

# Exploring a High-Dimensional Galaxy Dataset

Nicole Atherly\*

Mahir Kothary†

Grace Telford‡

## ABSTRACT

We present a tool that facilitates interactive exploration of correlations between parameters in a high-dimensional astronomy dataset. This dataset contains measurements of 25 properties of 10,000 galaxies derived from spectra from the Sloan Digital Sky Survey. Our web-based tool, implemented using D3 and JavaScript, allows the user to rapidly generate one and two-dimensional orthogonal projections of the dataset and dynamically change these projections. It also enables brushing and linking between plots so that the user can search for variations in the distribution of galaxies in different regions of parameter space. We obtain feedback on this tool from astronomers and other scientists, many of whom express interest in applying our system to enable exploration of other multidimensional datasets.

## 1 INTRODUCTION

In astronomy, studying correlations between various galaxy properties (e.g., the mass in stars, the rate of star formation, or the chemical composition of the gas between stars) for large sample sizes can be used to make inferences about the physical processes that govern galaxy evolution. However, measurements of such properties are derived from modeling the distribution of light emitted by a galaxy, and therefore are model-dependent and likely biased. Before drawing conclusions about physical mechanisms based on the correlations between such galaxy properties, it is important to understand the systematic errors that may be affecting the observed correlations.

The process of looking for biases in a high-dimensional dataset can be time consuming, hindering progress in research. It is difficult to choose useful and interesting low-dimensional projections of the data to investigate, given the very large number of possibilities. Further, just looking at the distributions of the data in one or two-dimensional projections does not help in identifying systematics; it is necessary to compare the distribution of the full dataset to subsamples of the data in different regions of parameter space. This can be accomplished by color coding points meeting some filtering criterion, but manually generating many such figures that efficiently sample parameter space using traditional analysis methods (e.g. Python or MATLAB) is prohibitively slow.

A tool that facilitates rapid, interactive exploration of the distributions of a dataset in many one and two-dimensional orthogonal projections would help astronomers to better understand the biases in multidimensional datasets. To be useful,

such a tool must allow the user to quickly change the parameter(s) shown in a given projection. Further, the user must be able to compare data points in a small region of parameter space to the full dataset and quickly switch to a different region of interest in parameter space.

This paper describes our design of a tool that enables detailed exploration of a high-dimensional astronomy dataset containing measurements of 25 different properties of 10,000 star-forming galaxies. Parameters obtained from the MPA/JHU Catalog [4] were measured from spectra taken by the Sloan Digital Sky Survey [10]. The Python module `pyqz` [1] was used to calculate some parameters.

## 2 RELATED WORK

The problem of effectively visualizing high-dimensional data is an active area of research, especially as large and complex datasets are becoming more common. Visual representations are limited to two or three dimensions, and as the number of dimensions increases, the time required to manually search for an informative low-dimensional projection of the data increases rapidly.

Various techniques have been devised to speed up the process of finding useful projections of the data. The technique of parallel coordinates is used to display many parameters of a dataset in two dimensions and allow the user to see structure in the data [3]; however, the ordering of the parameter axes strongly affects the utility of the resulting visualization. TOPCAT [8] is an astronomy-specific tool that allows rapid manipulation of large data tables; however, this tool does not allow for easy and rapid exploration of changing distributions of data in different regions of parameter space. Ggobi [6] is an interactive visualization tool that “tours” the dataset by rotating through many low-dimensional projections of the data so that the user can rapidly see many different views and quickly identify interesting projections. It also allows the user to use brushing/linking between multiple graphs to explore structure in the data.

Several authors have devised methods of automating the choice of useful projections, using orthogonal [7] and non-orthogonal [9] projections. Studies that analyze the relative effectiveness of different methods of choosing low-dimensional projections have found that the optimal dimensionality reduction technique is highly dependent on the nature and size of the dataset [5, 2]. At present, there exist many methods of finding potentially useful low-dimensional projections of multidimensional datasets, but the success of each technique in identifying interesting structure in the data varies greatly with dataset and specific task.

Since this dataset is known to be noisy (both due to Poisson distributed photon counting errors and systematic effects) and the errors are usually non-Gaussian, we choose to avoid machine-learning approaches to finding useful non-orthogonal projections, as such projections would likely be dominated by

---

\*athern@uw.edu

†mahirk@uw.edu

‡otelford@uw.edu

noise. Given that the best method for automatically selecting useful low-dimensional projections is so task-dependent, we design our own system that allows the user to efficiently explore this multidimensional dataset of galaxy parameters and quickly find useful projections of the data. This enables us to include the features tailored to the specific problem of searching for biases in the dataset.

## 3 METHODS

### 3.1 Software

To create the visualization we used D3 (data-driven documents), an open-source JavaScript library to generate svg elements along with HTML and CSS using a data source (here, the galaxy parameters from the Sloan Digital Sky Survey). The program is about 95% in Javascript using the D3 library and the remaining is CSS and HTML. This allows the base page to be loaded quickly while D3 loads and creates the various svg elements, which tends to be slow.

### 3.2 Data

The dataset we use contains measurements of 25 different properties of 10,000 star-forming galaxies. All variables in this dataset are quantitative. Parameters obtained from the MPA/JHU Catalog [4] were measured from spectra taken by the Sloan Digital Sky Survey [10]. The Python module `pyqz` [1] was used to calculate some parameters.

### 3.3 Design Features

#### 3.3.1 Figures & Layout

To provide the user with as much information about the distributions of the data in the different parameters as possible, we provided a dense scatter plot manipulated to leave out several sigma outliers to show the most interesting data. The display also includes four histograms below the scatter plot, each of them with the ability to view a different parameter.

We have ensured that we fit the scatter plot in the first view and the four corresponding histograms in a second view below such that a single scroll allows the user to switch between the two views. This design makes it easy for the user to draw connections between the views, as well as allows the user to efficiently take screenshots of any interesting plots that are generated during the course of exploration.

#### 3.3.2 Visual Encodings

Our design makes use of scatter plots and histograms, using length and position as the primary visual encodings. This is a natural choice in the context of visualizing quantitative astronomy data, in that the intended user (an astronomer) will be familiar with these chart types. These visual encodings are also insensitive to issues of nonlinear perception.

We chose simple and consistent colors to avoid distracting the user and facilitate quickly making connections between the charts. Our color choices also cater to color blind users; see Section 4.2.2.

#### 3.3.3 Dropdown Menus & Tool Tips

We added tool tips to each of the points as well as the overplotted bars representing the selected data on the histograms. This gives the user exact values of the plotted variables and summary statistics for the bars on the histograms which may be

useful for research processes. Simple and dynamically loaded dropdown menus minimize the loading time.

### 3.3.4 Brushing & Linking

The main feature of our tool is brushing and linking. The viewer is able to select any rectangular subsection of the scatter plot to be observed more closely. Upon selection, the points of the scatter plot will turn orange and new orange bars will appear on the histograms below, overplotted on top of the blue bars showing the distribution of the full dataset in that parameter. The orange bars signify how many galaxies from the selection fit into a particular bin in each of the histograms. Viewers are able to quickly see how the distributions of various parameters changes for different selected regions in the scatter plot.

Importantly, the selected region remains when the user changes the parameters shown on the histograms, allowing for a rapid search for interesting distributions in any of the measured parameters. The user may also drag the selected region around on the scatter plot, keeping the shape fixed, to facilitate exploring the entire two-dimensional space in the histogram quickly.

## 4 RESULTS

### 4.1 The Galaxy Explorer Tool

Our tool allows users to easily explore our large set of data. Any two dimensions can be used as the axes for the main scatter plot, as shown in Figure 1. Each point on the plot represents a single galaxy, and its exact values can be viewed by hovering over the point. Four more dimensions can be selected to be displayed in each of the four histograms on the lower half of the page (Figure 2). Through brushing and linking, a selection made on the scatter plot will update each histogram with the distribution of those galaxies within the selected dimensions. This allows for easy comparisons between variables and subsections of the data.

### 4.2 Evaluation

#### 4.2.1 Load Time

Based on tests on several machines, the time to load our visualization is somewhat variable. To load the basic page takes roughly 10 milliseconds and between 1 and 3 seconds to load the data and render the visualization.

#### 4.2.2 Color

Using Spectrum, a Chrome extension, we checked for a range of color blindness and how different interactions may be affected by the blindness. By analyzing different color palettes across interactions we verify that our chosen color scheme does not cause any distortion in how interactions are carried out and does not negatively impact the viewing experience for color blind individuals.

## 5 DISCUSSION

### 5.1 Design Tradeoffs

#### 5.1.1 Number of Histograms

We have approximately 25 parameters, and the intended user must be able to check for interesting changes in the distributions of the galaxies in any of these parameters in different regions of the two-dimensional parameter space shown on the

scatter plot. We chose four histograms, somewhat arbitrarily, to allow the user to look at distributions in several different parameters at once. We felt that more histograms would either overwhelm the user with data or make the visualization too large to be useful; we already need to spread the figures across two views to avoid them becoming too small to be useful.

### 5.1.2 Scatter Plot

We chose to show the data in the form of a scatter plot in order for each individual galaxy to be easily viewed. This allows for users to hover over points and see the exact values for many galaxies. Many points have very close values, and so not all points can be seen easily, but many more are revealed on hover with the tooltips. However, the occlusion proves to be problematic especially on smaller windows where there is less room for the points to spread out. We attempted to change the opacity of the points in order for more densely populated areas to be more obvious, but it slowed down the visualization to an unacceptable level. Fortunately, the dense regions in the scatter plot become immediately obvious when the selection area is moved across the plot, since the orange bars showing the selected data on the histograms are much larger in the denser regions.

### 5.1.3 Layout & Figure Sizes

Though sizing the scatter plot and histograms to each fit on the size of one screen is great for creating larger graphs, it also makes it harder to make comparisons quickly. In order for the viewer to observe the effects of a new brushing, they would have to scroll down to view how each histogram has changed. While it would be nice for the entirety of the visualization to be viewed on one screen, the resulting graphs would be much smaller and perhaps less useful to the viewer especially when selecting smaller groups of points.

### 5.1.4 Selection Updating

Due to memory and browser constraints, we only update the histograms after a brush selection has been fully completed. Turning the selected points orange is trivial and so can be done for the user to better see exactly what they are selecting, but simultaneously drawing the orange histogram bars in the middle of selection would have caused too much lag. However, since most users will only view the histograms after scrolling down once a selection has been completed, this should not prove to be much of an issue.

## 5.2 Informal Feedback

We obtain feedback from two sources: (1) people who viewed our poster describing this project, and (2) astronomy graduate students who were asked to play with our galaxy dataset exploration tool and give their impressions and suggestions.

Overall, both the astronomers and poster session attendees had positive comments about this tool, and many of the people with whom we spoke expressed interest in adapting our code to explore other multi-dimensional datasets in fields ranging from astronomy to biology to materials testing. Most people who used the tool said that our tool would be useful for searching for interesting correlations and biases in a high-dimensional dataset.

Of course, even though people generally liked the tool, we received some useful suggestions for improvement. The following is a compilation of some suggestions from both astronomers and poster session attendees:

- Change appearance of drop down menus to be more modern, perhaps including a mini-view of histogram next to each variable name
- Allow for non-rectangular brushing area (e.g., triangular or freeform)
- Switch to a density map to avoid obscuration issues and better show distribution of data points in scatter plot

## 5.3 Usefulness for Intended Task

Co-author Grace Telford uses this dataset in her research. In the course of playing with this tool during development, she has already seen some interesting trends in plots that she had never thought to make before (e.g., bifurcation in one of the scatter plots generated). Before this project began, a few biases in the dataset had already been found, and those biases are indeed noticeable in the plots generated using our tool. Therefore, from the point of view of the intended user, our design is indeed useful for searching for bias in this dataset.

## 6 FUTURE WORK

### 6.1 Design Alterations

We would like to incorporate many of the suggestions listed in Section 5.2, particularly switching from a scatter plot to a heat map or adding a toggle between the two views. This would provide a more informative view of the distribution of galaxies in the various two-dimensional projections. Of course, we would have to alter the brushing area to snap to the edges of the bins used to generate the heat map.

Generally, this tool can be made even more flexible by allowing the user more control over what types of plots are generated; e.g., each histogram could have a toggle to switch between a histogram and a scatter or density plot.

We could also move the drop down menus for selecting the parameters used in the scatter plot to the locations of the axis labels. This would free up space next to the scatter plot, which could then be devoted to providing more details-on-demand. If a user clicks on a point on the scatter plot (or bin in the heat map), then a box could appear next to the plot showing more information about that galaxy (or statistics of the bin).

### 6.2 Extending to Other Datasets

Given that so many people have already expressed interest in applying this tool to their data, we would like to make it as easy as possible for users to adapt our code to analyze other multi-dimensional datasets. We plan to add a license to our Github repository and advertise this project on personal webpages so that others can fork our repository and use our code. Of course, this will require thorough documentation and a description of how to change the axis labels, which are currently hard-coded.

Because D3 is used to render our figures (and the tool contains a scatter plot with one circle rendered per data point), we are limited in the amount of data that our tool can handle. Currently, it is working well on a dataset of 10,000 galaxies with 25 measured parameters. However, this is a subsample of the

full dataset of about 140,000 galaxies. The clear next step is to add a database backend to this system to speed up binning and selecting data. Our goal is to scale up this tool to handle hundreds of thousands of data points.

## REFERENCES

- [1] Dopita, M. A. et al., 2013, *New Strong-line Abundance Diagnostics for H II Regions: Effects of  $\kappa$ -distributed Electron Energies and New Atomic Data*, The Astrophysical Journal Supplement, 208, 1
- [2] Etemandpour, R., et al., 2014, *Perception-Based Evaluation of Projection Methods for Multidimensional Data Visualization*, IEEE Transactions on Visualization and Computer Graphics, 21, 81 - 94
- [3] Inselberg, A., 1997, *Multidimensional Detective*, IEEE Symposium on Information Visualization, 100 - 107
- [4] Kauffmann, G. et al., 2003, *Stellar masses and star formation histories for  $10^5$  galaxies from the Sloan Digital Sky Survey*, Monthly Notices of the Royal Astronomical Society, 341, 33
- [5] Keim, D. & Krigel, H., 1996, *Visualization Techniques for Mining Large Databases: A Comparison*, IEEE Transactions on Knowledge and Data Engineering, 8, 923 - 938
- [6] Swayne, D., Lang, D. T., Buja, A., and Cook, D., 2003, *GGobi: evolving from XGobi into an extensible framework for interactive data visualization*, Computational Statistics & Data Analysis, 43, 423 - 444
- [7] Tatu, A., et al., 2009, *Combining Automated Analysis and Visualization Techniques for Effective Exploration of High-Dimensional Data*, IEEE Symposium on Visual Analytics Science and Technology, 59 - 66
- [8] Taylor, M., 2005, *TOPCAT & STIL: Starlink Table/VOTable Processing Software*, Astronomical Data Analysis Software and Systems XIV ASP Conference Series, 347, 25
- [9] Yang, J., Peng, W., Ward, M. O., and Rundensteiner, E. A., 2003, *Interactive Hierarchical Dimension Ordering, Spacing, and Filtering for Exploration of High-Dimensional Datasets*, IEEE Symposium on Information Visualization, 105 - 112
- [10] York, D. G. et al., 2000, *The Sloan Digital Sky Survey: Technical Summary*, The Astronomical Journal, 120, 1579

# Galaxy Explorer

Made by Nicole Atherly, Mahir Kothary, and Grace Telford

An interactive tool for exploring correlations between galaxy properties derived from Sloan Digital Sky Survey spectra. Choose parameters to plot on the scatter plot and histograms. Select a region of interest in the scatter plot, then scroll down to see how the distributions for the full dataset and your selected sample compare.

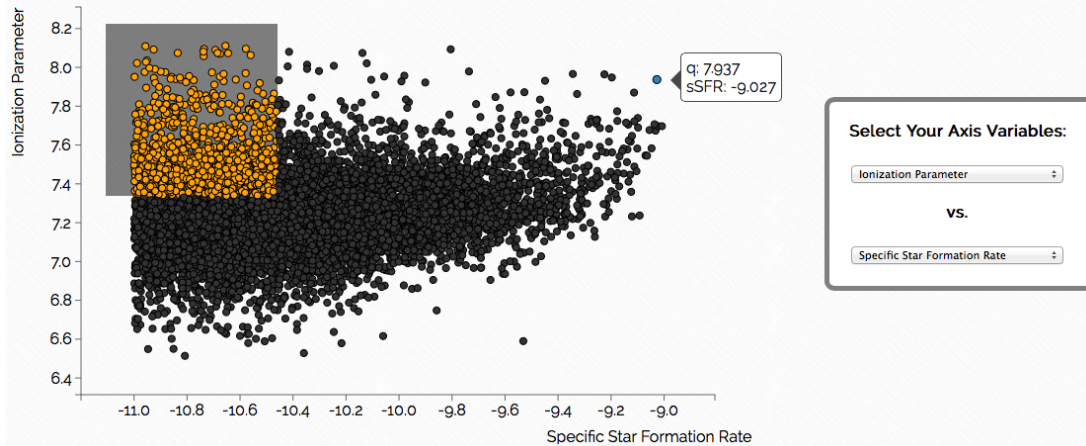


Figure 1: The main view of our design. The user is able to select any two parameters out of 25 available parameters to plot against each other on a scatter plot, where each point shows the values of those parameters for one galaxy. A tooltip shows the exact values of those quantities on mouseover. The user may then draw a rectangular brushing area to highlight points in a region of interest; these selected points are linked to histograms below (see Figure 2. )

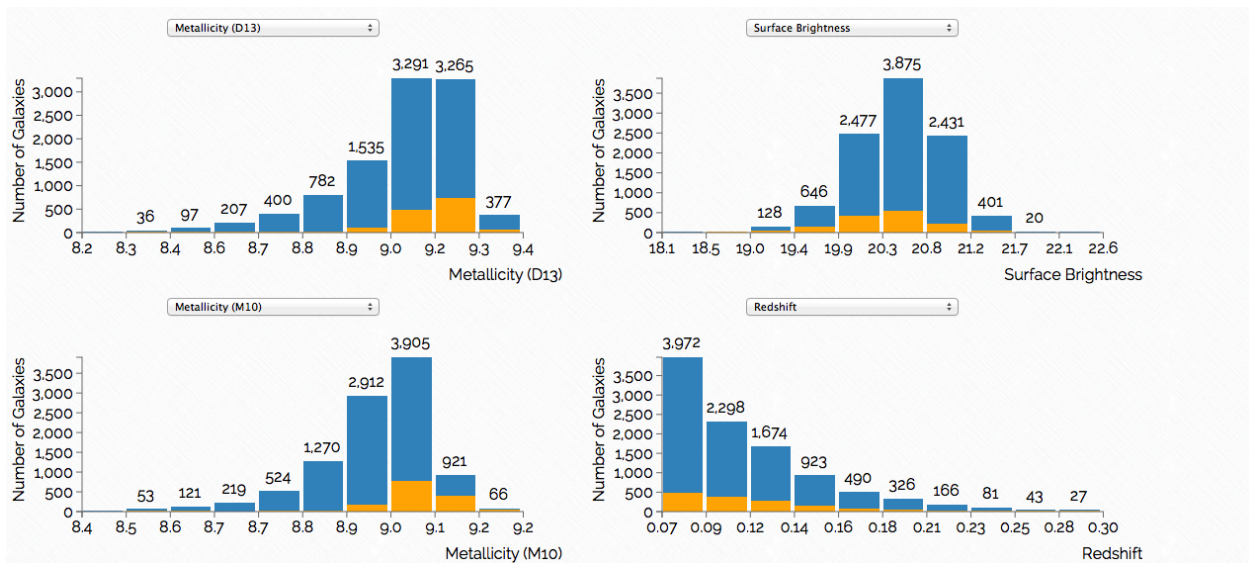


Figure 2: The lower half of our design. Four histograms show the distribution of the full dataset in four parameters, which the user can choose. The selected points inside the brushing area in the scatter plot above (see Figure 1) are over plotted as orange bars, showing how the distribution of that subsample of the data compares to the full sample. A tooltip shows the number of galaxies in the orange bars and the median value of the plotted parameter in each bin on scroll over.