Machine Learning for photometric redshift estimation



Markus Michael Rau Ben Hoyle Kerstin Paech Stella Seitz



Overview

- What are photometric redshifts (photoZ)?
- Why are accurate photoZs important for cosmology?
- How can we quantify the photoZ error distribution?
- How can we avoid systematic errors?

The cosmos in 3D





0.0 3000 4000 5000 6000 7000 8000 Observed Wavelength (Angstroms)

9000

10000

11000



PhotoZs can be obtained for all other galaxies of the photometric survey



From images to photometry



Photometry is challenging!

- bright objects
- crowded regions
- pixel failures

From images to photometry



Photometry is challenging!

Alternative: Apply a Convolutional Neural Network directly to the image cutouts produced by traditional tools



Hoyle 2015 arXiv:1504.07255

Consistent with state-ofthe art "conventional" method

In the following we use traditional photometry as inputs.

Photometry provides incomplete information about redshift



Sample distribution from different methods



Rau et al. 2015 arXiv:1503.08215

Relative Bias $(C_{\ell}) = \frac{C_{\ell}^{\text{phot}} - C_{\ell}^{\text{spec}}}{C_{\ell}^{\text{spec}}}$

Traditional ML methods (e.g. ANNz) determine a single ,best fit' prediction

This can lead to **biased** observables **We need to estimate the error distribution!**



How can we estimate the photoZ error distribution?



Basic idea: View regression as a classification problem. The probabilities for class membership determine the height of the redshift bin.

Problem: Computationally and storage expensive

Alternative: The Highest Weight Element (HWE)



Idea: Find nearest neighbor in color-magnitude space

Remember: We are primarily interested in the sample PDF

The HWE is extremely **fast** to compute and extremely **storage efficient**.

Bandwidth selection



Density estimate evaluated on spectroscopic calibration data!

Large samples difficult to obtain (here 17000 objects)

Bandwidth selection with Silvermans "rule of thumb" bw = 0.031Oversmoothed = 1.2*bw Undersmoothed = 0.8*bw Bandwidth selection alone can create a bias at 1-2% level

Higher if lensing weights are included

Sample selection bias

apply magnitude cuts to the calibration sample



How do we deal with this problem?



Remove galaxies with incomplete spectral calibration

We loose data



Augment training set: put galaxies artificially at a higher redshift

> Model guided Extrapolation

Data augmentation



Summary

- Inaccurate photoZ can bias cosmological observables
- We can accurately and efficiently estimate redshift distributions with representative calibration data
- We are challenged by various sources of systematic error:
 - Missing spectra for faint objects
 - Bandwidth selection
 - Variations between fields

These effects are **not understood** well enough!



