

Data science project

Klay Matthew Kulik*, Deepak Sharma[†] and Linxiao Wang[†]

*Department of Physics and Astronomy, [†]Department of Computer Science

Email: kkulik2@uwo.ca, dsharm55@uwo.ca, lwang739@uwo.ca

Abstract

Summarise your report.

Index Terms

Data science

I. INTRODUCTION

Motivate the problem and the general approaches you will use to solve it. Explain in enough detail so that a non-expert can understand why the subject-area problem is important and how you propose to solve it.

Introduce ugriz, and discuss the differences with RGB¹

II. DATA

A. Data Source

The data we used for this project come from the Sloan Sloan Digital Sky Survey (SDSS) Data Release (DR)7 [1] and the Galaxy Zoo (GZ) project first data release [2]. SDSS is a major multi-spectral imaging and spectroscopic redshift survey using a dedicated wide-field 2.5 m telescope located at Apache Point Observatory (APO) near Sacramento Peak in Southern New Mexico. DR7 contains five-band photometry for 357 million distinct objects [1]. The GZ project collected simple morphological classifications of nearly 900,000 galaxies drawn from SDSS DR7, contributed by hundreds of thousands of volunteers [2]. We used the galaxy photometry data from SDSS and the classification labels from GZ.

B. Pre-processing

Since the classifications in GZ are from volunteer votes, we only used the data entries that have high confidence classifications, that is, we choose galaxies with debiased probability (given in [2]) greater than 0.985 for spiral galaxies, and 0.926 for elliptical galaxies, respectively. We choose these thresholds to ensure that: (i) the galaxies used for training the neural network have highly accurate classifications; and (ii) the number of data for both classes in the training and test sets are balanced [3].

The data we got from SDSS are in Flexible Image Transport System (FITS) format. For each galaxy five FITS files are provided by SDSS, one for each band (u, g, r, i and z). Each FITS file contains a header part and a data part. We removed the header part and resized the images to 200×200 pixels. As we discussed in Section I, 'ugriz' covers a broader band than RGB images. Common ways to map 'ugriz' files to RGB images would take 3 or 4 channels (out of 5) of 'ugriz', and do a linear transformation on them, thus would definitely lose some information. In order to keep a more complete data for each galaxy we used 'ugriz' files rather than RGB images as the input to our classifier. Figure 1 shows the gray-scale images of the **u, g, r, i and z band photometry of a spiral galaxy**² as well as the RGB image of the same galaxy. We can see that there is information on each band but the RGB image only uses 3 or 4 of them which would cause information loss.

¹ Introduce ugriz

² is this the correct way to say it?

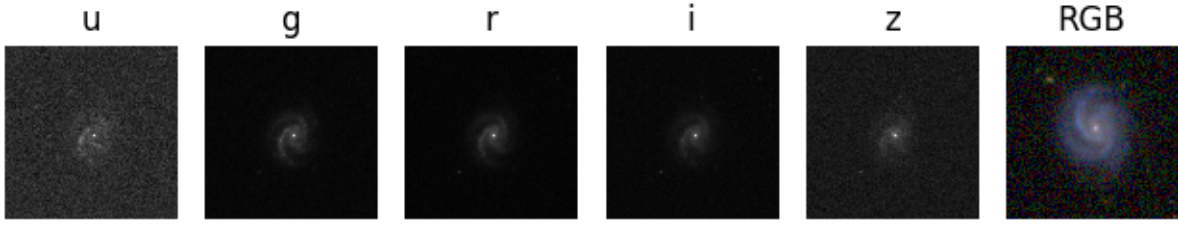


Fig. 1: The 'ugriz' images and the RGB image of the same galaxy

	Training	Validation	Testing	Total
Spiral				
Elliptical				

TABLE I: Numbers of data points in each class for training, validation and testing

Then we **split**³ the data into three sets for training, validation and testing. The number of each class for each set is shown in Table I.

III. METHOD

Describe the data science methods you applied, why you applied them, and how you applied them. Assume that your reader has similar background in data science methods as you do.

IV. RESULTS

Describe the results of applying your methods.

V. DISCUSSION

Explain the meaning of the results you obtained. E.g., did you find good/bad performance? Could your work be used in practice? What would help improve it?

REFERENCES

- [1] K. N. Abazajian, J. K. Adelman-McCarthy, M. A. Agüeros, S. S. Allam, C. A. Prieto, D. An, K. S. Anderson, S. F. Anderson, J. Annis, N. A. Bahcall *et al.*, "The seventh data release of the sloan digital sky survey," *The Astrophysical Journal Supplement Series*, vol. 182, no. 2, p. 543, 2009.
- [2] C. Lintott, K. Schawinski, S. Bamford, A. Slosar, K. Land, D. Thomas, E. Edmondson, K. Masters, R. C. Nichol, M. J. Raddick *et al.*, "Galaxy zoo 1: data release of morphological classifications for nearly 900 000 galaxies," *Monthly Notices of the Royal Astronomical Society*, vol. 410, no. 1, pp. 166–178, 2010.
- [3] A. Khan, E. Huerta, S. Wang, R. Gruendl, E. Jennings, and H. Zheng, "Deep learning at scale for the construction of galaxy catalogs in the dark energy survey," *Physics Letters B*, 2019.

³ [explain further more about how we split the data](#)