

# Application of Machine Learning Methods on Astronomical Databases

Apostolos Kiraleos

## Abstract

Galaxy redshift is a crucial parameter in astronomy that provides information on the distance, age, and evolution of galaxies. This dissertation investigates the application of machine learning for predicting galaxy redshifts. It involves the development and training of a neural network to analyze galaxy spectra sourced from the European Space Agency’s Gaia mission and showcases the practical implementation of machine learning in astronomy.

## Contents

Chapter 1: Introduction . . . . .	2
1.1 Motivation and Objectives . . . . .	2
1.2 Gaia space obvservatory . . . . .	2
1.3 Galaxy redshift . . . . .	3
1.4 Convolutional Neural Networks . . . . .	3

# Chapter 1: Introduction

## 1.1 Motivation and Objectives

Galaxy redshift estimation is pivotal in astronomy, enabling us to grasp the universe's structure and evolution. It has broad applications, including cosmology, galaxy evolution, large-scale structure mapping, and exoplanet research. However, it can be resource-intensive and error-prone.

This dissertation is driven by the practical need for precise galaxy redshift estimation in astronomy. The primary objective is to train a Convolutional Neural Network (CNN) using pre-processed galaxy spectra data.

The main focus of this dissertation is the training of the neural network. This involves developing and training a CNN model using a high-quality dataset of galaxy spectra. The process includes data preparation, architecture design, and optimizing training parameters to ensure the best performance.

Our aim is to evaluate the accuracy of the trained CNN in predicting galaxy redshifts. We will conduct a comprehensive evaluation using appropriate metrics and techniques to assess the model's performance, and also compare it against Gaia mission's Unresolved Galaxy Classifier (UGC) which is another machine learning approach to predicting galaxy redshifts using support vector machines (SVMs) instead of neural networks.

## 1.2 Gaia space observatory

The Gaia mission is a European Space Agency (ESA) space observatory that has been in operation since 2013. The primary goal of the mission is to create a three-dimensional map of our galaxy by measuring the positions, distances, and motions of over two billion stars.

At its core, Gaia is equipped with two optical telescopes accompanied by three scientific instruments, which collaborate to accurately ascertain the positions and velocities of stars. Additionally, these instruments disperse the starlight into spectra, facilitating detailed analysis.

Throughout its mission, the spacecraft executes a deliberate rotation, systematically scanning the entire celestial sphere with its two telescopes. As the detectors continuously record the positions of celestial objects, they also capture the objects' movements within the galaxy, along with any alterations therein.

Over the duration of its mission, Gaia conducts approximately 14 observations annually for each of its designated stars. Its primary objectives include the precise mapping of stellar positions, distances, motion patterns, and variations in luminosity. Gaia's mission anticipates the revelation of an extensive array of novel celestial entities, encompassing exoplanets and brown dwarfs, alongside the thorough examination of hundreds of thousands of asteroids located within our own Solar System. Furthermore, the mission encompasses an investigation into more than 1 million distant quasars while subjecting Albert Einstein's General Theory of Relativity to rigorous new assessments.

### 1.3 Galaxy redshift

Galaxy redshift is a fundamental astronomical property that describes the relative motion of a galaxy with respect to Earth. The redshift of a galaxy is measured by analyzing the spectrum of light emitted by the galaxy, which appears to be shifted towards longer wavelengths due to the Doppler effect. Redshift is a crucial parameter in astronomy, as it provides information about the distance, velocity, and evolution of galaxies.

The redshift of a galaxy is measured in units of “ $z$ ”, which is defined as the fractional shift in the wavelength of light emitted by the galaxy. Specifically, the redshift “ $z$ ” is defined as:

$$z = \frac{\lambda_{obsv} - \lambda_{emit}}{\lambda_{emit}}$$

where  $\lambda_{obsv}$  is the observed wavelength of light from the galaxy, and  $\lambda_{emit}$  is the wavelength of that same light as emitted by the galaxy. A redshift of  $z = 0$  corresponds to no shift in the wavelength (i.e., the observed and emitted wavelengths are the same), while a redshift of  $z = 1$  corresponds to a shift of 100% in the wavelength (i.e., the observed wavelength is twice as long as the emitted wavelength).

Accurate and efficient estimation of galaxy redshift is essential for a wide range of astronomical studies, including galaxy formation and evolution, large-scale structure of the universe, and dark matter distribution. However, measuring galaxy redshifts can be a challenging task due to various factors such as observational noise, instrumental effects, and variations in galaxy spectra.

### 1.4 Convolutional Neural Networks

Convolutional neural networks (CNNs) are a type of artificial neural network (ANN) that has proven to be highly effective for tasks involving image and video analysis, such as object detection, segmentation, and classification. CNNs are inspired by the structure and function of the visual cortex in animals, which contains specialized cells called neurons that are tuned to detect specific visual features. In a similar way, CNNs are designed to learn and extract meaningful visual features from raw image data.

At the core of a CNN are individual processing units called neurons, which are organized into layers. Each neuron receives input from other neurons in the previous layer, applies a mathematical operation to that input, and passes the output to the next layer. The output of the final layer of neurons is the predicted output of the network for a given input.

The key innovation of CNNs is their use of convolutional layers, which enable the network to automatically learn and extract local spatial features from raw input data. In a convolutional layer, each neuron is connected only to a small, localized region of the input data, known as the receptive field. By sharing weights across all neurons within a receptive field, the network can efficiently learn to detect local patterns and features, regardless of their location within the input image.

CNNs typically also include pooling layers, which downsample the output of the previous layer by taking the maximum or average value within small local regions.

This helps to reduce the dimensionality of the input and extract higher-level features from the local features learned in the previous convolutional layer.

The final layers of a CNN are fully connected layers, which take the outputs of the previous convolutional and pooling layers and use them to make a prediction. In the case of image classification, for example, the output of the final fully connected layer might be a vector of probabilities indicating the likelihood of each possible class. In our case, instead of class probabilities, the output of the final fully connected layer will yield a single numeric value representing the predicted redshift of the observed galaxy.