given purpose (Lindeberg 1998; Guyon & Elisseeff 2003, 2006; Brescia 2012). This process is usually called *feature selection* or *pruning*, and consists of finding a compromise between the number of features (and therefore the computational time) and the required accuracy of the final results. In order to do so, we extracted from the main catalog several subsets containing different groups of variables (features). Each one of these subsets was then analyzed separately in specific runs of the method (runs that in the data mining community are usually called experiments), in order to allow for comparison and evaluation. We wish to stress that our main concern was not only to disentangle which bands carry the most information but also, for a given band, which type of measurements (e.g., point-spread function (PSF), Petrosian, or isophotal magnitude) are more effective.

We performed a series of regression experiments to evaluate the performances obtained by the pruning of photometric quantities on the small data set described in Section 2. The pruning experiments consisted of several combinations of surveys and their features:

1. a *full* feature experiment to be used as a benchmark for all the other experiments;
2. some *service* experiments used to select the best combination of input features in order to eliminate redundancies in the flux measurements (i.e., Petrosian magnitudes against isophotal magnitudes);
3. *three-survey* experiments for all possible combinations of three (out of four) surveys;
4. *two-survey* experiments with all possible combinations of two (out of four) surveys;
5. *single-survey* experiments.

Table 1
Summary of the Data Extracted from the Databases of the Four Surveys and Merged to Form Our Final Catalog

Survey	Bands	Name of Feature	Synthetic Description
<i>GALEX</i>	<i>nuv, fuv</i>	mag, mag_iso mag_Aper_1 mag_Aper_2 mag_Aper_3 mag_auto and kron_radius	Near and Far UV total and isophotal mags phot. through 3, 4.5, and 7.5 arcsec apertures magnitudes and Kron radius in units of A or B
SDSS	<i>u, g, r, i, z</i>	psfMag	PSF fitting magnitude in the <i>u, g, r, i, z</i> bands
UKIDSS	<i>Y, J, H, K</i>	PsfMag AperMag3, AperMag4, AperMag6 HallMag, PetroMag	PSF fitting magnitude in the <i>Y, J, H, K</i> bands aperture photometry through 2, 2.8 and 5''.7 circular aperture in each band Calibrated magnitude within circular aperture r_hall and Petrosian magnitude in the <i>Y, J, H, K</i> bands
<i>WISE</i>	<i>W1, W2, W3, W4</i>	W1mpro, W2mpro, W3mpro, W4mpro	W1: 3.4 μm and 6''.1 angular resolution; W2: 4.6 μm and 6''.4 angular resolution; W3: 12 μm and 6''.5 angular resolution; W4: 22 μm and 12'' angular resolution. Magnitudes measured with profile-fitting photometry at the 95% level. Brightness upper limit if the flux measurement has $S/N < 2$
SDSS	...	z_{spec}	Spectroscopic redshift

Notes. Even though most names of the parameters are self-explanatory, we wish to remind that the various *psfMag* are magnitudes derived by integrating fluxes over the best fitting point-spread function. The aperture sizes refer to the radii.

The output of the experiments consisted of lists of photometric redshift estimates for all objects in the KB. All pruning experiments were performed using ~ 3000 objects in the training set and ~ 800 in the test set. In Table 2, we list the outcome of the experiments for the feature selection. Both $bias(\Delta z)$ and $\sigma(\Delta z)$ were computed using the objects in the test set alone. As can be seen, among the various types of magnitudes available for *GALEX* and UKIDSS, the best combination is obtained using the isophotal magnitudes for *GALEX* and the calibrated magnitudes (*HallMag*) for UKIDSS.

Therefore, at the end of the pruning phase, the best combination of features turned out to be the five SDSS *psfMag*, the two isophotal magnitudes of *GALEX*, the four *HallMag* for UKIDSS, and the four magnitudes for *WISE*.

4.3. Magnitudes versus Colors

Once the most significant features had been identified, we had to check which types of flux combinations were more effective, in terms of magnitudes or related colors. Experiments were performed on all five cross-matched data sets listed in Section 2.

As could be expected, the optimal combination turned out to be always the mixed one, i.e., the one including colors and one reference magnitude for each of the included surveys (*r* for SDSS, *nuv* for *GALEX*, *K* for UKIDSS, and *W4* for *WISE*). From the data mining point of view, this is rather surprising since the amount of information should not change by applying linear combinations between features. However, from the physical point of view, this can be easily understood by noticing that even though colors are derived as a subtraction of magnitudes, the content of information is quite different, since an ordering relationship is implicitly assumed, thus increasing the amount of information in the final output (gradients instead of fluxes). The additional reference magnitude instead removes the degeneracy in the luminosity class for a specific galaxy type.

4.4. MLPQNA Network Topology

The final check was about the hierarchical complexity of the network in terms of the number of internal layers, whether *shallow* or *deep* according to the definitions in Bengio & LeCun (2007), where *deep* is referred to as a feed-forward network with more than one hidden layer. The above-quoted cross-matched data sets were therefore processed through both a three-layer (input + hidden + output) and a four-layer (input + two hidden layers + output) network. In all cases, the four-layer network performed significantly better, thus confirming the performance enhancement with *deep* networks in the case of particularly complex nonlinear regression cases, i.e., in the case of highly multi-variate distributions of the input PS.

The experiments with the best results have been obtained using a four-layer network trained on the mixed (colors + reference magnitudes) data sets, and their statistics are reported in Tables 3–6.

5. DISCUSSION AND CONCLUSIONS

In 2002, we began to explore the usage of MLPs for the evaluation of photo-*z* both for *normal* galaxies and quasars (Tagliaferri et al. 2002). Several years later, D'Abrusco et al. (2007) used a combination of two MLPs to correct for the degeneracy introduced by the inhomogeneities in the KB. Then, Laurino et al. (2011) demonstrated that the subtleties in the mapping function could be more easily captured using the so-called Weak Gated Experts (WGE) method, a hierarchical combination of MLPs each specialized in a specific partition of the PS, whose individual outputs were combined by an additional MLP.

Furthermore, Bengio & LeCun (2007) published a seminal paper that somehow has disproved the Haykin pseudoteorem (Haykin 1998), pointing out that problems with a large amount of distribution irregularities in the PS are better treated by what they defined as *deep* networks, i.e., networks with more than

Table 2
Experiments for the Feature Selection Phase

<i>GALEX</i>	SDSS	UKIDSS	<i>WISE</i>	Bias (Δz)	σ (Δz)
Service experiments					
X	X	X	X	0.0033	0.174
X ^{1,2}	X	X ⁶	X	-0.0001	0.152
X ³	X	X ⁶	X	-0.0016	0.165
X ¹	X	X ⁶	X	0.0054	0.151
X ²	X	X ⁶	X	-0.0026	0.151
X ^{4,5}	X	X ⁶	X	-0.0008	0.152
X ^{1,2,3}	X	X ⁶	X	0.0041	0.163
X ^{2,3}	X	X ⁶	X	-0.0033	0.155
		X ^{6,7}		-0.0091	0.299
		X ⁷		0.0111	0.465
		X ⁶		-0.0081	0.294
Four survey experiment					
X ²	X	X ⁶	X	-0.0026	0.151
Three survey experiment					
X ²	X	X ⁶		-0.0046	0.152
X ²	X		X	0.0025	0.162
	X	X ⁶	X	-0.0032	0.179
X ²		X ⁶	X	0.0110	0.203
Two survey experiment					
		X ⁶	X	0.0045	0.236
X ²			X	0.0175	0.288
	X	X ⁶		-0.0027	0.210
	X		X	-0.0039	0.197
X ²	X			-0.0055	0.240
X ²		X ⁶		0.0133	0.238
One survey experiment					
			X	0.0165	0.297
X ^{1,2}	X			-0.0162	0.338
				0.0550	0.419
		X ⁶		-0.0081	0.294

Notes. Columns 1–4: surveys used for the experiment, where the superscript index indicates the used magnitudes: ¹ *mag*; ² *mag_iso*; ³ *magnitudes through 3, 4.5, and 7.5 arcsec apertures, mag_Aper 1, 2, and 3*; ⁴ *mag_auto*; ⁵ *kron_radius*; ⁶ *HallMag*; ⁷ *PetroMag*. A cross in a column means that the survey corresponding to that column was used for the experiment.

one computational (hidden) layer. In this paper, we exploited the Bengio & LeCun (2007) findings by using the supervised machine learning based method MLPQNA to evaluate photometric redshifts of quasars using multi-band data obtained from the cross-matching of the *GALEX*, SDSS, UKIDSS, and *WISE* surveys.

In Tables 3–5, we compare our best results to those presented by other authors (Ball et al. 2008; Richards et al. 2009; Laurino et al. 2011; Bovy et al. 2012), in terms of a homogeneous set of statistical indicators, defined in Section 3.3. Unfortunately, the whole set of indicators was not available for all bibliographical sources and in several cases we could only use a few quantities. Results are listed according to the combinations of surveys used in the experiment.

The best experiment, which makes use of a selected combination of parameters drawn from the four cross-matched surveys, leads to a bias = 0.004 and a median absolute deviation, MAD = 0.02. The fraction of catastrophic outliers, i.e., of objects with photo- z deviating more than 2σ from the spectroscopic value is $<3\%$, leading to $\sigma(\Delta z_{\text{norm}}) = 0.035$ after

Table 3
Comparison Among the Performances of the Different References

Exp	Bias (Δz)	σ (Δz)	MAD (Δz)	rms (Δz)
SDSS				
MLPQNA	0.007	0.25	0.102	0.26
Bovy et al.	...	0.46
Laurino et al.	0.210	0.28	0.110	0.35
Ball et al.	...	0.35
Richards et al.	...	0.52
SDSS + <i>GALEX</i>				
MLPQNA	0.003	0.21	0.060	0.22
Bovy et al.	...	0.26
Laurino et al.	0.13	0.21	0.061	0.25
Ball et al.	...	0.23
Richards et al.	...	0.37
SDSS + UKIDSS				
MLPQNA	0.001	0.25	0.066	0.26
Bovy et al.	...	0.28
SDSS + <i>GALEX</i> + UKIDSS				
MLPQNA	0.0009	0.18	0.043	0.19
Bovy et al.	...	0.21
SDSS + <i>GALEX</i> + UKIDSS + <i>WISE</i>				
MLPQNA	0.002	0.15	0.040	0.15

Notes. MLPQNA is related to our experiments, based on a four-layer network, trained on the mixed (colors + reference magnitudes) data sets. In some cases, the comparison references are not reported due to the missing statistics. Column 1: reference; Column 2–5, respectively: bias, standard deviation, MAD, and rms, calculated on $\Delta z = (z_{\text{spec}} - z_{\text{phot}})$ related to the test sets. For the definition of the parameters and for discussion see the text.

their removal (as reported in Table 6). The larger the number of surveys (bands) used, the more accurate the results are. This result, which might seem evident, is not obvious at all, since the higher amount of information carried by the additional bands implies also a strong decrease in the number of objects that are contained in the training set and should therefore cause a decrease in the generalization performances of the network.

This result, together with the fact that MLPQNA also performs well with small KBs (Cavuoti et al. 2012b), seems particularly interesting, since it has far reaching implications on ongoing and future surveys: a better photometric coverage is much more relevant than an increase of spectroscopic templates in the KB.

Concerning the performance evaluation in terms of photometric redshift reconstruction, all statistical results reported throughout this paper refer to test data sets only. In fact, it is a good practice to evaluate the results on data (i.e., the test set) that have never been presented to the network during the training and/or validation phases. The use of *test plus training* data might introduce an obvious positive systematic bias which could mask reality.

More generally, empirical methods, such as MLPQNA, have the advantage that the training set is made up of real sky objects. Hence, they do not suffer from the uncertainty of having accurate templates. In this sense, any empirical method intrinsically includes effects such as the filter bandpass and flux calibrations. In fact, as discussed in depth by Collister & Lahav (2004), one of the main drawbacks of these methods is the difficulty in

Table 4
Comparison Among the Performances of the Different References

Exp	Bias (Δz_{norm})	σ (Δz_{norm})	MAD (Δz_{norm})	rms (Δz_{norm})	NMAD (Δz_{norm})
SDSS					
MLPQNA	0.032	0.15	0.039	0.17	0.058
Laurino et al.	0.095	0.16	0.041	0.19	...
Ball et al.	0.095	0.18
Richards et al.	0.115	0.28
SDSS + GALEX					
MLPQNA	0.012	0.11	0.029	0.11	0.043
Laurino et al.	0.058	0.29	0.029	0.11	...
Ball et al.	0.06	0.12
Richards et al.	0.071	0.18
SDSS + UKIDSS					
MLPQNA	0.008	0.11	0.027	0.11	0.040
SDSS + GALEX + UKIDSS					
MLPQNA	0.005	0.087	0.022	0.088	0.032
SDSS + GALEX + UKIDSS + WISE					
MLPQNA	0.004	0.069	0.020	0.069	0.029

Notes. MLPQNA is related to our experiments, based on a four-layer network, trained on the mixed (colors + reference magnitudes) data sets. In some cases, the comparison references are not reported, due to the missing statistics. Column 1: reference; Columns 2–6, respectively: bias, standard deviation, MAD, rms, and NMAD calculated on $\Delta z_{\text{norm}} = (z_{\text{spec}} - z_{\text{phot}})/(1 + z_{\text{spec}})$ related to the test sets. For the definition of the parameters and for discussion see the text.

Table 5
Comparison in Terms of Outliers Percentages Among the Different References

Exp	Outliers ($ \Delta z $)		Outliers ($ \Delta z_{\text{norm}} $)	
	$> 2\sigma(\Delta z)$	$> 4\sigma(\Delta z)$	$> 2\sigma(\Delta z_{\text{norm}})$	$> 4\sigma(\Delta z_{\text{norm}})$
SDSS				
MLPQNA	7.68	0.38	6.53	1.24
Bovy et al.		0.51		
SDSS + GALEX				
MLPQNA	4.88	1.61	4.57	1.37
Bovy et al.		1.86		
SDSS + UKIDSS				
MLPQNA	4.00	1.73	3.82	1.38
Bovy et al.		1.92		
SDSS + GALEX + UKIDSS				
MLPQNA	2.86	1.47	3.05	0.23
Bovy et al.		1.13		
SDSS + GALEX + UKIDSS + WISE				
MLPQNA	2.57	0.87	2.88	0.91

Notes. In some cases, the comparison references are not reported, due to the missing statistics. Column 1: reference; Columns 2 and 3 are fractions of outliers at different σ based on $\Delta z = (z_{\text{spec}} - z_{\text{phot}})$; Columns 4 and 5 are the fractions of outliers at different σ based on $\Delta z_{\text{norm}} = (z_{\text{spec}} - z_{\text{phot}})/(1 + z_{\text{spec}})$. Column 4 reports our catastrophic outliers, defined as $|\Delta z_{\text{norm}}| > 2\sigma(\Delta z_{\text{norm}})$.

extrapolating to regions of the input PS that are not well sampled by the training data. Therefore, the efficiency of empirical methods degrades for objects at fainter magnitudes than those included in the training set, as this would require an extrapolation capability for data having properties, such as redshift and photometry, not included in the learned sample. In fact, another strong requirement of such methods is that the training set must be large enough to properly cover the PS in

terms of colors, magnitudes, object types, and redshift. In this case, the calibrations and corresponding uncertainties are well known, and only limited extrapolations beyond the observed locus in color–magnitude space are required. In conclusion, under the conditions described above about the consistency of the training set, a realistic way to measure photometric uncertainties is to compare the photometric redshift estimation with spectroscopic measures in the test samples.

Table 6
Catastrophic Outliers Evaluation and Comparison between the Residual $\sigma_{\text{clean}}(\Delta z_{\text{norm}})$ and NMAD(Δz_{norm})

Exp	No. of Obj.	σ (Δz_{norm})	% Catas. Outliers	σ_{clean} (Δz_{norm})	NMAD (Δz_{norm})
SDSS	41431	0.15	6.53	0.062	0.058
SDSS + <i>GALEX</i>	17876	0.11	4.57	0.045	0.043
SDSS+UKIDSS	12438	0.11	3.82	0.041	0.040
SDSS+ <i>GALEX</i> +UKIDSS	5836	0.087	3.05	0.040	0.032
SDSS+ <i>GALEX</i> +UKIDSS+ <i>WISE</i>	5716	0.069	2.88	0.035	0.029

Notes. The reported number of objects for each cross-matched catalog refers to the test sets only. Catastrophic outliers are defined as objects where $|\Delta z_{\text{norm}}| > 2\sigma(\Delta z_{\text{norm}})$. The standard deviation $\sigma_{\text{clean}}(\Delta z_{\text{norm}})$ is calculated after removing the catastrophic outliers, i.e., on the data sample for which $|\Delta z_{\text{norm}}| \leq 2\sigma(\Delta z_{\text{norm}})$.

As can be seen in Tables 3–5, in all cases MLPQNA obtains very relevant results. Only in the SDSS+*GALEX* case do the non-normalized quantities (i.e., those referred to the error $\Delta z = z_{\text{spec}} - z_{\text{phot}}$) show a substantial agreement between our results and those by Laurino et al. (2011). The better performances of MLPQNA in the normalized indicators (i.e., those referred to the error $\Delta z_{\text{norm}} = (z_{\text{spec}} - z_{\text{phot}})/(1 + z_{\text{spec}})$) are a consequence of the better performance of the MLPQNA method in terms of the fraction of catastrophic outliers.

We wish to stress that both our four-layer MLPQNA and the WGE method discussed in Laurino et al. (2011) take advantage of a substantial improvement in complexity with respect to the traditional three-layer MLP networks used in the literature, and therefore deal better with the complexity of the multi-color PS. Average statistical indicators such as bias and standard deviation, however, provide only part of the information that allows us to correctly evaluate the performances of a method and, for instance, they provide only very little evidence of the systematic trends that are observed as a sudden increase in the residuals spread over specific regions of the redshift space (Laurino et al. 2011). In the worst cases, these regions correspond to degeneracies in the PS and, as could be expected, the relevance of such degeneracies decreases for an increasing number of bands.

As far as the analysis of the catastrophic outliers is concerned, according to Mobasher et al. (2007), the parameter $D_{95} \equiv \Delta_{95}/(1 + z_{\text{phot}})$ enables the identification of outliers in photometric redshifts derived through SED fitting methods (usually evaluated through numerical simulations based on mock catalogs). In fact, in the hypothesis that the redshift error $\Delta z_{\text{norm}} = (z_{\text{spec}} - z_{\text{phot}})/(1 + z_{\text{spec}})$ is Gaussian, the catastrophic redshift error limit would be constrained by the width of the redshift probability distribution, corresponding to the 95% confidence interval, i.e., with $\Delta_{95} = 2\sigma(\Delta z_{\text{norm}})$. In our case, however, photo- z are empirical, i.e., not based on any specific fitting model, and it is preferable to use the standard deviation value $\sigma(\Delta z_{\text{norm}})$ derived from the photometric cross-matched samples, although it could overestimate the theoretical Gaussian σ , due to the residual spectroscopic uncertainty as well as to the method training error. Therefore, we consider as catastrophic outliers the objects with $|\Delta z_{\text{norm}}| > 2\sigma(\Delta z_{\text{norm}})$. It is also important to note that for empirical methods, it is useful to analyze the correlation between the $\text{NMAD}(\Delta z_{\text{norm}}) = 1.48 \times \text{median}(|\Delta z_{\text{norm}}|)$ and the standard deviation $\sigma_{\text{clean}}(\Delta z_{\text{norm}})$ calculated for the data sample for which $|\Delta z_{\text{norm}}| \leq 2\sigma(\Delta z_{\text{norm}})$. In fact, the quantity NMAD would be comparable to the value of the σ_{clean} .

As is shown in Table 6, in our data the $\sigma_{\text{clean}}(\Delta z_{\text{norm}})$ is always slightly larger than the corresponding $\text{NMAD}(\Delta z_{\text{norm}})$, which is

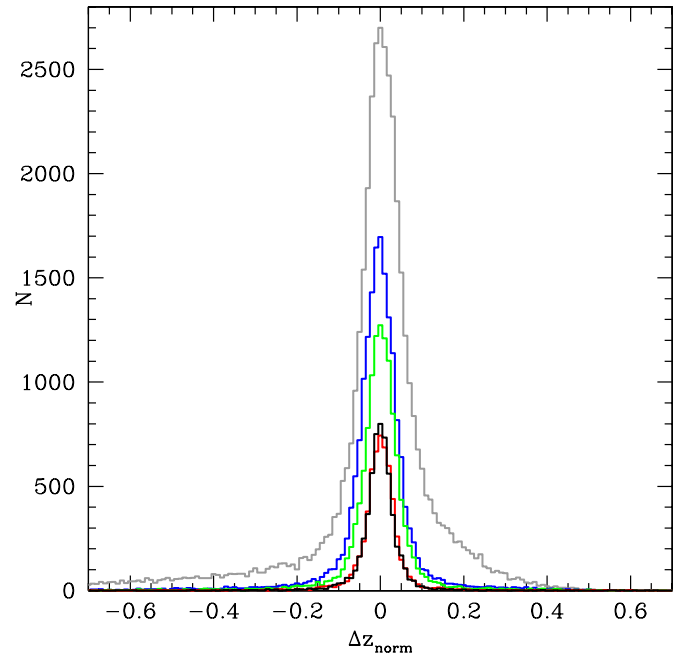


Figure 4. Δz_{norm} distributions for all five cross-matched test data sets. The lines referred to are, respectively, SDSS (gray), SDSS+*GALEX* (blue), SDSS+UKIDSS (green), SDSS+*GALEX*+UKIDSS (red), and SDSS+*GALEX*+UKIDSS+*WISE* (black).

(A color version of this figure is available in the online journal.)

exactly what is expected due to the overestimate induced by the above considerations (see also Figures 4 and 5).

Finally, we would like to stress that the difficulties encountered by us and by other teams in comparing different methods, especially in light of the crucial role that photo- z play in the scientific exploitation of present and future large surveys (cf. The Dark Energy Survey Collaboration 2005, Chambers 2011, Refregier et al. 2010), confirm that it would be desirable to re-propose an upgraded version of the extremely useful PHAT contest (Hildebrandt et al. 2010, Cavuoti et al. 2012b), which allowed a direct, effective, and non-ambiguous comparison of different methods applied on the same data sets and through the same set of statistical indicators. This new contest should be applied to a much larger data set, with a more practical selection of photometric bands, and should also take into account other parameters such as the scalability and robustness of the algorithms, as well as the degeneracy characterization.

The authors thank the anonymous referee for comments and suggestions that helped to improve the paper, and the whole

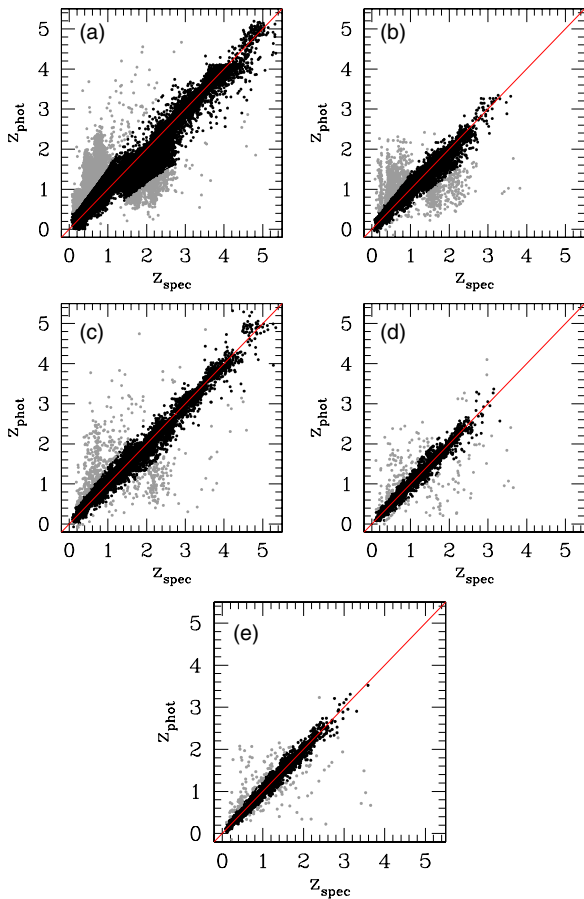


Figure 5. Scatter plots (z_{spec} vs. z_{phot}); (a) SDSS, (b) SDSS+GALEX, (c) SDSS+UKIDSS, (d) SDSS+GALEX+UKIDSS, and (e) SDSS+GALEX+UKIDSS+WISE. All diagrams refer to results on test sets. Gray points are catastrophic outliers (defined in Table 5). The red line is the dot-to-dot straight line passing through photometric and spectroscopic redshift limits in the available knowledge base.

(A color version of this figure is available in the online journal.)

DAMEWARE team⁶ for many useful discussions. The authors also thank the financial support of Project F.A.R.O., third call by the University Federico II of Naples, and of the PRIN-MIUR 2011 for Euclid Mission. G.L. thanks Professor G. S. Djorgovski and the whole Department of Astronomy at the California Institute of Technology in Pasadena for their hospitality. A.M. and M.B. acknowledge the financial support of PRIN-INAF 2010, *Architecture and Tomography of Galaxy Clusters*. R. D’A. acknowledges the financial support of the US Virtual Astronomical Observatory, which is sponsored by the National Science Foundation and the National Aeronautics and Space Administration.

APPENDIX

THE QUASI NEWTON LEARNING RULE

Most Newton methods use the Hessian of the function to find the stationary point of a quadratic form. It needs to be stressed, however, that the Hessian of a function is not always available and in many cases is far too complex to be computed in an analytical way. More often, it is easier to compute the function gradient, which can be used to approximate the Hessian via N consequent gradient calculations. In order to better understand

why QNAs are so powerful, it is convenient to start from the classical and quite common Gradient Descent Algorithm (GDA) used for back propagation (Bishop 2006). In GDA, the direction of each updating step for the MLP weights is derived from the error descent gradient, while the length of the step is determined from the learning rate. In the case of particularly complex problems, this method is inaccurate and ineffective, and therefore may get stuck in local minima. A more effective approach is to move toward the negative direction of the gradient (*line search direction*) not by a fixed step, but by moving toward the minimum of the function along that direction. This can be achieved by first deriving the descent gradient and then by analyzing it with the variation of the learning rate (Brescia 2012). Let us suppose that at step t , the current weight vector is $w^{(t)}$, and let us consider a search direction $d^{(t)} = -\nabla E^{(t)}$. If we select the parameter λ in order to minimize $E(\lambda) = E(w^{(t)} + \lambda d^{(t)})$, then the new weight vector can be expressed as

$$w^{(t+1)} = w^{(t)} + \lambda d^{(t)} \quad (\text{A1})$$

and the problem of *line search* becomes a one-dimensional minimization problem that can be solved in many different ways. Simple variants are (1) to move $E(\lambda)$ by varying λ by small intervals, then evaluate the error function at each new position, and stop when the error begins to increase; or (2) to use the parabolic search for a minimum and compute the parabolic curve crossing pre-defined learning rate points. The minimum d of the parabolic curve is a good approximation of the minimum of $E(\lambda)$ and it can be derived by means of the parabolic curve which crosses the fixed points with the lowest error values.

Another approach instead makes use of *trust region* based strategies that minimize the error function by iteratively growing or contracting the region of the function by adjusting a quadratic model function that best approximates the error function. In this sense, this technique can be considered as a dual to line search, since it tries to find the best size of the region by fixing the step size (while the line search strategy always chooses the step direction before selecting the step size; Celis et al. 1985). All these approaches, however, rely on the assumption that the optimal search direction is given at each step by the negative gradient: an assumption that not only is not always true, but can also lead to serious wrong convergence. In fact, if the minimization is done along the negative gradient direction, then the subsequent search direction (the new gradient) will be orthogonal to the previous one: in fact, note that when the line search finds the minimum, it is

$$\frac{\partial E}{\partial \lambda}(w^{(t)} + \lambda d^{(t)}) = 0, \quad (\text{A2})$$

and hence

$$g^{(t+1)T} d^{(t)} = 0, \quad (\text{A3})$$

where $g \equiv \nabla E$. The iteration of the process therefore leads to oscillations of the error function which slow down the convergence process. The method implemented here relies on selecting other directions so that the gradient component, parallel to the previous search direction, would remain unchanged at each step. Suppose that you have already minimized with respect to the direction $d^{(t)}$ starting from the point $w^{(t)}$ and reaching the point $w^{(t+1)}$, where Equation (A3) becomes

$$g(w^{(t+1)})^T d^{(t)} = 0 \quad (\text{A4})$$

⁶ http://dame.dsf.unina.it/project_members.html

by choosing $d^{(t+1)}$ so as to preserve the gradient component parallel to $d^{(t)}$ equal to zero. It is possible to build a sequence of directions d in such a way that each direction is conjugated to the previous one on the dimension $|w|$ of the search space (this is known as a conjugate gradient method; Golub & Ye 1999). In the presence of a squared error function, the updated weight algorithm is

$$w^{(t+1)} = w^{(t)} + \alpha^{(t)} d^{(t)} \quad (\text{A5})$$

with

$$\alpha^{(t)} = -\frac{d^{(t)T} g^{(t)}}{d^{(t)T} H d^{(t)}}. \quad (\text{A6})$$

Furthermore, d can be obtained for the first time via the negative gradient, and in the subsequent iterations as a linear combination of the current gradient and of the previous search directions:

$$d^{(t+1)} = -g^{(t+1)} + \beta^{(t)} d^{(t)} \quad (\text{A7})$$

with

$$\beta^{(t)} = \frac{g^{(t+1)T} H d^{(t)}}{d^{(t)T} H d^{(t)}}. \quad (\text{A8})$$

This algorithm finds the minimum of a square error function at most in $|w|$ steps but at a high computational cost, since in order to determine the values of α and β it makes use of that *Hessian matrix* H which, as we already mentioned, is very demanding in terms of computing time: a fact that puts serious constraints on the application of this family of methods to medium/large data sets. Excellent approximations for the coefficients α and β can, however, be obtained from analytical expressions that do not use the Hessian matrix explicitly. For instance, β can be calculated through any one of the following expressions (respectively, Hestenes & Stiefel 1952; Fletcher & Reeves 1964; Polak & Ribiere 1969):

$$\text{Hestenes–Stiefel : } \beta^{(t)} = \frac{g^{(t+1)T} (g^{(t+1)} - g^{(t)})}{d^{(t)T} (g^{(t+1)} - g^{(t)})} \quad (\text{A9})$$

$$\text{Fletcher–Reeves : } \beta^{(t)} = \frac{g^{(t+1)T} g^{(t+1)}}{g^{(t)T} g^{(t)}} \quad (\text{A10})$$

$$\text{Polak–Ribiere : } \beta^{(t)} = \frac{g^{(t+1)T} (g^{(t+1)} - g^{(t)})}{g^{(t)T} g^{(t)}}. \quad (\text{A11})$$

These expressions are all equivalent if the error function is square-typed, otherwise they assume different values. Typically, the Polak–Ribiere equation obtains better results because, if the algorithm is slow and subsequent gradients are quite alike between them, its equation produces values of β such that the search direction tends to assume the negative gradient direction (Vetterling & Flannery 1992).

Concerning the parameter α , its value can be obtained by using the line search method directly. The method of conjugate gradients reduces the number of steps to minimize the error up to a maximum of $|w|$ because there could be almost $|w|$ conjugate directions in a $|w|$ -dimensional space. In practice, however, the algorithm is slower because, during the learning process, the property *conjugate* of the search directions tends to deteriorate. In order to avoid the deterioration, it is useful to restart the algorithm after $|w|$ steps by resetting the search direction with the negative gradient direction.

By using a local square approximation of the error function, we can obtain an expression for the minimum position. The gradient in every point w is in fact given by

$$\nabla E = H \times (w - w^*), \quad (\text{A12})$$

where w^* corresponds to the minimum of the error function, which satisfies the condition

$$w^* = w - H^{-1} \times \nabla E. \quad (\text{A13})$$

The vector $-H^{-1} \times \nabla E$ is known as a Newton direction and it is the base for a variety of optimization strategies, such as, for instance, the QNA, which instead of calculating the H matrix and then its inverse, uses a series of intermediate steps of lower computational cost to generate a sequence of matrices that are more and more accurate approximations of H^{-1} . From the Newton formula (Equation (A13)) we note that the weight vectors on steps t and $t+1$ are correlated with the correspondent gradients by the formula

$$w^{(t+1)} - w^{(t)} = -H^{(-1)}(g^{(t+1)} - g^{(t)}), \quad (\text{A14})$$

which is known as the *Quasi Newton Condition*. The approximation G is therefore built in order to satisfy this condition. The formula for G is

$$G^{(t+1)} = G^{(t)} + \frac{pp^T}{p^T v} - \frac{(G^{(t)}v)v^T G^{(t)}}{v^T G^{(t)}v} + (v^T G^{(t)}v)uu^T, \quad (\text{A15})$$

where the vectors are

$$p = w^{(t+1)} - w^{(t)}; v = g^{(t+1)} - g^{(t)}; u = \frac{p}{p^T v} - \frac{G^{(t)}v}{v^T G^{(t)}v}. \quad (\text{A16})$$

Using the identity matrix to initialize the procedure is equivalent to considering, step by step, the direction of the negative gradient while, at each next step, the direction $-Gg$ is definitely a descent direction. The above expression could carry the search out of the interval of validity for the squared approximation. The solution is hence to use the *line search* to find the minimum of the function along the search direction. By using such a system, the weight updating expression (Equation (A5)) can be formulated as follows

$$w^{(t+1)} = w^{(t)} + \alpha^{(t)} G^{(t)} g^{(t)}, \quad (\text{A17})$$

where α is obtained by the *line search*.

One of the main advantages of QNA, compared with conjugate gradients, is that the *line search* does not require the calculation of α with a high precision because it is not a critical parameter. Unfortunately, however, it again requires a large amount of memory to calculate the matrix G ($|w| \times |w|$) for large $|w|$. One way to reduce the required memory is to replace at each step the matrix G with a unitary matrix. With such replacement and after multiplying by g (the current gradient), we obtain

$$d^{(t+1)} = -g^{(t)} + A p + B v. \quad (\text{A18})$$

Note that if the line search returns exact values, then the above equation produces mutually conjugate directions. A and B are scalar values defined as

$$A = -\left(1 + \frac{v^T v}{p^T v}\right) \frac{p^T g^{(t+1)}}{p^T v} + \frac{v^T g^{(t+1)}}{p^T v}$$

$$B = \frac{p^T g^{(t+1)}}{p^T v}. \quad (\text{A19})$$

REFERENCES

- Abazajian, K. N., Adelman-McCarthy, J. K., Agüeros, M. A., et al. 2009, *ApJS*, **182**, 543
- Aihara, H., Allende Prieto, C., An, D., et al. 2011, *ApJS*, **193**, 29
- Ball, N. M., Brunner, R. J., Myers, A. D., et al. 2008, *ApJ*, **683**, 12
- Bengio, Y., & LeCun, J. 2007, *Large-Scale Kernel Machines* (Cambridge, MA: MIT Press)
- Bishop, C. M. 2006, *Pattern Recognition and Machine Learning* (Berlin: Springer)
- Bovy, J., Myers, A. D., Hennawi, J. F., et al. 2012, *ApJ*, **749**, 41
- Brescia, M. 2012, in *New Trends in E-Science: Machine Learning and Knowledge Discovery in Databases. In Horizons in Computer Science Research*, ed. Thomas S. Clary (Series Horizons in Computer Science, Vol. 7; Hauppauge, NY: Nova Science Publishers)
- Brescia, M., Cavuoti, S., Paolillo, M., Longo, G., & Puzia, T. 2012, *MNRAS*, **421**, 1155
- Broyden, C. G. 1970, *JIMIA*, **6**, 76
- Budavari, T., Csabai, I., Szalay, A. S., et al. 2001, *AJ*, **122**, 1163
- Byrd, R. H., Nocedal, J., & Schnabel, R. B. 1994, *MatPr*, **63**, 129
- Cavuoti, S., Brescia, M., Longo, G., Garofalo, M., & Nocella, A. 2012a, *Science—Image in Action* (Singapore: World Scientific), **241**
- Cavuoti, S., Brescia, M., Longo, G., & Mercurio, A. 2012b, *A&A*, **546**, 1
- Celis, M., Dennis, J. E., & Tapia, R. A. 1985, in *Numerical Optimization*, ed. P. Boggs, R. Byrd, & R. Schnabel (Philadelphia, PA: SIAM), 71
- Chambers, K. C. 2011, *BAAS*, **43**, 222.02
- Collister, A. A., & Lahav, O. 2004, *PASP*, **116**, 345
- D’Abrusco, R., Longo, G., & Walton, N. A. 2009, *MNRAS*, **396**, 223
- D’Abrusco, R., Staiano, A., Longo, G., et al. 2007, *ApJ*, **663**, 752
- Davidon, W. C. 1968, *CompJ*, **10**, 406
- Fernandez-Soto, A., Lanzetta, K. M., Chen, H. W., Pascarelle, S. M., & Yahata, N. 2001, *ApJS*, **135**, 41
- Fletcher, R. 1970, *CompJ*, **13**, 317
- Fletcher, R., & Reeves, C. M. 1964, *CompJ*, **7**, 149
- Floudas, C. A., & Jongen, H. T. 2005, *J. Glob. Optim.*, **32**, 409
- Geisser, S. 1975, *J. Am. Stat. Assoc.*, **70**, 320
- Giannantonio, T., Scranton, R., Crittenden, R. G., et al. 2008, *PhRvD*, **77**, 123520
- Goldfarb, D. 1970, *MaCom*, **24**, 23
- Golub, G. H., & Ye, Q. 1999, *SIAM J. Sci. Comput.*, **21**, 1305
- Guyon, I., & Elisseeff, A. 2003, *J. Mach. Learn. Res.*, **3**, 1157
- Guyon, I., & Elisseeff, A. 2006, in *Feature Extraction, Foundations and Applications*, ed. I. Guyon, S. Gunn, M. Nikravesh, & L. A. Zadeh (Studies in Fuzziness and Soft Computing Series; New York: Springer), 207
- Haykin, S. 1998, *Neural Networks: A Comprehensive Foundation*, Vol. 2 (Englewood Cliffs, NJ: Prentice-Hall)
- Hennawi, J. F., Strauss, M. A., Oguri, M., et al. 2006, *AJ*, **131**, 1
- Hestenes, M. R., & Stiefel, E. 1952, *JRNBS*, **49**, 409
- Hildebrandt, H., Arnouts, S., Capak, P., et al. 2010, *A&A*, **523**, 21
- Kearns, M. 1996, in *A Bound on the Error of Cross Validation Using the Approximation and Estimation Rates, with Consequences for Training-Test Split*, *Neural Information Processing 8*, ed. D. S. Touretzky, M. C. Mozer, & M. E. Hasselmo (San Mateo, CA: Morgan Kaufmann Publishers), 183
- Lanzetta, K. M., Yahil, A., & Fernandez-Soto, A. 1998, *AJ*, **116**, 1066
- Laurino, O., D’Abrusco, R., Longo, G., & Riccio, G. 2011, *MNRAS*, **418**, 2165
- Lawrence, A., Warren, S. J., Almaini, O., et al. 2007, *MNRAS*, **379**, 1599
- Lindeberg, T. 1998, *Int. J. Comput. Vis.*, **30**, 77
- Marlin, B. M. 2008, PhD thesis, Univ. of Toronto
- Martin, D. C., Fanson, J., Schiminovich, D., et al. 2005, *ApJ*, **619**, L1
- Mobasher, B., Capak, P., Scoville, N. Z., et al. 2007, *ApJS*, **172**, 117
- Myers, A. D., Brunner, R. J., Richards, G. T., et al. 2006, *ApJ*, **638**, 622
- Polak, E., & Ribiere, G. 1969, *Revue Fr. Inf. Rech. Oper.* **3**, 35
- Refregier, A., Amara, A., Kitching, T. D., et al. 2010, *Euclid Imaging Consortium Science Book* (arXiv:1001.0061)
- Richards, G. T., Fan, X., Schneider, D. P., et al. 2001a, *AJ*, **121**, 2308
- Richards, G. T., Myers, A. D., Gray, A. G., et al. 2009, *ApJS*, **180**, 67
- Richards, G. T., Weinstein, M. A., Schneider, D. P., et al. 2001b, *AJ*, **122**, 1151
- Ripley, B. D. 1996, *Pattern Recognition and Neural Networks* (Cambridge: Cambridge Univ. Press)
- Schneider, D. P., Richards, G. T., Hall, P. B., et al. 2010, *AJ*, **139**, 2360
- Scranton, R., Ménard, B., Richards, G. T., et al. 2005, *ApJ*, **633**, 589
- Shanno, D. F. 1970, *MaCom*, **24**, 647
- Shanno, D. F. 1990, *Recent Advances in Numerical Techniques for Large-scale Optimization, Neural Networks for Control* (Cambridge, MA: MIT Press)
- Sylvain, A., & Celisse, A. 2010, *Stat. Surv.*, **4**, 40
- Tagliaferri, R., Longo, G., Andreon, S., et al. 2002, *Neural Networks and Photometric Redshifts* (Berlin: Springer), arXiv:astro-ph/0203445
- The Dark Energy Survey Collaboration 2005, *The Dark Energy Survey*, White Paper, 42 (arXiv:0510346)
- Veron-Cetty, M.-P., & Veron, P. 2000, *A Catalogue of Quasars and Active Nuclei*, ESO Scientific Report, No. 19
- Vetterling, T., & Flannery, B. P. 1992, *Conjugate Gradients Methods in Multidimensions. Numerical Recipes in C—The Art of Scientific Computing*, ed. W. H. Press & S. A. Teukolsky (2nd ed.; Cambridge: Cambridge Univ. Press)
- Wolf, C., Meisenheimer, K., Kleinheinrich, M., et al. 2004, *A&A*, **421**, 913
- Wright, E. L., Eisenhardt, P. R. M., Mainzer, A. K., et al. 2010, *AJ*, **140**, 1868