

Quantifying correlations between galaxy emission lines and stellar continua

Róbert Beck^{1,2*}, László Dobos^{1,2}, Ching-Wa Yip^{2,3},
Alexander S. Szalay² and István Csabai¹

¹*Department of Physics of Complex Systems, Eötvös Loránd University, 1117 Budapest, Hungary*

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

³*Wolfram Research, Somerville, MA 02144, USA*

Accepted 2015 December 22. Received 2015 November 21; in original form 2015 June 4

ABSTRACT

We analyse the correlations between continuum properties and emission line equivalent widths of star-forming and active galaxies from the Sloan Digital Sky Survey. Since upcoming large sky surveys will make broad-band observations only, including strong emission lines into theoretical modelling of spectra will be essential to estimate physical properties of photometric galaxies. We show that emission line equivalent widths can be fairly well reconstructed from the stellar continuum using local multiple linear regression in the continuum principal component analysis (PCA) space. Line reconstruction is good for star-forming galaxies and reasonable for galaxies with active nuclei. We propose a practical method to combine stellar population synthesis models with empirical modelling of emission lines. The technique will help generate more accurate model spectra and mock catalogues of galaxies to fit observations of the new surveys. More accurate modelling of emission lines is also expected to improve template-based photometric redshift estimation methods. We also show that, by combining PCA coefficients from the pure continuum and the emission lines, automatic distinction between hosts of weak active galactic nuclei (AGNs) and quiescent star-forming galaxies can be made. The classification method is based on a training set consisting of high-confidence starburst galaxies and AGNs, and allows for the similar separation of active and star-forming galaxies as the empirical curve found by Kauffmann et al. We demonstrate the use of three important machine learning algorithms in the paper: k -nearest neighbour finding, k -means clustering and support vector machines

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

1 INTRODUCTION

Stellar population synthesis models are very successful in explaining the spectral energy distribution of galaxies in the optical (Fioc & Rocca-Volmerange 1997; Bruzual & Charlot 2003; Maraston & Strömbäck 2011; Vazdekis et al. 2012) but they do not account for the characteristic emission lines originating from the excited interstellar gas. Starburst galaxies and galaxies with an active nucleus can produce emission lines so strong that can reach 60 per cent of the continuum flux in certain bands, or as much as 1 mag (Atek et al. 2011). As a result, pure population synthesis models are

not enough to account for observations made with broad-band photometric filters. Since future large sky surveys will make photometric observations only, accurate modelling of the emission lines will be essential to estimate physical properties (including photometric redshifts) of galaxies precisely.

The purpose of this paper is to empirically quantify correlations between properties of the stellar continuum of galaxy spectra, and the strengths of emission lines. We also propose a recipe for generating realistic emission lines in the optical regime for stellar continua taken from population synthesis models. Moreover, we present a novel classification method to differentiate between starburst and active galaxies.

Results presented in the paper are obtained with the help of three important, widely used machine learning tech-

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

niques that have just started to gain popularity in astronomical data analysis. Local linear regression using nearest neighbours (Csabai et al. 2007; Kerekes et al. 2013) has been used for physical parameter estimation based on broad-band photometry. *k*-means clustering, an automatic, unsupervised classification algorithm has been applied successfully, for instance, to classify gamma-ray bursts (Chattopadhyay et al. 2007; Veres et al. 2010). Support vector machines (SVM), a supervised classification algorithm has been used for star-galaxy separation (Kovács & Szapudi 2015) and transient detection (Wright et al. 2015). We will briefly introduce these methods later in the paper. For a detailed introduction to the field, refer to Ivezic et al. (2014).

As with all training set-based empirical methods, the validity of our results is limited to the training set’s coverage of the parameter space (in our case the redshift, metallicity, luminosity, continuum and line properties). Extrapolation capabilities of empirical techniques to parameter ranges outside the coverage is usually poor compared to theoretical models. While this certainly constrains the applicability of our results to strong emission line galaxies of the Sloan Digital Sky Survey (SDSS), the method itself can be easily extended to galaxies outside the investigated sample by augmenting the training set.

The structure of the paper is as follows. In Sec. 2, we explain the sample selection and data reduction methods. Sec. 3 describes the line reconstruction methods we investigated. An empirical method for star-forming–active galactic nucleus (AGN) separation is given in Sec. 4. We present a stochastic procedure to generate realistic emission lines for continuum models in Sec. 5. We summarize our findings and outline future work in Sec. 6.

Wavelengths are generally quoted in vacuum. More information on the data used for this study, additional figures and program source code are available on the web site of the paper¹. Colour versions of the figures are available in the online version of the paper.

1.1 Earlier work

Thanks to the large amount of flux-calibrated optical galaxy spectra accumulated by the SDSS, precision of galaxy spectrum modelling has been improved significantly during the last decade. Many software tools and libraries exist to generate realistic stellar continua from a prescribed star formation history and various libraries of single stellar population spectra with a wide range of metallicities and initial mass function choices (Fioc & Rocca-Volmerange 1997; Bruzual & Charlot 2003; Maraston & Strömbäck 2011; Bressan et al. 2012; Vazdekis et al. 2012). Models have also been extended with descriptions of interstellar extinction, the UV–IR balance (Silva et al. 1998; Charlot & Fall 2000; da Cunha et al. 2010) and the chemical evolution of the gas from which stars can form (Davé, Finlator & Oppenheimer 2011).

Emission lines of galaxy spectra carry a large amount of information about the abundance and ionization states of elements in the interstellar gas. Based on ionization ratios of the various elements, the source of primary radiation responsible for the excitation of the interstellar medium (ISM)

can be characterized (Baldwin, Phillips & Terlevich 1981; Kewley et al. 2001; Kauffmann et al. 2003b) The two primary radiation sources are young, hot, massive stars and active galactic nuclei (AGNs). Their different spectra (thermal and power law, respectively) cause different ionization states and ratios of the most common elements which, in turn, produce well measurable, strong, often broad emission lines: the Balmer series of hydrogen, [OII], [OIII], [NII], [SII], etc. Population synthesis models do not account for the emission of the ISM.

Photoionization models (Stasińska 1984; Ferland et al. 2013) yield accurate line ratios for any primary radiation spectrum and gas composition. To couple stellar population synthesis with models of photoionization, shock-heating of the interstellar gas, emission of the dust etc., the star formation history, several interactions between the stellar populations, the active nucleus, the dust and gas content need to be accounted for. For instance, AGNs are very likely to be responsible for quenching rapid star formation following starburst periods in the galaxy but they also emit ionizing radiation that excites gas, evaporates dust and produces shock waves that heat the ISM. Also, starburst periods are followed by high supernova activity that enriches the ISM with metals, leading to significant chemical evolution which must be reflected in the models of emission lines. Additionally, a recent advancement in stellar population synthesis is the inclusion of stellar rotation and binary evolution effects which have been shown to noticeably influence the strengths of some emission lines (Eldridge & Stanway 2012; Stanway et al. 2014; Leitherer et al. 2014; Topping & Shull 2015). Taking everything into account is not possible without detailed hydrodynamic simulation of the galaxies (Jonsson, Groves & Cox 2010; Kewley et al. 2013) or without making significant simplifications to the models. Various software, notably PÉGASE and BPASS (Fioc & Rocca-Volmerange 1997; Eldridge & Stanway 2012), can be used to generate emission lines on top of stellar continua computed from stellar population synthesis. The photoionization part of these softwares, however, introduces a large set of free parameters that describe the distribution and composition of the ISM. A frequently used way of reducing the number of free parameters is to make theoretical or empirical assumptions. Typical theoretical simplifications include the assumption of spherical symmetry or the use of a common ionization spectrum for all gas clouds (Stasińska 1984; Fioc & Rocca-Volmerange 1997; Ferland et al. 2013). If no strict physical considerations can be made, to generate realistic emission lines on top of modelled continua using any photoionization code, one has to estimate the *a priori* distribution of model parameters by comparing large ensembles of models with observations. For instance, the code Le Phare (Ilbert et al. 2006) uses the relations of Kennicutt (1998) to parametrize emission lines.

Another route to take to generate realistic emission lines is to work on an entirely empirical basis. Yip et al. (2004) demonstrated that stellar continua of SDSS galaxies form a 1D sequence and thus, can be characterized by a single numerical value, the `eclass`. The value of `eclass` for each galaxy spectrum is obtained by expressing the continuum on a basis derived from principal component analysis (PCA). Gyóry et al. (2011) showed that strong correlations between the `eclass` (i.e. the stellar continua) of starburst and AGN

¹ <http://www.vo.elte.hu/papers/2015/emissionlines>

galaxies exist. They applied PCA to expand emission line equivalent widths (EWs) of SDSS galaxies on a 3D basis and correlate the principal components with the `eclass` of the continua. We take their approach a step further: based on the correlations, we give a recipe to automatically generate emission lines with realistic distribution and refine star-forming–AGN separation using the principal components and SVM.

2 DATA REDUCTION

We started with the entire spectroscopic galaxy sample of the Sloan Digital Sky Survey Data Release 7 (?) which we later filtered by signal-to-noise ratio and line strength. As one of our goals was to accurately fit broad AGN lines, we measured line parameters ourselves.

2.1 Continuum fitting and line measurements

PCA is widely used to derive a representative basis from optical spectra of galaxies. When performing PCA on emission line galaxies, the eigenspectra are primarily sensitive to the variations in emission line strengths and only secondly to continuum features (Connolly et al. 1995; Yip et al. 2004). Obviously, the slope of the continuum is correlated with emission lines but the variance of the lines is bigger. To run PCA on the pure continua, one has to mask the regions of emission lines, or eliminate the lines completely by subtracting line models from the measured spectra. Line fits have to be precise enough so that the line-subtracted continua contain minimal residuals. We reprocessed the entire set of SDSS DR7 galaxy spectra according to these requirements with our own implementation of the algorithm detailed in this section.

One frequent method of fitting continua in the optical band is to express the spectrum as a non-negative linear combination of template spectra (Tremonti et al. 2004) while also accounting for the intrinsic attenuation and velocity dispersion. Although more advanced, Bayesian and PCA-based methods exist (Kauffmann et al. 2003a; Chen et al. 2012) to derive physical properties from the continuum, as we were mainly interested in the emission lines, we retained the former technique for continuum subtraction. First, we corrected for galactic extinction, masked emission lines and fitted the continuum using the templates from Bruzual & Charlot (2003) by also fitting the velocity dispersion and intrinsic extinction in parallel. Intrinsic extinction was modelled following Charlot & Fall (2000). Metallicity was taken into account by fitting four sets of templates of differing metallicities and choosing the one with minimal reduced χ^2 . Thus, the fitted metallicity can take one of four values: $Z = 0.004, 0.008, 0.02$ or 0.05 . We did not take the nebular continuum emission into account, which, in the case of young starburst galaxies, can contribute a non-negligible flux to the near-infrared part of the spectrum (Leitherer & Heckman 1995). Since the entire continuum was fitted with stellar templates only, we expect a slight overestimation of absorption lines, and therefore the overestimation of emission lines for starburst galaxies. On the other hand, within the wavelength coverage of SDSS spectroscopy, nebular continuum emission is significant only in the case of stellar pop-

ulations younger than 10 Myr or at very low metallicities of $Z \sim 0.0001$ (Mollá, García-Vargas & Bressan 2009), and only about 0.5 per cent of our sample potentially fall into this parameter range.

Due to discrepancies between continuum models and SDSS spectra (Maraston et al. 2009), the continuum-subtracted spectrum consists of three components: the emission lines, the noise and a slowly changing background that originates from the imperfect models. Since the emission lines and noise are high-frequency components, one can easily eliminate the background by a high-pass filter. For this purpose, we used a 50 \AA wide rolling median filter. This was wide enough to leave broad AGN lines almost intact, yet remove any residuals of the incorrect background subtraction. Fig. 1 illustrates this procedure.

Once the low-frequency background has been removed, lines are fitted using a technique we call *noise-limited fitting*. To precisely fit all strong emission lines, including those of active galaxies, we use three increasingly complex line models.

- A single Gaussian:

$$F(\lambda) = A \cdot e^{-\frac{(\lambda-\lambda_0)^2}{\sigma^2}}$$

- Two Gaussians centred on the same wavelength but with different variance

$$F(\lambda) = A \cdot e^{-\frac{(\lambda-\lambda_0)^2}{\sigma_a^2}} + B \cdot e^{-\frac{(\lambda-\lambda_0)^2}{\sigma_b^2}}$$

- Two Gaussians allowing for a small offset $\Delta\lambda < 5 \text{ \AA}$ between the centres, different variance

$$F(\lambda) = A \cdot e^{-\frac{(\lambda-\lambda_a)^2}{\sigma_a^2}} + B \cdot e^{-\frac{(\lambda-\lambda_b)^2}{\sigma_b^2}}$$

While the first model is enough to fit emission lines with typical velocity dispersion, the second model is necessary for lines with broad wings and the third model for asymmetric lines. Our objective is to find the simplest, yet well-fitting model. Overlapping emission lines are – obviously – fitted together, but we do not enforce any correlation on the EWs of lines from the same ion. Also, the velocity dispersions of the lines, even of those from the same ion, are fitted independently. First, we fit the lines with the simplest model, subtract it from the measurement and compare the residual within the region of the emission line with the noise in wavelength ranges without lines. If rms of the residual inside the region of the line is at least two times than elsewhere, we reject the model and attempt to fit the line with a more complex one. Fig 2 illustrates how this technique works on asymmetric broad AGN lines.

Tab. 2.1 summarizes the fitted and subtracted emission lines. Line model fit parameters are available online. In panel (a) of Fig. 3, we plot the Baldwin–Phillips–Terlevich (BPT) diagram of the sample using a colour coding we are going to use throughout the paper.

2.2 Comparison with other work

It is interesting to compare our line fits to those of Brinchmann et al. (2004). In the cited work, the authors used a simpler technique of fitting nebular emission lines of SDSS galaxies with the primary focus on the signal-to-noise ratio

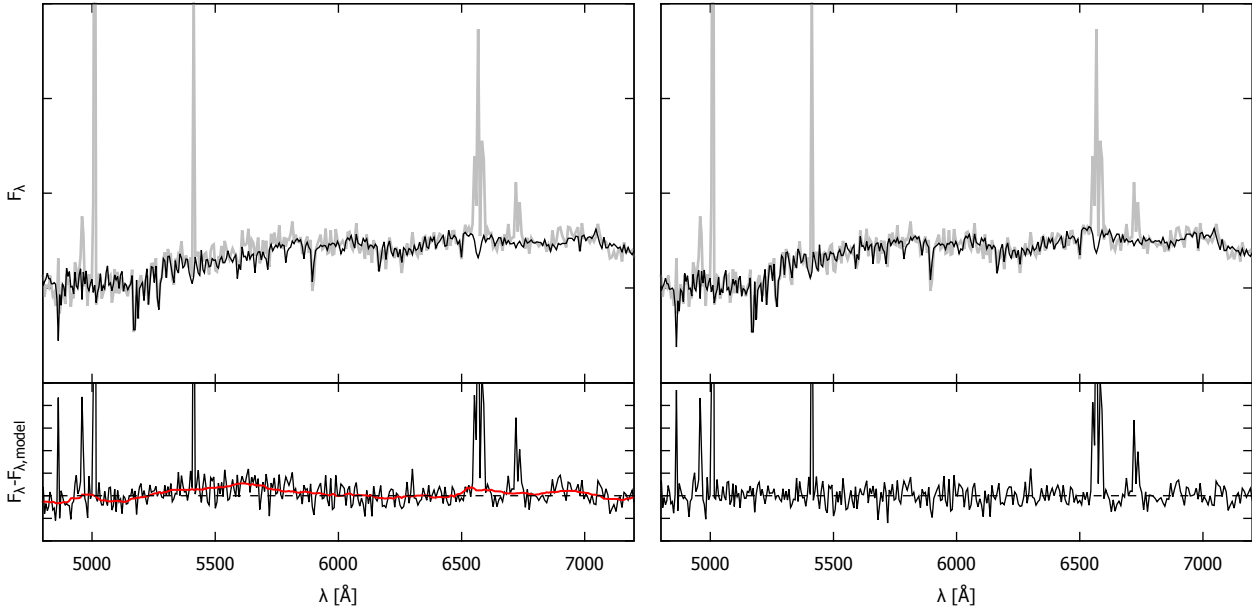


Figure 1. Illustration of fitting the stellar continuum. In the top left panel the best non-negative least square fit from 10 Bruzual–Charlot templates is plotted, the residual is visible in the bottom left panel. The effect of the low-pass filter on the residual is drawn with a red curve in the bottom left panel; we subtract this curve from the noisy residual prior to fitting emission lines. The top right panel illustrates the best-fitting continuum model, corrected for discrepancies by adding back the low-pass-filtered residual to the stellar population synthesis spectrum. The top right panel shows the high-pass-filtered residual used for fitting the lines.

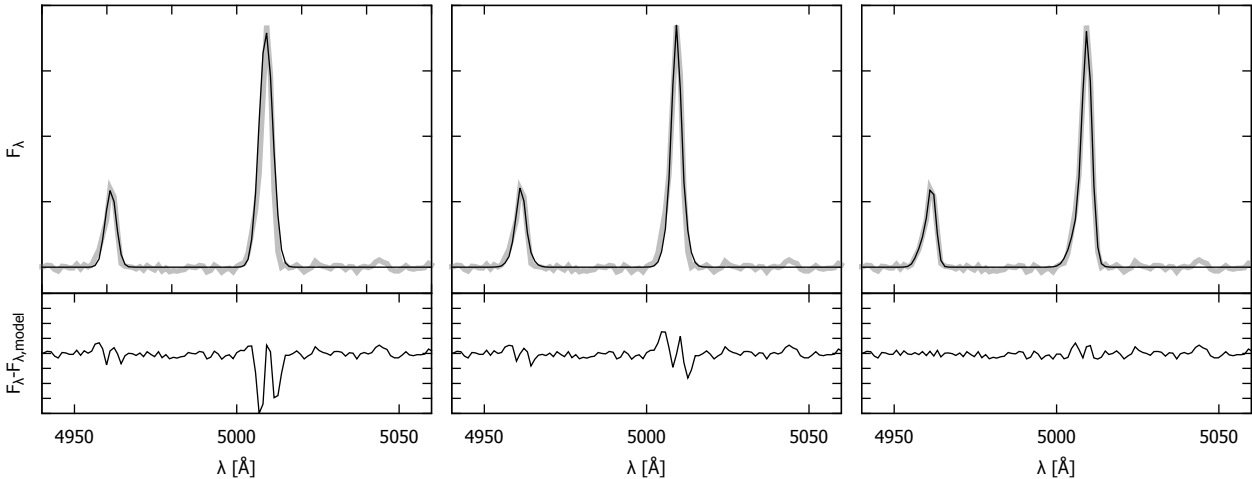


Figure 2. Illustration of noise-limited fitting of asymmetric emission lines with increasingly complex models. The top panels show the original continuum-subtracted spectrum in grey and the best-fitting models in black. The bottom panels show the residuals. The left-hand panel corresponds to a single Gaussian fit, the middle panel to two Gaussians centred on the same mean wavelength but with different variance while the right-hand panel shows the results from fitting two Gaussians with slightly different centre wavelengths. In this case, the most complex model is accepted as the line residuals are higher than the average noise for both simpler models, whereas the line residual is comparable to the average noise in the third case.

of line measurements and not on the minimization of the residuals after line subtraction. As a result, their line models cannot directly be used to get a pure continuum due to the high residuals of the fitting.

In Fig. 4, we compare the EWs of the most prominent emission lines as derived with our technique and with the method of Brinchmann et al. (2004). In the case of strong emission lines, our measurements of line strengths are very similar to the results of Brinchmann et al. (2004), but

we estimate weak emission lines significantly higher. This is very likely due to the high-pass filtering applied to the continuum-subtracted spectrum, cf. Sec. 2.1. Yet, between 97.9% and 99.1% of our line, EWs are within 3σ of Brinchmann et al. (2004), with the exception of $H\alpha$ and $H\beta$ where only 93.9% and 87.7%, respectively, of the measurements are within 3σ . Also, Brinchmann et al. (2004) measured weak lines by fitting them together with stronger lines of the same ion, imposing a constraint on line ratios, whereas we fitted

Line	λ_{vac} (Å)	Line	λ_{vac} (Å)	Line	λ_{vac} (Å)
OII	3727.09	H γ	4341.68	OI	6365.54
OII	3729.88	OIII	4364.44	NI	6529.03
H θ	3798.98	H β	4862.68	NII	6549.86
H η	3836.47	OIII	4932.60	H α	6564.61
H ζ	3890.16	OIII	4960.30	NII	6585.27
H ϵ	3971.20	OIII	5008.24	SII	6718.29
SII	4072.30	HeI	5877.65	SII	6732.67
H δ	4102.89	OI	6302.05		

Table 1. List of the fitted nebular emission lines.

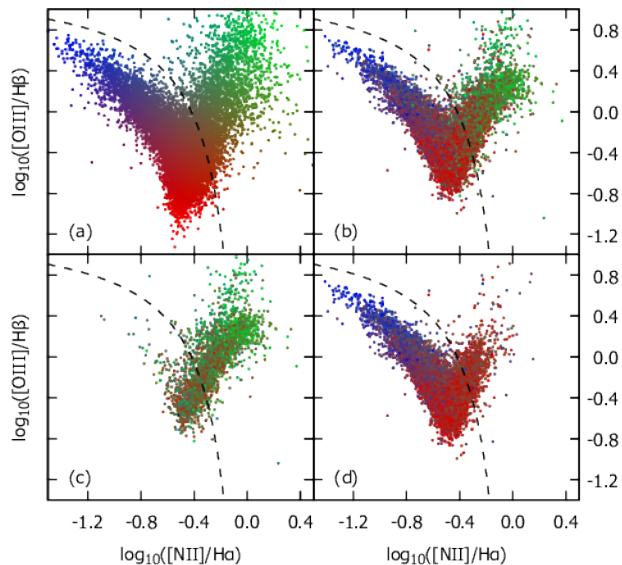


Figure 3. Original and reconstructed BPT diagrams of strong emission line galaxies sampled from SDSS DR7. In each panel, the dashed curved shows the empirical segregation line between star-forming galaxies and AGN as defined by Kauffmann et al. (2003b). Panel (a) is plotted from directly measured line EWs. Galaxies are colour coded based on their loci in the BPT plane: blue galaxies are star-forming, green ones are AGNs and red ones are the intermediate weak AGNs in the bottom corner of the distribution. This colour coding based on directly measured emission lines is used in all BPTs throughout the paper. Panel (b) displays the BPT of line log EWs reconstructed from continuum principal components using the local linear regression method with the 30 nearest neighbours in PCA space. Panels (c) and (d) show galaxies only that were originally classified as (c) AGNs, (d) star-forming using directly measured line EWs. While lines of strong AGNs and extreme starburst galaxies can be reconstructed well, there is significant ‘cross-talk’ in the quiescent region.

these lines independently. Weaker lines can easily become undetectable in noisy regions, hence our fitting method introduces some selection bias.

2.3 Galaxy sample selection

We selected a smaller sample of $N = 13788$ galaxies from the entire set of continuum and line-fitted spectra that met the following criteria:

- observed at a signal-to-noise ratio $S/N > 5$,
- all 11 emission lines listed in Tab. 2 are measured and non-zero. These lines are the same as in Györy et al. (2011).

The sample size was further limited to an easily manageable number by choosing a section of the sky (right ascension between 220° and 230°).

The requirement that all 11 emission lines should be measurable results in a sample containing galaxies with ongoing star formation or possessing an active nucleus only. Fig. 5 shows the selection effects on the distribution of the apparent and absolute r -band magnitudes, the redshift and the metallicity. While the cut in signal-to-noise ratio did produce a cutoff around $r = 19$ apparent magnitude and a relative increase of objects towards smaller redshifts, the absolute magnitude distribution shows that our selection method prefers fainter, smaller and younger galaxies. Galaxies with lower metallicities are also selected with higher probability, presumably due to correlations between ongoing strong star formation, metallicity and age. Nevertheless, galaxies with solar and above solar metallicities are still present in the sample.

The SDSS DR7 main galaxy sample, which makes up the majority of our training set, was not selected for morphological type or colour, and thus includes a wide variety of galaxies (Strauss et al. 2002). At larger redshifts, however, different environments, e.g. harder radiation fields and higher ionization parameters (Steidel et al. 2014), can lead to significantly different emission line characteristics. Certainly, the validity of our results is constrained by parameter ranges covered by the sample. By using a data set that goes beyond the types of galaxies observed by the SDSS, one can easily apply our method to a broader range of galaxies.

2.4 Continuum principal components

Principal components of the stellar continuum were derived from the fitted model spectra instead of the measurements directly. Although the precise line modelling would make it possible to subtract emission lines from the original spectra or run PCA directly on the measurements by masking out emission lines, due to the limited size of the sample which would make eigenspectra noisy, we choose to use the models instead. Fitted continuum models were taken at rest frame, convolved with the best-fitting velocity dispersion kernel and normalized to have equal flux in the following featureless rest-frame wavelength ranges: 4250–4300 Å, 4600–4800 Å, 5400–5500 Å, 5600–5800 Å. PCA was done in the 3722–6761 Å range with 0.6 Å binning. The average continuum was subtracted from the individual spectra prior to calculating the covariance matrix.

Eigenspectra were determined using the Lanczos singular value decomposition (SVD) algorithm from PROPACK (Larsen 2005). The algorithm calculates only a given number of singular vectors with the largest corresponding singular values. This was very useful in our case as the spectra consisted of 5065 data points, whereas we were interested in the first five principal components only.

The average spectrum and the resulting eigenspectra are plotted in Fig. 6. As the average was subtracted, the first eigenspectrum corresponds to the colour of the galaxy. The following two basis vectors are rather similar at first sight but the third one shows more prominent absorption lines. They are together very likely to determine the age and metallicity of the galaxy as the 4000 Å-break is very strong in both of them. The fourth vector probably corre-

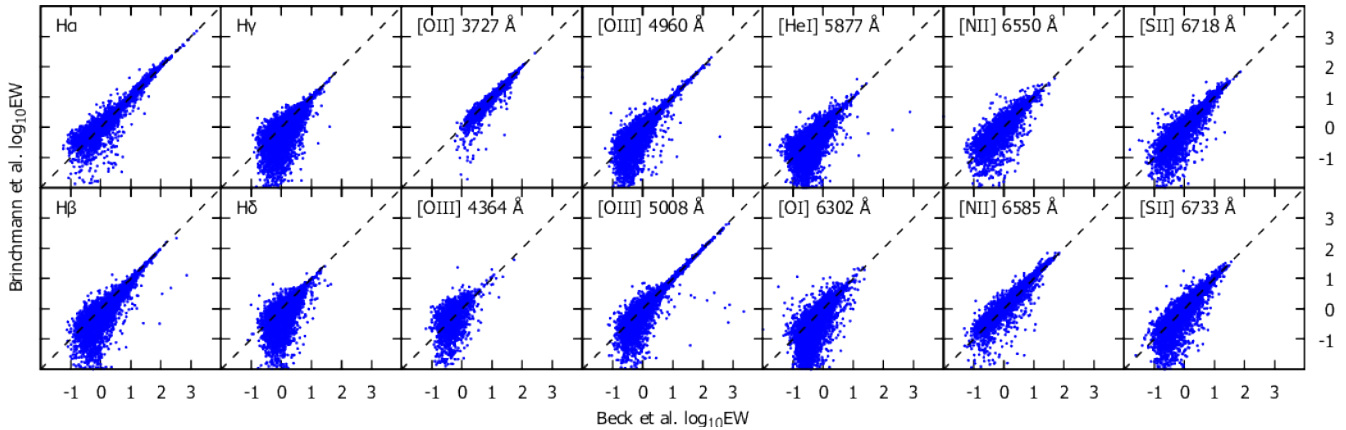


Figure 4. Comparison of emission line equivalent width of Brinchmann et al. (2004) (y -axis) with ours (x -axis). EWs are expressed in angstroms; scales on both axes are the same. Density plots are normalized for EW bins in such a way that stronger lines are also visible. Our estimate on weak lines is systematically higher, which appears more pronounced due to the log scale, but is not that significant compared to errors.

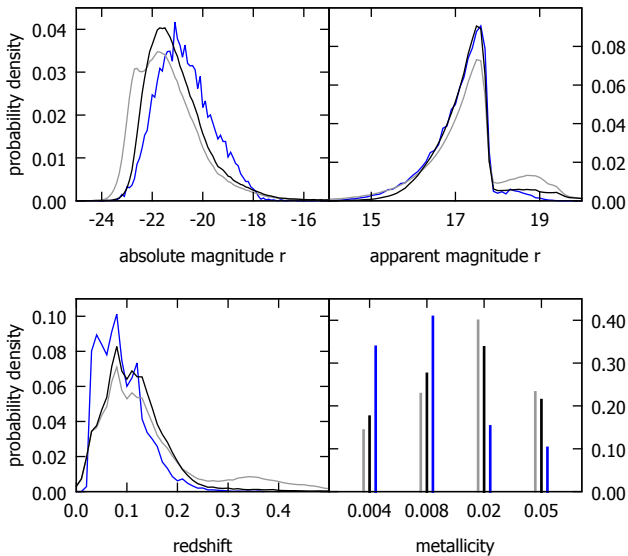


Figure 5. Normalized histograms of galaxy properties. Our sample is plotted in blue, the grey lines correspond to the entire DR7 spectroscopic galaxy sample, while the black lines show the DR7 sample excluding the deeper LRG sub-sample. The latter provides a better comparison to our data set, since we selected predominantly from the main galaxy sample.

sponds to the width of absorption lines thus correlates with velocity dispersion. The magnitude of the fourth and fifth eigenvalues is similar, and they already mark the start of the plateau in the distribution of eigenvalues, therefore taking more eigenspectra into account does not significantly increase the variance explained by them.

2.5 Emission line principal components

In contrast to what was done by Gyóry et al. (2011), we calculate principal components of the *logarithm* of emission line EWs. Fig. 7 shows the resulting singular vectors. Taking the logarithm is more useful when one is interested in

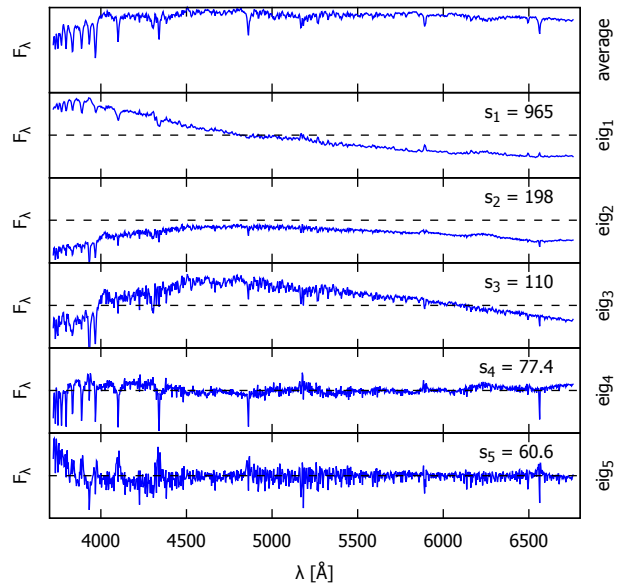


Figure 6. The average and the first five eigenvectors of the principal component analysis of galaxy continua, ordered by the corresponding singular values (as displayed in each panel). See the text for the physical interpretation of the eigenspectra.

line ratios instead of absolute line strengths and uses linear methods for the analysis. We have to mention, however, that using the logarithm of the EWs also means that the results presented in the rest of the paper will be valid in the *logarithmic sense* only.

3 RECONSTRUCTING EMISSION LINES

Our goal was to empirically estimate emission line EWs from continuum principal components. If there exists any correlation between the continuum and emission lines of galaxy

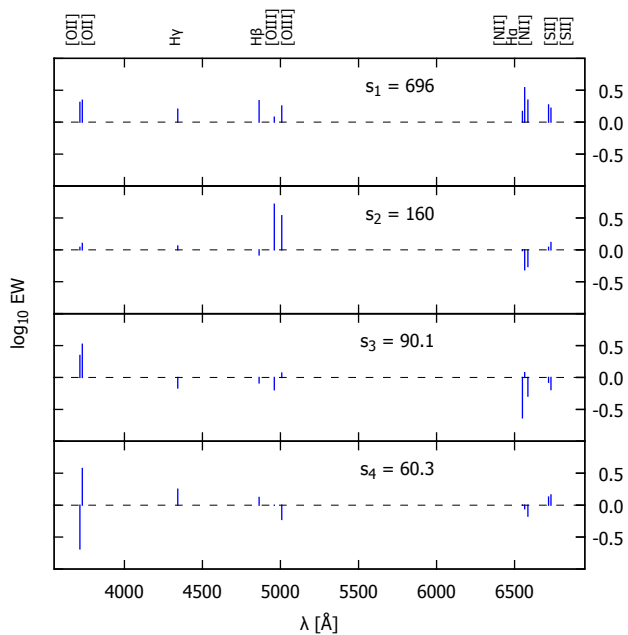


Figure 7. The first four singular vectors of the correlation matrix of the logarithm of emission line equivalent widths, ordered decreasingly by the corresponding singular values (as displayed in each panel). The fourth vector is very likely to be just noise as [OII] lines should not have different signs.

spectra, it is clearly non-linear. Global linear methods to analyse the correlations are not useful in this case, yet *locally linear* methods still can be used.

3.1 Local linear regression

Let us consider an ensemble of measurements where measured values are split into two sets $D = \{\mathbf{d}_i\}$ and $R = \{r_i\}$, i indexing the individual measurements. For the sake of simplicity, r_i are taken to be scalars whereas \mathbf{d}_i are vectors, thus D forms a metric space of dimension N . The Euclidean metric is often used to measure distances among data vectors of D even though it might lack any physical interpretation. Our objective is to characterize known, or predict unknown r_i from the always known \mathbf{d}_i vectors. To estimate r_i from \mathbf{d}_i , first we find the k -nearest neighbours of \mathbf{d}_i in D . Let us denote the set of indices of these nearest neighbours with $NN(\mathbf{d}_i, D, k)$, where $i \notin NN$ by definition. Then we express r_i in the following form

$$r_i \approx c_i + \mathbf{a}_i \mathbf{d}_i. \quad (1)$$

Note, that both \mathbf{a}_i and \mathbf{d}_i are vectors and their dot product is taken in the formula above. The c_i constants and the \mathbf{a}_i coefficients need to be determined individually for every (\mathbf{d}_i, r_i) measurement using standard linear regression by minimizing

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

where i is still the index of the measurement, j runs on the nearest neighbours and w_j is a weight. The expression of χ^2

is similar if r_i are vectors instead of scalars but the \mathbf{a} coefficients become matrices. Errors in r_i and the components of \mathbf{d}_i can be incorporated into the value of w_i . Similarly, neighbours in NN can be ordered by distance from \mathbf{d}_i and the inverse of (the square of) the distance can be used as a weight in Eq. 2.

Local linear regression has many advantages over global non-linear modelling. First of all, global models are usually either too simple to describe the data or prone to overfitting. Local linear models, on the other hand, are simple and can be used to characterize the local estimation errors. For instance, one can measure the goodness of the estimation of r_i by the χ^2 of the local linear fit. The challenge in local linear fitting is to find the k -nearest neighbours quickly in large data sets. Spatial indexing, most often a kD -tree index is used for this purpose (Csabai et al. 2007).

3.2 Emission line reconstruction from the continuum

We applied the local linear regression technique to estimate emission line EWs from the stellar continuum principal components. To test whether continuum PCs carry more information regarding the emission lines than broad-band SDSS magnitudes, we will also perform the regression analysis directly on the photometric magnitudes in Sec. 3.3. Further tests are done with randomized samples (Sec. 3.5) to get a picture of the performance of our method.

By using the notation of Sec. 3.1, \mathbf{d}_i became the first five continuum principal components and r_i became the log EWs. Emission line log EWs were fitted individually based on the log EWs of the $k = 30$ nearest neighbour galaxies in the continuum PCA space. The χ^2 of the fitting was weighted by the inverse-square distance of the neighbours from the query point. The value of $k = 30$ was chosen as a rule of thumb: we are fitting $5 + 1$ parameters and the number of data points must be large enough to adequately determine that many parameters but small enough to preserve locality. Modifying this parameter within reasonable limits (e.g. 25 – 40) does not significantly impact the results.

Fig. 8 shows the reconstructed log EWs of emission lines as functions of the directly measured EWs. EWs reconstructed from the continuum are in reasonably good agreement with directly measured log EWs. The relative flux error σ_r of the line reconstruction is Gaussian but a systematic shift δ is visible in the case of [OII], [OIII] and [NII] ($\delta \approx 0.1, 0.15, 0.15$, respectively). The typical value of the relative error is $\sigma_r \approx 0.3$ for hydrogen and sulfur, $\sigma_r \approx 35\%$ for [OII] and [NII], and $\sigma_r \approx 45\%$ for [OIII]. Col. 3-4 of Tab. 2 list the outcome of the correlation analysis for the 11 investigated lines using the local linear regression technique. Pearson’s product-moment correlation coefficient ρ and the rms error σ were calculated for each line. These numbers also show that fits are most accurate for the hydrogen and sulfur lines ($\rho > 0.8$) whereas oxygen and nitrogen lines are significantly less correlated with direct line measurements.

Fig. 10 shows the dependence of the error of the reconstruction on galaxy properties (cf. Fig. 5 for histograms of these) for select lines. The brighter and higher metallicity galaxies have a larger fraction of AGNs and are estimated with higher errors, especially in the case of oxygen and sulfur lines. Objects at higher redshifts generally exhibit de-

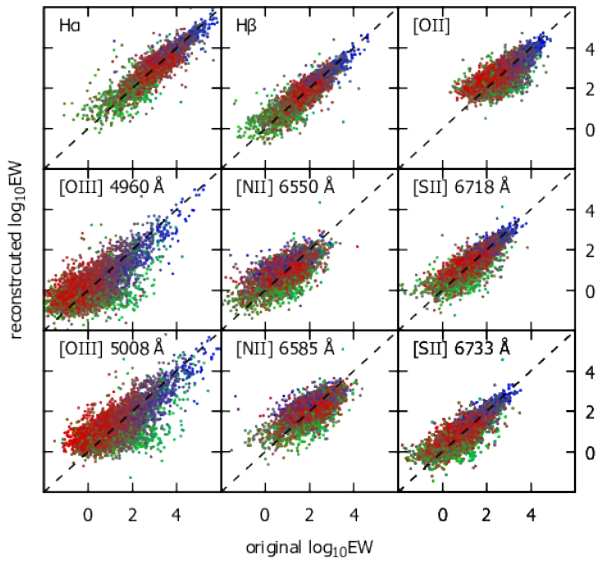


Figure 8. Reconstructed log EWs from continuum principal components. Estimated log EWs are plotted as functions of the directly measured log EWs for the 11 emission lines we used. Colour coding of data points is the same as in panel (a) of Fig. 3 and reflects the activity class of galaxies.

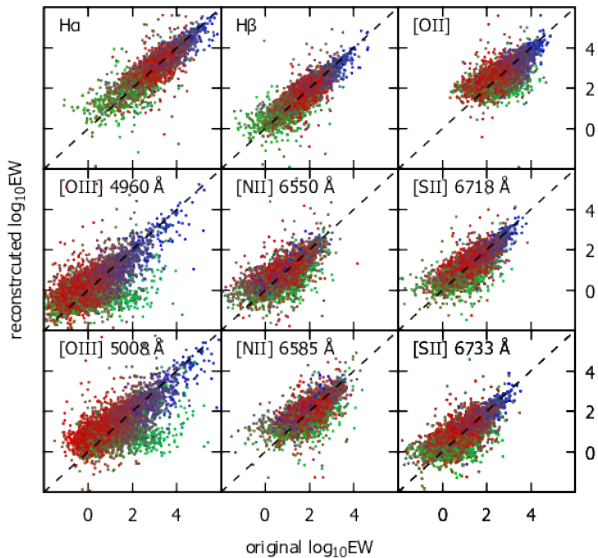


Figure 9. Reconstructed log EWs from broad-band SDSS magnitudes. Estimated log EWs are plotted as functions of the directly measured log EWs for the 11 emission lines we used. Colour coding of data points is the same as in panel (a) of Fig. 3 and reflects the activity class of galaxies.

creasing accuracy. The error of line reconstruction visibly increases towards the limits of our training set. This is due to the fact that near the edges of the training set there are fewer galaxies and the nearest neighbours used to estimate the emission lines are generally less similar to each other and to the galaxy whose lines are being fitted.

As we expected, emission lines can be much better reconstructed from the continuum of star-forming galaxies due to the strong connection between the young stellar popula-

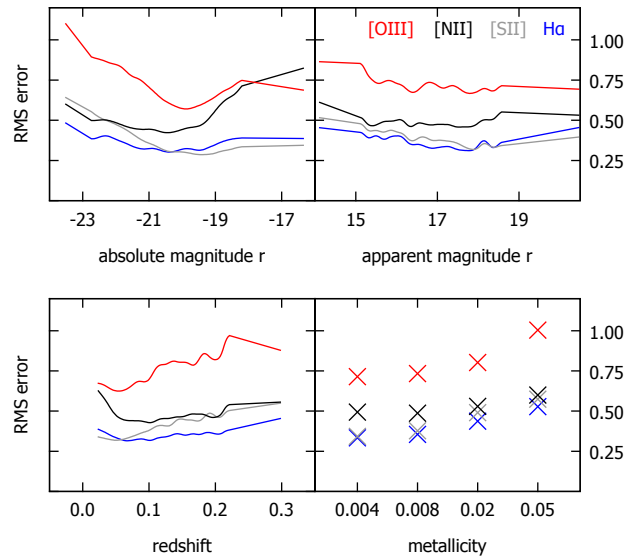


Figure 10. The rms error of emission line log EW reconstruction as the function of various galaxy properties. The colours correspond to the following emission lines: red – [OIII] 5008 Å, black – [NII] 6585 Å, grey – [SII] 6718 Å, and blue – H α . See the text for a discussion.

tion and the ISM: young massive stars are responsible for the excitation of interstellar gas clouds. Nevertheless, [OII] and [NII] lines show a significant scatter even in the star-forming case. Interestingly, sulfur lines can be reconstructed much better.

One intriguing result is that, while [OIII] is an important indicator of nuclear activity, its reconstruction from continuum properties in case of AGNs seems rather problematic. It is understandable as AGN activity correlates much less with the properties of the stellar populations than in the star-forming case. Yet, some connections exist as it is visible from [NII] and the hydrogen lines.

3.3 Emission line reconstruction from broad-band magnitudes

To see if continuum PCA is any better than directly estimating emission lines from broad-band magnitudes, we performed the above analysis using the SDSS photometric magnitudes instead of the principal components. For this purpose, we used dereddened model magnitudes without any K -correction. The lack of K -correction is not supposed to significantly affect the procedure as the redshift distribution of the sampled galaxies is rather sharp.

While broad-band magnitudes are strongly correlated with continuum principal components, it is still interesting to see how lines are reconstructed from them. First of all, magnitudes are highly correlated with each other, whereas PCA eliminates covariance. Also, observed magnitudes are already ‘contaminated’ with emission lines which might result in stronger correlations with EWs. Results of line reconstruction from magnitudes are plotted in Fig. 9. Compared with line reconstruction from PCA as plotted in Fig. 8, no clear difference can be seen in terms of scatter, perhaps with the exception of more outliers being visible in the photo-

line	λ [Å]	Principal components 30 NN galaxies		Magnitudes 30 NN galaxies		Principal components 30 random galaxies		Randomized 30 NN galaxies	
		ρ	σ	ρ	σ	ρ	σ	$\rho \times 10^2$	σ
H α	6565	0.898	0.388	0.842	0.481	0.803	0.561	-1.46	0.961
H β	4863	0.882	0.369	0.840	0.430	0.795	0.535	-1.52	0.854
SII	6718	0.839	0.416	0.773	0.492	0.798	0.484	-1.23	0.832
SII	6733	0.827	0.433	0.751	0.516	0.754	0.527	-1.36	0.840
H γ	4342	0.816	0.418	0.773	0.465	0.698	0.772	-2.81	0.790
OII	3727	0.749	0.498	0.700	0.547	0.556	0.716	-0.710	0.817
OIII	5008	0.743	0.784	0.673	0.884	0.673	0.877	-1.10	1.268
OIII	4960	0.721	0.773	0.659	0.858	0.628	0.890	-1.35	1.208
NII	6585	0.680	0.514	0.677	0.527	0.411	0.815	-0.318	0.757
NII	6550	0.664	0.570	0.669	0.579	0.367	0.880	0.306	0.820

Table 2. Numerical properties of the various line reconstruction methods, for all 11 emission lines. The four methods are as follows: (1) from continuum principal components, fitting the 30 nearest neighbours, (2) from broad-band magnitudes, fitting the 30 nearest neighbours, (3) from continuum principal components, but instead of using the 30 nearest neighbours we used 30 random galaxies, and (4) from continuum principal components, but with a randomized sample (as a cross-test). For each reconstruction, we calculated the Pearson product-moment correlation coefficient ρ and rms error σ . Emission lines are ordered by reconstructability using the first method.

metric case. Thus, log EWs can be reconstructed from magnitudes almost as well as from the principal components. Quantitative results are listed in Col. 5-6 of Tab. 2. We have to emphasize here that our sample contained strong emission line galaxies only, thus the strong correlation between magnitudes and log EWs exists only for our sample and cannot be generalized to all galaxies.

3.4 Non-local line reconstruction from the continuum

To test whether a single global linear model is sufficient to reproduce the lines, we repeated the procedure of estimating log EWs from the continuum principal components as described in Sec. 3.2 but instead of using the 30 nearest neighbour galaxies, we randomly selected 30 galaxies from the entire sample. Another difference was that the χ^2 of the fit was not weighted by the inverse-square of the distance from the query point to relax the effect of locality. By looking at Col. 7-8 of Tab. 2, it is somewhat surprising that correlation coefficients and rms errors of the individual lines did not get much worse. By looking at panel (c) of Fig. 11, one can clearly see, however, that the star-forming branch of the BPT diagram cannot be reconstructed this way, and the AGN sequence is also greatly distorted. The conclusion is that emission line log EWs cannot be explained by a simple, global linear relationship with continuum principal components. Thus, local fitting from nearest neighbours is necessary to reconstruct the BPT from either continuum principal components or broad-band magnitudes.

3.5 Cross-tests with randomized data

As another test, we shuffled the sample and randomly paired continuum principal components with emission line vectors of other galaxies to break the locality in the PCA space. Results are listed in Col. 9-10 of Tab. 2. It is not surprising that correlations almost entirely disappear in the shuffled case (note that ρ values are multiplied by 10^2). This supports that information about the emission lines is indeed encoded in the continuum of a galaxy spectrum.

3.6 Reconstructing the BPT

As the BPT diagram is generally used to classify emission line galaxies into star-forming and AGN, it is very informative to see how well the various methods can recover it solely from the continuum. The difference among the line estimation methods introduced above is obvious once the BPT diagram is plotted from the reconstructed log EWs, as it was done in Fig. 11. Local linear fitting of EWs using the nearest neighbours and reconstructing lines from broad-band magnitudes give similarly fair, qualitatively correct BPT diagrams while the randomization of the sample disrupts the diagram entirely.

To further analyse the properties of a reconstructed BPT diagram, we will stick to local linear regression based on the continuum principal components. In Fig. 3, we plot the original BPT for reference, the reconstructed diagram for all galaxies, and two diagrams showing AGNs and star-forming galaxies only, as classified by Kauffmann et al. (2003b).

The first thing to see in panels (c) and (d) of Fig. 3 is the mixing of weak star-forming galaxies with weak AGNs in the bottom corner of the reconstructed BPT diagram. The mixing is caused by the bad reconstructability of the [OIII] line which is most likely due to lack of a strong correlation between AGN activity level and the stellar continuum.

4 REVISITING STAR-FORMING/AGN SEPARATION

The mixed, low activity – low star formation rate region is located at the bottom corner of the [NII]/H α –[OIII]/H β BPT diagram. Empirically drawn BPT diagrams are noisy enough to smear pure star-forming galaxies and mixed star-forming/AGNs together in this part of the BPT so it is an interesting question whether it is possible to segregate galaxies into two distinct classes or not by incorporating information on the stellar continuum into classification model. By visually inspecting the projections of the 5D continuum PCA space, one can see that while AGNs and star-forming galaxies occupy different loci, they cannot be clearly separated into two disjoint sets by cuts in any principal component dimensions, nor the distribution of galaxies is bimodal.

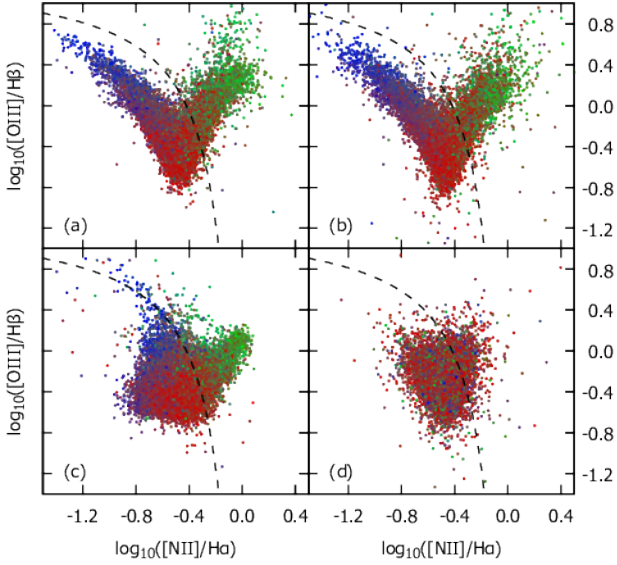


Figure 11. BPT diagrams with log EWs reconstructed from the stellar continuum using different methods and cross-tests. Panel (a) shows reconstructed log EWs from continuum principal components using local linear regression from the 30 nearest neighbours. Panel (b) is the reconstruction of lines from broad-band magnitudes by local linear regression from the 30 nearest neighbours. Panel (c) was drawn from lines estimated using the global linear regression technique from 30 randomly selected galaxies. Panel (d) is the cross-test using local linear regression but with shuffled continuum principal components. The colour coding of the data points is based on the original BPT as in panel (a) of Fig. 3. See the text for discussion.

We turned to SVM, a machine learning algorithm, to determine an empirical segregation plane between the two classes in the continuum-PCA–line-PCA space.

4.1 Support vector machines

SVM are supervised learning algorithms which can be trained to automatically classify multidimensional data vectors into two disjoint sets (Vapnik 1998; Karatzoglou, Meyer & Hornik 2006). The training phase starts with compiling a *training set* of data vectors that are tagged as either belonging to class *A* or class *B*. During learning, the model will find a hyperplane in the space of data vectors which separates the elements of *A* and *B* with the largest possible margin². Once the model is trained, it can be used to classify any query point into one of the two classes.

4.2 Automatic star-forming/AGN classification

We compiled the training set from our emission line galaxy sample by selecting galaxies on the BPT that could be clas-

² Since this is often not possible in the original space of data vectors because the distributions of the two halves of the training set are non-convex, kernel functions are used to map training set vectors into a higher dimensional space where linear segregation is possible (Schölkopf et al. 2000). Another option to handle non-convex situations is to find a *best possible* segregation plane which minimizes the overlap.

sified with high confidence either as pure star-forming or AGN. To select high-confidence AGNs only, we picked galaxies above the theoretical maximum starburst line of Kewley et al. (2001):

$$\log_{10} \left(\frac{[\text{OIII}]}{H\beta} \right) > 0.61 \left[\log_{10} \left(\frac{[\text{NII}]}{H\alpha} \right) - 0.47 \right]^{-1} + 1.19.$$

Star-forming galaxies were selected to fall below the empirical starburst line of Kauffmann et al. (2003b):

$$\log_{10} \left(\frac{[\text{OIII}]}{H\beta} \right) < 0.61 \left[\log_{10} \left(\frac{[\text{NII}]}{H\alpha} \right) - 0.05 \right]^{-1} + 1.3,$$

and at the same time be above the following line, defined by us:

$$\log_{10} \left(\frac{[\text{OIII}]}{H\beta} \right) > 3 \log_{10} \left(\frac{[\text{NII}]}{H\alpha} \right) + 1.55.$$

The line was drawn empirically to cut out the most reliably identifiable part of the star-forming population. Curves on Fig. 12 illustrate these cuts.

We used the first five continuum principal components and the first four log EW principal components of the training set galaxies as input data vectors to SVM. By combining information from the continuum into the training, we might hope a better separation of the two galaxy types in the mixed lower corner of the BPT than simply from the emission lines. As SVM is a strictly empirical model, we shall not, however, draw far-reaching theoretical conclusions from its outcome. Since our training set was not containing the mixed region, it was directly separable into two disjoint classes by a linear cut. Consequently, data points did not need to be projected into any higher dimensional space by a kernel function, like in most applications of SVM, we simply ran it on the 5 + 4 dimensional vectors of the continuum + line PCA space.

We plot the results of the SVM-based classification in Fig. 12. Panels (a) and (b) show the two training set classes (star-forming and AGN, respectively) with log EWs of the originally measured emission lines. Panels (c) and (d) show the outcome of the SVM classification. Even though the mixed region was not included in the training set at all, SVM reproduced the empirical segregation line of Kauffmann et al. (2003b) surprisingly well, with only about 6 per cent of the sample scattered into the opposite region.

5 GENERATING REALISTIC RANDOM EMISSION LINES

5.1 The stochastic recipe

Based on our findings, we propose a simple stochastic recipe to generate a realistic distribution of emission lines for stellar population synthesis models that provide the continuum only. The algorithm works by expressing the model continuum as a linear combination of the basis vectors derived from PCA of the continua of SDSS galaxies. According to these principal components, the model spectrum is classified into one of the 60 continuum classes. We used *k*-means clustering to define the continuum classes, as described in Sec. 5.3.

Let us denote the average continuum vector with $e_{0,\lambda}$ and the PCA basis vectors with $e_{i,\lambda}$, where i indexes the

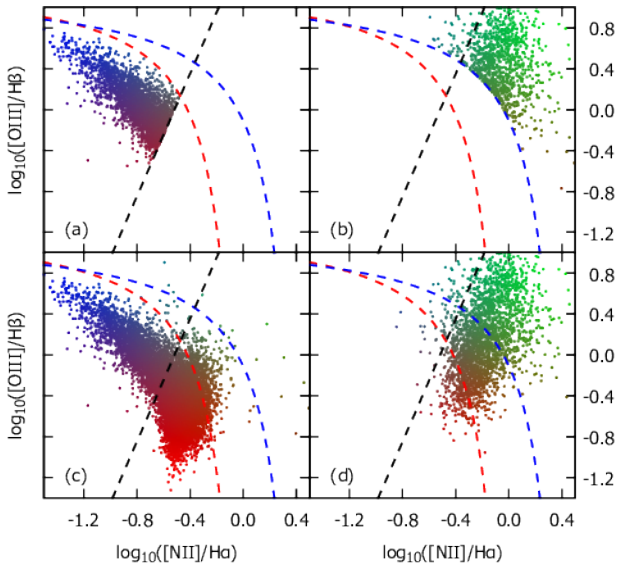


Figure 12. BPT diagrams resulting from the SVM-based star-forming/AGN separation. Panels (a) and (b) show star-forming and AGN galaxies used to train the algorithm. Panels (c) and (d) display the outcome of the automatic classification. The theoretical segregation line of Kewley et al. (2001) is drawn in blue and the empirical one of Kauffmann et al. (2003b) in red. Our star-forming/low-activity line is in black. Colour coding of the data points is based on the original BPT as in panel (a) of Fig. 3. See the text for discussion.

five dimensions of the PCA space and λ goes over the wavelength bins. Continuum classes are given by the centre of mass vectors $c_{n,i}$ where n indexes the 60 classes. Within each class, model lines are randomly generated from a multivariate Gaussian distribution. The mean line log EWs \mathbf{m}_n and the covariance matrices \mathbf{C}_n of the distributions are pre-calculated from the real galaxy sample and provided for each of the 60 continuum classes.

The detailed recipe for generating realistic emission lines given a stellar continuum model spectrum is the following.

(i) Rebin the rest-frame model spectrum s_λ to the grid of the basis vectors and normalize it as described in Sec. 2.4 to get \tilde{s}_λ .

(ii) Subtract the average continuum $e_{0,\lambda}$ from the normalized spectrum.

(iii) Express the continuum as a linear combination of the provided basis by calculating the dot products $a_i = \sum_\lambda [e_{i,\lambda} \cdot (\tilde{s}_\lambda - e_{0,\lambda})]$

(iv) Find the class centre $c_{n,i}$ in the continuum PCA space that is the closest (in Euclidean distance) to the vector a_i of the linear coefficients.

(v) Take the covariance matrix \mathbf{C}_n and mean line log EW vector \mathbf{m}_n of the line distribution within the closest class and generate a random vector of line log EWs from the corresponding multivariate Gaussian distribution.

Data necessary to generate random lines are published on the paper’s web site.

5.2 k -means clustering

k -means clustering is a machine learning algorithm that classifies data points based on their distances from cluster centres: points belonging to a cluster must be closer to the centre of mass of that particular cluster than any other clusters³. This implicit definition of a cluster makes finding the best exact solution a hard problem, but heuristic, randomized algorithms exist that can find a reasonable clustering relatively fast (Forgy 1965; MacQueen 1967; Hartigan & Wong 1979; Lloyd 1982). The only inputs of k -means clustering are the data vectors and k , the number of clusters wanted. The output is the centres of mass of the k clusters. Once the latter are known, new points can be classified simply by measuring their distances from the cluster centres and putting them into the one with the closest centre.

5.3 Automatic classification of emission line galaxies

To construct our stochastic model of emission lines, we started from the 5 + 4-dimensional vector space of continuum and log EW principal components of our high signal-to-noise ratio SDSS galaxy sample. First, we classified galaxies into continuum-log EW classes using the k -means clustering function of R and the algorithm of MacQueen (1967).

To choose the right number of clusters, one has to consider the variance of emission line log EWs as functions of the number of the clusters. The variance in each cluster is supposed to be decreasing as the number of clusters is growing since clusters are becoming smaller. The minimum variance is limited by the noise in the data. The $\sigma(k)$ curves for all emission lines are plotted in panel (a) of Fig. 13. Hence, to minimize the variance of line strengths within each class, we chose $k = 60$ as all curves get essentially flat above this value. For training sets of different sizes and characteristics, a similar analysis of the variance is advisable to determine the input parameter of k -means clustering.

5.4 Modelling the emission line distributions

Once k -means clusters in the continuum-PCA-line-PCA space are determined in the way described in Sec. 5.3, we have to model the distribution of emission line log EWs within each cluster. If the number of clusters is sufficiently high, clusters will become small enough that the distribution of emission lines within them can be well modelled by a multivariate Gaussian distribution parametrized with \mathbf{m}_n and \mathbf{C}_n . We note that a multivariate Gaussian distribution, when its entire covariance matrix is known, does not only account for individual line strengths but also for line ratios, including ratios from the same line series. It is also important to mention that, while we did the k -means classification of galaxies in the 5 + 4-dimensional continuum + log EW space, stellar population synthesis models yield the continuum coefficients only. As a result, when classifying model continua, we measure distances from cluster centres in the

³ The k -means algorithm basically constructs a Voronoi tessellation from the data vectors, with seeds being the centres of mass of the clusters.

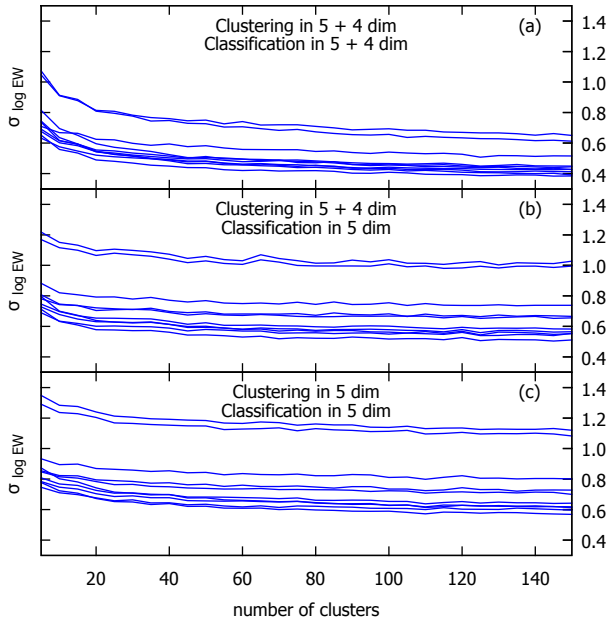


Figure 13. Variance of log EW of the 11 strong emission lines, averaged over all clusters, as a function of the number of clusters. Panel (a) refers to the case of considering both the continuum and line principal components for clustering and classification. Variance results from the sum of line measurement errors and uncertainty due to the finite size of the clusters. Panel (b) shows the effect of misclassification when only the continuum principal components are used to classify galaxies, with the clustering done in both the continuum and line log EW PCA space. Misclassification will add extra scatter to the randomly generated log EWs, cf. Sec. 5.5 and Fig. 14. Panel (c) illustrates the case of performing both the clustering and classification in the continuum principal component space only. Even with the additional variance due to misclassification, using both the continuum and line log EW principal components for the clustering is favourable.

5D continuum-PCA subspace only. This will introduce some mixing among clusters as determined by k -means and cause somewhat larger scatter in the randomly generated log EWs than what it would be based solely on the \mathbf{C}_n covariance matrices. This effect is shown in panel (b) of Fig. 13. It is still worth using the entire 5 + 4-dimensional space to run the k -mean classification because the resulting variances are still lower than using the continuum principal components only, cf. panel (c) of Fig. 13. Also, because the covariance of the lines is treated stochastically, there will be random scatter in line ratios as well.

A direct test of the algorithm is to take the galaxies of our SDSS sample, generate emission lines based on their fitted continua and see how well the BPT diagram can be reconstructed. Results of this procedure are plotted in Fig. 14 where we also show the original BPT for reference in panel (a) next to the stochastically generated BPT in panel (b). While the curve of star-forming galaxies and the AGN mixing sequence is reproduced reasonably well, there are also a large number of red data points corresponding to the bottom corner of the original BPT visible in all regions of the plot.

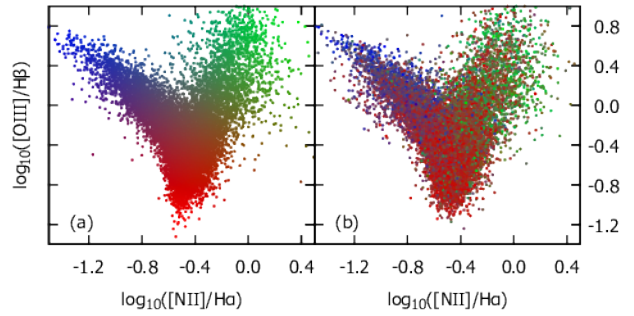


Figure 14. Panel (a) shows the original diagram plotted from the directly measured emission line for reference. Panel (b) is the BPT diagram plotted from stochastically generated emission lines based on the recipe described in Sec. 5. While lines were generated from a multivariate distribution randomly, based on the location of the continuum in the PCA space, the resulting BPT diagram resembles the original one remarkably well, although more scatter and significant mixing of galaxy types is visible. Colour coding of the data points is based on the original BPT as in panel (a) of Fig. 3.

5.5 Shortcomings of the method

The recipe outlined above yields line EWs only, and is sufficient to reproduce the flux excess caused by emission lines but not line widths. In general, line widths should be taken to be equal to the velocity dispersion, at least in the case of star-forming galaxies. Width distributions of broad lines would need to be investigated to generate AGN lines with realistic breadth distribution.

As we pointed out in Sec. 5.4, using only the continuum principal components to generate log EWs introduces additional variance due to the mixing of the classes as defined in the continuum + line space. Additionally, as lines are randomly generated based on a multivariate Gaussian distribution, EWs and line ratios are not guaranteed to be correct for individual galaxies, but will be for the entire ensemble of mock galaxies.

If the goal is to generate realistic emission lines for individual model continua, we suggest using the local linear regression method as described in Sec. 3.2. While that technique yields more accurate emission line estimates, it also requires a much larger input data set and more heavyweight algorithms.

6 SUMMARY AND FUTURE WORK

We have measured the emission lines of galaxies from SDSS DR7 to analyse the correlation between the emission lines and the stellar continua in the optical wavelength range. We have developed an algorithm, noise limited fitting, to accurately measure the parameters of broad and asymmetric emission lines, yet avoid overfitting of narrow, symmetric lines. We have also demonstrated how to correct for discrepancies between theoretical stellar continuum modelling and real measurements by low-pass filtering the residual before emission line fitting.

In Sec. 3, we have shown that optical emission line log EWs can be reasonably well reconstructed from both the optical stellar continuum and broad-band magnitudes of galaxies.

The main practical use case of our method is to generate emission lines for stellar continua from stellar population synthesis models, provided that the models fall into the wavelength and physical parameter coverage of our training set. Since our sample contained strong emission line galaxies only (with all 11 prominent lines measured), the results cannot be generalized to any type of galaxy without extending the training set, but the algorithm still applies. Also, further research is necessary to use our line reconstruction method for galaxies with fewer and weaker lines: correlations between the stellar continuum and the probability of the very presence of weak emission lines need to be taken into account.

Another application of our method is to estimate emission lines of photometric galaxies. The technique readily works for the SDSS *ugriz* filter set, but the existing training set can be adapted to other filter systems as well. While one simple way to do this is to compute synthetic photometry from the spectra, building a new training set by cross-matching the photometric measurements made with the other filter set to our spectroscopic sample is a better option (provided that the survey overlaps with the SDSS), since it would automatically account for the unknown systematics in spectrophotometry. With the outlined modifications, our technique will be of great value for analysing data from large photometric surveys like PanSTARRS and the LSST.

Additionally, by correcting for the contributions of strong emission lines to broad-band magnitudes, our method can be useful in improving template-based photometric redshift estimation algorithms to narrow the performance gap between the theoretical and the empirical approach.

In Sec. 4, we used a supervised machine learning algorithm, SVM, to verify the empirical demarcation line between star-forming galaxies and AGNs defined by Kauffmann et al. (2003b). Even though we used only extreme starburst galaxies and strong AGNs to train the algorithm, SVM yielded a result very similar to the analytical segregation curve, only about 6 per cent of galaxies in the bottom corner of the BPT diagram got classified differently. A future application to SVM would be to revisit the Seyfert/LINER separation as it was done in Kewley et al. (2006).

Finally, in Sec. 5, we gave a very simple recipe to generate random emission lines with realistic EWs on top of stellar continua generated by stellar population synthesis modes. We have demonstrated that, despite its simplicity, the method can qualitatively reconstruct the BPT. Our model has its application when the objective is not the accurate modelling of the emission lines of individual galaxies but rather generating stochastic mock catalogues with more realistic broad-band magnitudes.

ACKNOWLEDGEMENTS

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

REFERENCES

Abazajian K. N. et al., 2009, *ApJS*, 182, 543

- Atek H. et al., 2011, *ApJ*, 743, 121
 Baldwin J. A., Phillips M. M., Terlevich R., 1981, *PASP*, 93, 5
 Bressan A., Marigo P., Girardi L., Salasnich B., Dal Cero C., Rubele S., Nanni A., 2012, *MNRAS*, 427, 127
 Brinchmann J., Charlot S., White S. D. M., Tremonti C., Kauffmann G., Heckman T., Brinkmann J., 2004, *MNRAS*, 351, 1151
 Bruzual G., Charlot S., 2003, *MNRAS*, 344, 1000
 Charlot S., Fall S. M., 2000, *ApJ*, 539, 718
 Chattopadhyay T., Misra R., Chattopadhyay A. K., Naskar M., 2007, *ApJ*, 667, 1017
 Chen Y.-M. et al., 2012, *MNRAS*, 421, 314
 Connolly A. J., Szalay A. S., Bershadsky M. A., Kinney A. L., Calzetti D., 1995, *AJ*, 110, 1071
 Csabai I., Dobos L., Trencsényi M., Herczegh G., Józsa P., Purger N., Budavári T., Szalay A. S., 2007, *Astronomische Nachrichten*, 328, 852
 da Cunha E., Eminian C., Charlot S., Blaizot J., 2010, *MNRAS*, 403, 1894
 Davé R., Finlator K., Oppenheimer B. D., 2011, *MNRAS*, 416, 1354
 Eldridge J. J., Stanway E. R., 2012, *MNRAS*, 419, 479
 Ferland G. J. et al., 2013, *Revista Mexicana de Astronomía y Astrofísica*, 49, 137
 Fioc M., Rocca-Volmerange B., 1997, *A&A*, 326, 950
 Forgy E. W., 1965, *Biometrics*, 21, 768
 Gyóry Z., Szalay A. S., Budavári T., Csabai I., Charlot S., 2011, *AJ*, 141, 133
 Hartigan J. A., Wong M. A., 1979, *Applied Statistics*, 28, 100
 Ilbert O. et al., 2006, *A&A*, 457, 841
 Ivezić Ž., Connolly A., VanderPlas J., Gray A., 2014, *Statistics, Data Mining, and Machine Learning in Astronomy*. Princeton University Press, Princeton, NJ
 Jonsson P., Groves B. A., Cox T. J., 2010, *MNRAS*, 403, 17
 Karatzoglou A., Meyer D., Hornik K., 2006, *Journal of Statistical Software*, 15, 9
 Kauffmann G. et al., 2003a, *MNRAS*, 341, 33
 Kauffmann G. et al., 2003b, *MNRAS*, 346, 1055
 Kennicutt, Jr. R. C., 1998, *ARA&A*, 36, 189
 Kerekes G., Csabai I., Dobos L., Trencsényi M., 2013, *Astronomische Nachrichten*, 334, 1012
 Kewley L. J., Dopita M. A., Leitherer C., Davé R., Yuan T., Allen M., Groves B., Sutherland R., 2013, *ApJ*, 774, 100
 Kewley L. J., Dopita M. A., Sutherland R. S., Heisler C. A., Trevena J., 2001, *ApJ*, 556, 121
 Kewley L. J., Groves B., Kauffmann G., Heckman T., 2006, *MNRAS*, 372, 961
 Kovács A., Szapudi I., 2015, *MNRAS*, 448, 1305
 Larsen R. M., PROPACK - Software for large and sparse SVD calculations, <http://sun.stanford.edu/~rmunk/PROPACK/>
 Leitherer C., Ekström S., Meynet G., Schaerer D., Agienko K. B., Levesque E. M., 2014, *ApJS*, 212, 14
 Leitherer C., Heckman T. M., 1995, *ApJS*, 96, 9
 Lloyd S., 1982, *IEEE Transactions on Information Theory*, 28, 129
 MacQueen J., 1967, *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, Vol.

- 1: Statistics. Univ. California Press, Berkeley, CA, p. 281
Maraston C., Strömbäck G., 2011, MNRAS, 418, 2785
Maraston C., Strömbäck G., Thomas D., Wake D. A.,
Nichol R. C., 2009, MNRAS, 394, L107
Mollá M., García-Vargas M. L., Bressan A., 2009, MNRAS,
398, 451
Schölkopf B., Smola A. J., Williamson R. C., Bartlett P. L.,
2000, Neural Computation, 12, 1207
Silva L., Granato G. L., Bressan A., Danese L., 1998, ApJ,
509, 103
Stanway E. R., Eldridge J. J., Greis S. M. L., Davies
L. J. M., Wilkins S. M., Bremer M. N., 2014, MNRAS,
444, 3466
Stasińska G., 1984, A&AS, 55, 15
Steidel C. C. et al., 2014, ApJ, 795, 165
Strauss M. A. et al., 2002, AJ, 124, 1810
Topping M. W., Shull J. M., 2015, ApJ, 800, 97
Tremonti C. A. et al., 2004, ApJ, 613, 898
Vapnik V. N., 1998, Statistical Learning Theory. Wiley,
New York, NY
Vazdekis A., Ricciardelli E., Cenarro A. J., Rivero-
González J. G., Díaz-García L. A., Falcón-Barroso J.,
2012, MNRAS, 424, 157
Veres P., Bagoly Z., Horváth I., Mészáros A., Balázs L. G.,
2010, ApJ, 725, 1955
Wright D. E. et al., 2015, MNRAS, 449, 451
Yip C. W. et al., 2004, AJ, 128, 585

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