

# Deconstructing a galaxy: identifying components of M83 with photometric clustering<sup>\*</sup>

P. Barmby<sup>1†</sup> and A. K. Kiar<sup>1‡</sup>

<sup>1</sup>*Department of Physics and Astronomy and Centre for Planetary Science and Exploration, University of Western Ontario, London, ON, N6A 3K1*

## ABSTRACT

**Key words:** keywords here

## 1 INTRODUCTION

Galaxies are complex systems, comprised of numerous components with an enormous range of size, mass, density, and composition. These components can be divided into baryonic (stars and their remnants, nebulae, star clusters, nucleus) and non-baryonic (dark matter); cataloging the components and describing the interactions between them is a key step in elucidating the natural history of galaxies. Only in nearby galaxies can individual sub-components be resolved. As observational technology has advanced, the definition of “nearby” has changed and will continue to do so, from Milky Way satellites and Local Group galaxies, to a few Megaparsecs (distance at which stars can be resolved with HST), to XX Mpc (distance at which stars can be resolved with JWST), to the entire observable universe with potential future facilities ().

What is the most efficient way to survey the sub-components of a nearby galaxy? Here we are discussing components detectable in imaging at ultraviolet through infrared wavelengths, i.e. with effective temperatures in the range XX–XX K. Much cooler or hotter types of objects (molecular gas, accreting compact objects) are better-detected at other wavelengths. Particular stellar types, or star clusters, are often identified with broad-band colour-magnitude diagrams (e.g. ). Narrow-band filters can also isolate special stellar types (e.g. ) or objects prominent in emission lines such as planetary nebulae or supernova remnants (e.g. ). Observations are typically designed with detection of particular classes in mind and sometimes re-used for additional purposes (e.g. ). Spectroscopic follow-up is often required to confirm candidates. New observational facilities which provide spatially-resolved spectroscopy (, e.g.) may reduce the need for separate imaging and follow-up steps, but greatly increase the complexity of initial data analysis.

Multi-wavelength surveys are extremely common in studies of unresolved galaxies in the distant universe. While

these are often designed to select galaxies or active galactic nuclei with specific properties (e.g. ), sometimes they are pure blank-field surveys. Broadband ( $R = \Delta\lambda/\lambda < X$ ) filters are the most common imaging modality, although there have been a few attempts at narrow- or medium-band surveys as well (e.g. ?). Clustering in colour space can be used to select particular classes of objects from a survey, for example in selecting AGN via mid-infrared colours (e.g. ), or high-redshift galaxies via Lyman-break dropouts (e.g. ). **give some examples here of sophisticated analysis of colour spaces.**

The purpose of this work is to treat a nearby galaxy as if it were a blank field for surveys, and investigate the usefulness of different photometric colours for identifying sub-components. We make use of the Early Release Science (ERS) observations with the Wide-Field Camera 3 (WFC3) of the nearby spiral galaxy M83 ( ) and in particular the catalog of point sources produced by . We form colours from the photometric measurements in the catalog and apply several clustering techniques to two-colour datasets. In conjunction with published catalogs of galaxy components, we identify the optimum parameters for clustering such a photometric dataset, and the best choices of filter.

## 2 DATA

The dataset used for this study is the Wide-Field Camera-3 Early Release Science (ERS) observations of the nearby spiral galaxy Messier 83 (M83). M83 is a grand-design spiral of type SAB, located at a distance of 4.66 Mpc ( ) and the largest member of the M83 subgroup of the nearby Centaurus group of galaxies ( ). The galaxy’s apparent radius of  $\sim 12$  arcmin ( ) is reasonably well-matched to the camera’s field of view (XX true? XX) **And here we note some other interesting things about M83.**

The objective of the ERS observations as a whole was to probe star formation in galaxies. The observations of M83 were made in broad- and narrow-band filters in order to characterize both stellar and nebular properties. They

<sup>†</sup> E-mail: pbarmby@uwo.ca

<sup>‡</sup> E-mail: akiar@uwo.ca

**Table 1.**

Filter	Name	Exposure time
F225W	Wide UV	1800 s
F336W	<i>U</i> -band	1890 s
F438W	<i>B</i> -band	1180 s
F487N	H $\beta$	2700 s
F555W	V-band, South field	1203 s
F814W	<i>I</i> -band	1203 s

cover a  $3.6 \times 3.6$  kpc<sup>2</sup> region in the northern portion of the galaxy, including the nucleus, a portion of a spiral arm and an interarm region. The spatial resolution of the images is  $0''.0396$  arcsec pixel<sup>-1</sup>, corresponding to a linear scale of  $XX$  pc pixel<sup>-1</sup> at the 4.66 Mpc distance. A complete description of the observations and data processing is given by ?; our work here uses the observations in the UVIS channel, listed in Table 1. A number of previous studies have used the ERS M83 dataset for various purposes. These include studies of star clusters (???Bastian et al. 2011, 2012; ?; ?; Andrews et al. 2014; Chandar et al. 2014; Adamo et al. 2015; ?; ?; ?), H II regions (?), supernova remnants and the interstellar medium (??Blair et al. 2014, 2015), resolved stars (??), and a super-Eddington off-nuclear black hole (?).

We analyze the catalog produced by ? and made available via \*\*REF\*\*. The objects in this catalog were detected on a ‘white-light’ image produced by a weighted combination of the *UBVI* images. Photometry in 0.5- and 3-pixel radius apertures at the positions of the detected sources was performed on the broad- and narrow-band images and tabulated in the Vega magnitude system. We apply the correction to the F657N magnitude zeropoint (from 20.72 to 22.35) noted in the header of the catalog. ? discussed aperture corrections for this catalog, but since we are primarily concerned with colours, we omit any aperture corrections. The catalog contains about 68000 objects which are expected to include individual stars, star clusters, supernova remnants, *Hii* regions, planetary nebulae, stellar blends, and background galaxies. Completeness and reliability of the catalog are not discussed by ?, but a visual inspection of the the detected sources on the white-light image suggests that Nine objects are flagged in the catalog as being problematic and we remove them from our analysis.

Table 2 and Figure 1 characterize the catalog in terms of measurements in individual filters. Not all objects are detected in all filters; Table 2 gives the number of objects for which photometry is reported in a given filter, the number for which reported magnitude uncertainty is 0.2 mag or less, and the aperture magnitude at which the median magnitude uncertainty is 0.2 mag. Figure 1 shows the distributions of magnitudes and uncertainties in the individual filters.

Our analysis in this paper is primarily concerned with colours, rather than luminosities. Uncertainties in colours are computed as the quadrature sum of the relevant magnitudes. Observations in 10 bands allow the generation of 45 different colours, but not all of these colours are likely to be useful in characterizing components of the galaxy. As the F555W band has the most individual detections, we initially compute colours relative to this band. The last column of Table 2 gives the number of objects for which a ‘good’ colour (uncertainty < 0.2 mag) is available.

**Table 2.**

Filter	$N_{\text{obj}}$	$N_{\text{good}}$	$m_{\text{good}}$	$N_{\text{X}-555}$
F225W	57585	1196	m	1149
F336W	61787	3572	m	3493
F373N	55908	236	m	229
F438W	64692	5528	m	5490
F487N	60956	429	m	427
F502N	61715	461	m	460
F555W	66539	13556	m	—
F657N	67819	3285	m	2406
F673N	61979	682	m	679
F814W	63759	11976	m	7199

### 3 ANALYSIS

As the size of galactic surveys grows, the number of dimensions available for analysis increases. In this survey, 45 different colour combinations are possible, creating a space of 45 possible dimensions. Clustering methods provide an efficient way of finding structure in high dimensional data by searching for structure in the feature spaces that cannot be visually inspected. A feature space is a set of  $n$  features that are associated with measurable quantities. In this study, each feature space is defined as a set of colours.

which can lead to the discovery of new patterns in spaces that we are already comfortable with. The methods used in this study are unsupervised techniques. Unsupervised clustering takes the data independent of prior classification, and does not rely on prior labeling of the data.

#### 3.1 Principle Component Analysis

Images of objects at various wavelengths share similar structures. These structures can be built using a subset of only the most significant of those images. Principle Component Analysis (PCA) is a process of determining the most significant features of a high dimensional space in order to rebuild structures effectively ?. PCA projects a data set on its most significant basis, only keeping the dimensions that best explain the data. These dimensions are kept by imposing a maximum number of components on the algorithm, and it is forced to only select the most significant ones ? PCA has been applied to many astronomical surveys... In this study, PCA is conducted to determine the most significant filters, and colours for classifying objects in M83. All ten filters were used in the initial analysis...

#### 3.2 Mean Shift Clustering

Mean Shift is a non-parametric clustering technique that is based on probability density function estimates of each point in the data. Mean Shift is a very powerful algorithm, but has not been widely used in astronomy. At each point, the algorithm estimates the density around that point using a small sample of objects surrounding the point. The power of Mean Shift clustering is that the clusters are not confined to a particular shape. Because Mean Shift moves towards the local mode near the data on which it was initialized, it is useful for estimating the number of significant clusters in a dataset ? The algorithm is based on two components: kernel

**Figure 1.** Distribution of magnitudes and uncertainties for objects in the ? M83 ERS catalog.

density estimation, and density gradient estimation. We will highlight the major components of the algorithm, for a full description of the, see ?.

The first element of Mean Shift is kernel density estimation. The major parameter of Mean Shift is bandwidth,  $H$ , which is assumed to be proportional to the matrix  $H = h^2 I$ , with  $h > 0$  ?. The density estimator for a multivariate density kernel is given by:

$$\hat{f}(x) = \frac{c_{k,d}}{nh^d} \sum_{i=1}^n k\left(\left\|\frac{x-x_i}{h}\right\|^2\right) \quad (1)$$

Where  $h$  is the magnitude of the bandwidth matrix,  $k(x)$  is the profile of kernel  $K(x)$ , and  $c_{k,d}$  is a constant making  $K(x)$  integrate to one?.

The second element of Mean Shift is density gradient estimation. The density gradient is estimated from the gradient of equation 1?. The density gradient is given by:

$$\nabla \hat{f}_{h,K}(x) = \frac{2c_{k,d}}{nh(d+2)} \left[ \sum_{i=1}^n k'\left(\left\|\frac{x-x_i}{h}\right\|^2\right) \right] \left[ \frac{\sum_{i=1}^n x_i k'\left(\left\|\frac{x-x_i}{h}\right\|^2\right)}{\sum_{i=1}^n k'\left(\left\|\frac{x-x_i}{h}\right\|^2\right)} - x \right] \quad (2)$$

The second term of equation 2, is the Mean Shift; the difference between the weighted mean using  $k'$ , and  $x$  ?. Applying a normal kernel to the Mean Shift, the second term of equation 2 becomes:

$$m_{h,K}(x) = \frac{\sum_{i=1}^n x_i \exp\left(\left\|\frac{x-x_i}{h}\right\|^2\right)}{\sum_{i=1}^n \exp\left(\left\|\frac{x-x_i}{h}\right\|^2\right)} - x \quad (3)$$

$m_{h,K}$  is the Mean Shift, and always points in the direction of largest ascent through the estimated density function ?.

Mean Shift clustering involves the application of equation 3 to shift the points of a data set towards the direction of the Mean Shift vector ?. The points are shifted by:

$$x^{i+1} = x^i + m_{h,K}(x^i) \quad (4)$$

Shifting the data points by equation 4 ensures that when the points converge, the center is the area of highest local density, or density “mode”. The density mode can be interpreted as the center of a significant cluster in the data set, and is used to classify the objects that were shifted towards it. Equation 1 introduces the bandwidth parameter  $h$ . Estimating the bandwidth correctly is critical to determining the correct number of clusters. If the bandwidth is too low, the density estimate will be undersmoothed, and Mean Shift will produce many small clusters ?. This is a result of the large density gradient resulting from a low bandwidth, causing many data points to be interpreted as local modes. Conversely, if the bandwidth is too large, a small number of large clusters will be detected, resulting in groupings of data that may blur the underlying structure ?.

### 3.3 K-Means Clustering

K-Means clustering is one of the most widely used clustering methods and has been used to identify a wide range of interstellar and intergalactic objects. It is simple, robust, and easy to implement when analyzing high dimensional spaces, making it a powerful way to analyze galactic surveys. Generally, k-means begins by selecting  $k$  data points at random and deems these points cluster centers. Each object in the data set is then assigned to a cluster center by computing the least-squares distance to each center. K-Means aims to minimize the sum of squares within each cluster given by:

$$J = \sum_{n=1}^N \sum_{k=1}^K \min(\|x_n - \mu_k\|^2) \quad (5)$$

Each point,  $x$ , is then assigned to the cluster center with the lowest distance in equation 5?. Once all data points have been assigned, the centers are re-calculated by taking the average of all the points in each cluster. This process continues until the centers do not change after two consecutive iterations ?.

K-Means requires the number of clusters to be inputted by the user. This poses a challenge for high dimensional data, as the number of clusters cannot be estimated by visual inspection.

- (i) description of clustering and classification
- (ii) description of PCA
- (iii) description of Mean-Shift
- (iv) description of K-Means
- (v) experiments with how to apply the techniques
- (vi) final parameters used

## 4 RESULTS

Well, what did you learn?

## ACKNOWLEDGMENTS

This research has made use of the NASA/IPAC Extragalactic Database (NED) which is operated by the Jet Propulsion Laboratory, California Institute of Technology, under contract with the National Aeronautics and Space Administration. We acknowledge the efforts of WFC3 Science Oversight Committee in conducting the Early Release Science program.

## REFERENCES

- Adamo A., Kruijssen J. M. D., Bastian N., Silva-Villa E., Ryon J., 2015, MNRAS, 452, 246
- Andrews J. E., Calzetti D., Chandar R., Elmegreen B. G., Kennicutt R. C., Kim H., Krumholz M. R., Lee J. C., McElwee S., O’Connell R. W., Whitmore B., 2014, ApJ, 793, 4