

Various project. 3.1.  
object classification  
classification is an



Various data mining algorithms are used by astronomers in most of the applications in astronomy. However, studies and several projects have also been made by data mining experts utilizing astronomical data because astronomy has produced many large datasets that are flexible to the approach along with other fields such as medicine and high energy physics. Examples of such projects are the SKICAT-Sky Image Cataloging and Analysis System for catalog production and catalog analysis from digitized sky surveys particularly the scans of the second Palomar Observatory Sky Survey; the JAR Tool- Jet Propulsion Laboratory Adaptive Recognition Tool used for recognition of volcanoes in the over 30, 000 images of Venus returned by the Magellan mission; the subsequent and more general Diamond Eye and the Lawrence Livermore National Laboratory Sapphire project. 3. 1.

Object classification is an important preliminary step in the scientific process as it provides a method for organizing information in a way that can be used to make hypotheses and compare with models. The two useful concepts in object classification are the completeness and the efficiency, also known as recall and precision. They are defined in terms of true and false positives (TP and FP) and true and false negatives (TN and FN).

The completeness is the fraction of objects that are truly of a given type that are classified as that type: and the efficiency is the fraction of objects classified as a given type that are truly of that type. These two quantities are interesting astrophysically because, while one wants both higher completeness and efficiency, there is generally a tradeoff involved. The importance of each often depends on the application, for example, an

<https://assignbuster.com/various-project-31-object-classification-classification-is-an/>

investigation of rare objects generally requires high completeness while allowing some contamination (lower efficiency), but statistical clustering of cosmological objects requires high efficiency, even at the expense of completeness.

### 3. 1. 1. Star-Galaxy Separation

Due to their small physical size in comparison to their distance from us, almost all stars are unresolved in photometric datasets, and thus appear as point sources. Galaxies, however, despite being further away, generally subtend a larger angle, and thus appear as extended sources.

However, other astrophysical objects such as quasars and supernovae, also appear as point sources. Thus, the separation of photometric catalogs into stars and galaxies, or more generally, stars, galaxies, and other objects, is an important problem. The sheer number of galaxies and stars in typical surveys (of order  $10^8$  or above) requires that such separation be automated. This problem is a well studied one and automated approaches were employed even before current data mining algorithms became popular, for example, during digitization by the scanning of photographic plates by machines such as the APM and DPOSS. Several data mining algorithms have been employed, including ANN, DT, mixture modeling, and SOM, with most algorithms achieving over 95% efficiency. Typically, this is done using a set of measured morphological parameters that are derived from the survey photometry, with perhaps colors or other information, such as the seeing, as a prior. The advantage of this data mining approach is that all such information about each object is easily incorporated.

### 3. 1. 2. Galaxy Morphology

Galaxies come in a range of different sizes and shapes, or more collectively, morphology. The most well-known system for

<https://assignbuster.com/various-project-31-object-classification-classification-is-an/>

the morphological classification of galaxies is the Hubble Sequence of elliptical, spiral, barred spiral, and irregular, along with various subclasses. This system correlates to many physical properties known to be important in the formation and evolution of galaxies.

Because galaxy morphology is a complex phenomenon that correlates to the underlying physics, but is not unique to any one given process, the Hubble sequence has endured, despite it being rather subjective and based on visible-light morphology originally derived from blue-biased photographic plates. The Hubble sequence has been extended in various ways, and for data mining purposes the T system has been extensively used. This system maps the categorical Hubble types E, S0, Sa, Sb, Sc, Sd, and Irr to the numerical values -5 to 10. One can, therefore, train a supervised algorithm to assign T types to images for which measured parameters are available.

Such parameters can be purely morphological, or include other information such as color. A series of papers by Lahav and collaborators do exactly this, by applying ANNs to predict the T type of galaxies at low redshift, and finding equal accuracy to human experts. ANNs have also been applied to higher redshift data to distinguish between normal and peculiar galaxies and the fundamentally topological and unsupervised SOMANN has been used to classify galaxies from Hubble Space Telescope images, where the initial distribution of classes is not known. Likewise, ANNs have been used to obtain morphological types from galaxy spectra.

Photometric redshifts An area of astrophysics that has greatly increased in popularity in the last few years is the estimation of redshifts from photometric data (photo-zs). This is because, although the distances are less accurate than those obtained with spectra, the sheer number of objects with photometric measurements can often make up for the reduction in individual accuracy by suppressing the statistical noise of an ensemble calculation. The two common approaches to photo-zs are the template method and the empirical training set method. The template approach has many complicating issues, including calibration, zero-points, priors, multiwavelength performance (e. g., poor in the mid-infrared), and difficulty handling missing or incomplete training data. We focus in this review on the empirical approach, as it is an implementation of supervised learning. 3.

2. 1. Galaxies At low redshifts, the calculation of photometric redshifts for normal galaxies is quite straightforward due to the break in the typical galaxy spectrum at 4000Å.

Thus, as a galaxy is redshifted with increasing distance, the color (measured as a difference in magnitudes) changes relatively smoothly. As a result, both template and empirical photo-z approaches obtain similar results, a root-mean-square deviation of  $\sim 0.02$  in redshift, which is close to the best possible result given the intrinsic spread in the properties. This has been shown with ANNs SVM DT, kNN, empirical polynomial relations, numerous template-based studies, and several other methods. At higher redshifts, obtaining accurate results becomes more difficult because the 4000Å break is shifted redward of the optical, galaxies are fainter and thus spectral data are sparser, and galaxies intrinsically evolve over time.

<https://assignbuster.com/various-project-31-object-classification-classification-is-an/>

While supervised learning has been successfully used, beyond the spectral regime the obvious limitation arises that in order to reach the limiting magnitude of the photometric portions of surveys, extrapolation would be required. In this regime, or where only small training sets are available, template-based results can be used, but without spectral information, the templates themselves are being extrapolated. However, the extrapolation of the templates is being done in a more physically motivated manner. It is likely that the more general hybrid approach of using empirical data to iteratively improve the templates or the semi-supervised method described in will ultimately provide a more elegant solution. Another issue at higher redshift is that the available numbers of objects can become quite small (in the hundreds or fewer), thus reintroducing the curse of dimensionality by a simple lack of objects compared to measured wavebands.

The methods of dimension reduction can help to mitigate this effect.