

Galaxy Morphological Classification with Big Data: An Analysis of Large Spectral and Photometric Datasets*

Extended Abstract

Son Pham[†]

CU Boulder

Boulder, CO

son.pham-2@colorado.edu

ABSTRACT

CCS CONCEPTS

• **Mathematics of computing** → *Mathematical analysis*; • **Information systems** → *Data management systems*; • **Computing methodologies** → *Machine learning*; *Distributed computing methodologies*; • **Applied computing** → *Astronomy*.

KEYWORDS

ACM proceedings

ACM Reference Format:

Son Pham. 2022. Galaxy Morphological Classification with Big Data: An Analysis of Large Spectral and Photometric Datasets: Extended Abstract. In *Proceedings of CSCI 6502 - Big Data Analytics (Spring '22)*. ACM, New York, NY, USA, 4 pages. https://doi.org/10.475/123_4

*Produces the permission block, and copyright information

[†]Note

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).
Spring '22, March 2022, Boulder, CO, USA

© 2022 Copyright held by the owner/author(s).

ACM ISBN 123-4567-24-567/08/06.

https://doi.org/10.475/123_4

1 INTRODUCTION

Big data analytics is transforming the way scientists and astronomers study the universe. With each large scale telescope on the ground or in space, dozens of gigabytes of data is being generated per day - an amount that will take an increasingly amount of time to explore and process. For example, the Hubble Space Telescope generates about 120G of scientific data every week [2]. Sky surveys conducted by telescopes such as the Sloan Digital Sky Survey (SDSS) produces data releases each year that can be as high as 100's of TBs. SDSS provides a wide range of data types, such as optical spectra, infrared spectra, and imaging.

Recent discoveries, such as the presence of over 100 black holes in the center of our Milky Way, was realized using data from decades ago generated by the Chandra satellite. Scientific and technological advancement in combination led the way to this capability. With the increase in computational performance, astronomers now have the capabilities to explore these large datasets without the need to invest in large ground-based optical telescopes or work in a research lab. As big data analytics continue to grow, these discoveries will become more common as scientists continue to collect and process more of the data that is available. As technological advancements grow, they will have more readily available tools and platforms for their analysis. With the cummilation telescopes and technologies present today, the entire electromagnetic spectrum can be observed within a patch of sky of interest.

SDSS datasets, such provide many different data types, provides the opportunity to explore galaxies, stars, and

Parameter	Description	Unit
Total unique area covered	14,555	sq. deg
Total area of imaging (including overlaps)	31,637	sq. deg
Individual image field size	1361x2048 (0.0337)	pixels (sq. deg)
catalog objects	1,231,051,050	(-)
unique detections	932,891,133	(-)
Median PSF FWHM, r-band	1.3	arcsec
Pixel scale	0.396	arcsec

quasars too distant to view for amateur astronomers and conduct many different types of analysis on them. The analytics of interest to the authors include the use of optical imaging and spectral data. Legacy imaging generated from prior SDSS programs, if used with machine learning classification techniques, can provide automatic classification of morphological properties of these galaxies. Quantities measured from these images and spectra readings can also provide information such as magnitudes, redshifts, and object classifications. In particular, redshift can be a good indicator of galaxy morphology.

This paper will go over in detail the approach to perform morphological classification of galaxies using these available information from various sources. The main source of interest for the author is SDSS's data, which includes 100's of terabytes of data covering more than one-third of the entire celestial sphere [1].

2 LITERATURE SURVEY

2.1 External Studies

2.2 Studies Based on This Dataset

3 DATASETS

SDSS Imaging dataset.

4 APPROACH

4.1 Data Cleaning and Pre-Processing

4.2 Data Integration

There are many available sources that provide celestial sky survey data. In astronomy, many use a common file transfer type called Flexible Image Transport System (FITS), mainly used for image file transfer, but can also be used to store any type of data that can fit in a multi-dimensional array or table. Typical sky surveys contain equatorial coordinates, defined relative to the celestial sphere, that can be used to identify position of an object in the sky. These surveys also contain object ID's that is used to reference a particular object within a catalog or database. These object ID's are then used as keys to integrate various databases together.

To link the many objects from various data sources, the common method is to obtain the right ascension (RA) and declination (DEC) of the object in the celestial frame. RA is measured from the vernal equinox to the point of interest, going positive east along the celestial equator. The vernal equinox is the intersection of the celestial equator and the ecliptic where the ecliptic crosses the ascending node. The RA varies from 0 to 24 hours, and typically measured in hours-minutes-seconds (HMS). To get RA in decimal degrees, Equation 1 is used.

DEC is the angle from the celestial equator to the point of interest, positive going north (and negative going south). DEC varies from -90 to +90 degrees, with fractions of a degree measured in arcminutes and arcseconds. The notation for DEC is degrees-minutes-seconds, with a full circle providing 360 degrees, each degree having 60 arcminutes and each arcminute having 60 arcseconds. DEC in decimal degrees can be obtained using Equation 2. Some surveys, such as the SuperCOSMOS survey, already has their data in this format, but sources such as the Australia Telescope 20-GHz Survey contains coordinates in the standard form and must be converted.

$$RA_{deg} = 15(hours + \frac{minutes}{60} + \frac{seconds}{3600}) \quad (1)$$

$$DEC_{deg} = degrees + sign(degrees) \left(\frac{arcminutes}{60} + \frac{arcseconds}{3600} \right) \quad (2)$$

Once all sources format is converted in to a common decimal degrees format, the datasets can be read and integrated into a common database. To connect objects from different surveys together, we must look at the coordinates of the objects within each catalog and cross-match to find its closest counterpart, measured by the angular distance between them. Since the position of two given objects are measured as points on a sphere, the great-circle distance must be computed. This can be achieved using the haversine formula, which computes the distance of two points using a set of right ascension and declination angles, as shown in 3, where α and δ are the corresponding RA and DEC of the two points.

$$d_{\theta} = 2 \arcsin \cdot \sqrt{\sin^2 \frac{|\delta_1 - \delta_2|}{2} + \cos \delta_1 \cos \delta_2 \sin^2 \frac{|\alpha_1 - \alpha_2|}{2}} \quad (3)$$

Although this computation is useful in helping integrate multiple sources together, the algorithm does not scale well. Since the algorithm requires the calculation of the distance of each data object from one catalog to each data object in the second catalog, the time complexity is $O(n \cdot m)$, where n and m are the counts of data objects in the catalogs. This is not an issue for catalogs with a few hundred samples, but typical catalogs such as the SDSS contains over a million data points, which will be prohibitively expensive to compute.

An improved algorithm was used for this project's analysis. The first method was to use Python's numpy array structure to automatically loop through a catalog dataset. This provided a speed boost compared to standard Python standard library, but slowed down significantly when working with datasets greater than 10 million data objects. The second method was to constrain the search to within a given angular radius away from the desired celestial location, which will disregard points located outside of this radius. To find this index value within a catalog, a binary search is performed, where the catalog is repeatedly split in half until the value of interest is found. This provided a boost in performance for small and medium sized datasets, but slows down significantly for datasets with over 100

million data objects. The final method was to use a k-d tree to perform the crossmatching. Similar to the binary search algorithm, the k-d tree divides the k-dimensional space into two parts recursively until each segment is it's own leaf.

4.3 Statistical Analysis

One of the major issues in working with large datasets is that the standard analytical evaluation of statistical formulas requires an entire batch of data samples to be stored in memory for computation. Having 100's of thousands of samples can quickly exceed the memory limit of an average laptop or computer. This is especially the case when running analysis on embedded systems where memory is a driving constraint of performance. Stream processing allows a user to run calculations as the data comes in and releases memory that held samples of data from previous use.

Welford's method for statistical analysis allows the use of stream processing to compute some statistical parameters [3]. This method gives accurate estimates of the mean and variance without having to store all the data in memory. The standard process for computing standard deviation is to compute the mean of the data in one pass, then calculate the square deviation of values from the mean in the second pass. In crude methods of numerically computing deviation and means, one can compute the same standard deviation in one pass. Equation 4 shows this method by accumulating the sums of x_i and x_i^2 . This subtraction can result in loss of accuracy if the square of the mean is large while the variance is small.

$$\sigma = \sqrt{\left(n \sum_{1 \leq i \leq n} x_i^2 - \left(\sum_{1 \leq i \leq n} x_i \right)^2 \right) / n(n-1)} \quad (4)$$

Welford's method simply keeps a running sum of the data, number of samples, and deviation from data collected so far, and the user can view the sample variance and mean at anytime during the computational process to view their progression. Equation 5 show the recurrence formula from Welford which takes into account only the current and previous sample values. Equation 6 shows the computation of the mean, and Equation 7 for the standard deviation.

$$M_{2,n} = M_{2,n-1} + (x_n - \bar{x}_{n-1})(x_n - \bar{x}_n) \quad (5)$$

$$\sigma_n^2 = \frac{M_{2,n}}{n}$$

$$s_n^2 = \frac{M_{2,n}}{n-1}$$

4.4 Imaging

4.5 Spectroscopy

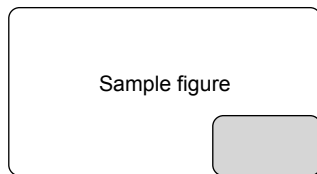


Figure 1: Sample figure

5 EVALUATION

(6) 5.1 Statistical Analysis

5.2 Imaging

(7) 5.3 Spectroscopy

6 RESULTS

7 APPLICATIONS

8 CONCLUSION

REFERENCES

- [1] Sloan Digital Sky Survey. 2022. Data Release 17 Scope. (2022). <https://www.sdss.org/dr17/scope/>
- [2] Nola Taylor Tillman. 2022. Hubble Space Telescope: Pictures, facts and history. (2022).
- [3] B. P. Welford. 1962. Note on a Method for Calculating Corrected Sums of Squares and Products. *J. ACM* 4, 3 (Aug. 1962). <https://www.jstor.org/stable/1266577>