# Exploring the Domain of Choice DJS Program

Sam Mendimasa, Danila Frolov, Celestine Wong, Cathy Poore, and Achuachua Tesoh-Snowsel

**Abstract**— In working with the University of Maryland, Baltimore County (UMBC)'s Choice Department of Juvenile Services (DJS) Intensive Advocacy Program, three visualizations were made from their collected data to help better the services of their AmeriCorps members as they work with the youths in the program. The problems addressed here include visualizing the geographical locations of the youths and teams, identifying optimal times to reach youth via phone, and comparing the effectiveness of each team. This presentation describes the methods and design to each visualization, its functionality and the findings concluded from each visualization. A dot distributed map, heat map and alluvial diagram were made to address each of these problems, respectively, and then presented to the client to assess and facilitate feedback for their program.

**Keywords:** Multidimensional Data, DJS, Dots Distribution Map, Heat Map, Alluvial Diagrams, Parallel Sets

---

## INTRODUCTION

The context of the project derives from the initiative set by the University of Maryland, Baltimore County (UMBC)'s Choice Program. The Choice Program provides 50 volunteer AmeriCorps members the resources to provide different services for youth living in Maryland's high-risk communities. One of these services includes the Choice Department of Juvenile Services (DJS) Intensive Advocacy Program, a community-based alternative to incarceration. The Choice DJS Program works to reduce recidivism with daily face-to-face contact, increase family engagement, and support youth with probation/legal requirements.

DJS members provide their services to a large number of young people in several regions of Maryland and work with a large amount of multidimensional data. This data contains participant demographics (name, race, age, etc.) and daily reports, i.e. the number of visits made by each team, type of visit (whether it is face-to-face or family contact), and the average time spent for each visit. The process of studying and processing this dataset can be difficult and time-consuming. So how to smooth this process by making it less exhausting and more exciting? The effective solution to this problem can be an interactive visual analysis, or simply a data visualization, in which users understand the significance of data by placing it in a visual context. Patterns, trends, and correlations that might go undetected in text-based data can then be exposed and recognized easier with data visualization software.

Thereby, this paper proposes three distinct visualizations displaying the heterogeneous and high-dimensional data provided by the members of the Choice Program. The main goal of this project is to extract the knowledge from the generated visualizations in order to answer the following questions: (1) how are the young participants of the Choice Program distributed within regions served by DJS, (2) what is the best time for DJS members to reach the youth during a day, and finally (3) how effective is the service of each of the DJS teams. We are interested in understanding how effective these visualizations are for users, both as observers of a presentation and as analysts.

## 2 RELATED WORK

This section presents the related work where we found valuable suggestions and techniques of how to make our visualizations effective and understandable.

Few [6] provides information on how to design effective maps. He explains that geographical information is physical and that when we display it, we represent physical characteristics of land masses, bodies of water, and terrain. In trying to deliver information about the geographical locations of the youths in the Choice Program, we decided that a map would be the most beneficial to our client.

Van Wijk and van Selow [10] present a method for the exploration and analysis of extensive time series data. They try to solve the problem of how to identify patterns and trends on multiple time scales (days, weeks, seasons) simultaneously. To achieve this, they use a combination of two methods: cluster analysis and the visualization of the result on a calendar. This is achieved by forming clusters of similar daily data patterns and visualizing the average patterns as graphs and the corresponding days on a calendar. Their method was found to be useful in exploring and visualizing large quantities of univariate time series data. We used van Wijk and van Selow's idea of a calendar to represent our time series data.

Bendix et al. [2] focus on the information visualization technique known as parallel sets which is optimized for categorical data. They describe the parallel sets technique as a method of combining a flexible layout of parallel coordinates with the idea of displaying frequencies as representatives for the categories. With this technique, the dimensions are displayed side by side, the frequency-based representation of categories and relations reveal even more complex information about the data, and meta information is provided. There are two parts of parallel sets: visual metaphor (that properly deals with categorical dimensions) and interaction concept (that facilitates the exploration of the data as well as the creation of new information about the data). We decided to employ this same technique to visualize our categorical data.

## 3 CHOICE PROJECT DOMAIN

In this section, we will discuss the specific demand for creating visualizations in this area. After carefully examining the dataset and needs of our client, we decided to develop three distinct visualizations that will represent our final product, where each of the views answers our client's questions and facilitates in identifying patterns and correlations in a clear and efficient way. The first visualization delivers the information about the geographical locations of the youth that receive support from the Choice Program. The second view allows users to answer the question, what is the best and worst time to contact the youth during the week. Ultimately, the last visual image will display the effectiveness of DJS teams, i.e. the ratio of number of home visits made by each team during the indicated period to the number of youth that successfully completed the program under the supporting team assigned to them. The need to answer three different questions about the data is what led to the creation of three separate visualizations.

### 3.1 Data

The data for this project is a CSV file consisting of over a hundred thousand entries. The 50 AmeriCorps members that volunteer for the Choice Program enter in data for each visit and interaction; these members are also divided into eight teams. The file contains raw data collected during the period of July 2016 to July 2017. The data contains information such as the regions that participate in the program, the teams that pursue the Choice Program mission, the number of visits for each team, the durations of each visit, the youth

demographics, the dismissal reasons, and the contact methods (phone call, family visit, or face-to-face visit). A term "average dosage" defines both numbers of contacts and length of contact. The daily log section, which contains specific information about each visit, provides the most important data that will be used as a base for creating a final product for our project.

## 3.2 Goals

Given our dataset, our client sought to answer the following questions, the distribution of the youths across Baltimore, Montgomery, and Prince George's County, as well as the eight Choice Teams, the best time of day to communicate with DJS youths by phone, and the effectiveness of the service of each of the Choice teams. We created three visualizations to address each of these questions and allows for further exploration of the Choice dataset, which could be used to determine trends and patterns that can improve the program.

## 3.3 Choice DJS Location Distribution

Since one of the requests of our client was to see the locations of the youth advocated by DJS teams, we came to conclusion that the best way to comply with this request was to create a dot distribution map where each dot represents an individual youth. The map is shown in Figure 1 and can be accessed via the interactive link under the image. The color of the dot corresponds to the team supporting that particular region. The key in the bottom of the image displays the names and corresponding colors of each of the eight teams. The visualization follows Shneiderman's mantra of "overview first, zoom and filter, details on demand" [10]. By quickly glancing at this visualization the viewer sees a geographical map with groups of colored dots representing all the youth in the database. Filtering this

data is achieved via an extremely simple mechanism. By clicking on the name of the team displayed in the bottom, a user can see the distribution of dots (youth) depending on the team selected. The search box on the right serves to filter participants by ID number, dismissal reason, date of birth, or race. Finally, hovering over a particular dot allows the user to view demographic information for that participant, including number of home visits by team members and program dismissal reason. This visualization easily conveys that majority of youths participating in the Choice DJS program during the indicated period are located in Baltimore City, whereas the rest of the participants are spread out nearly equal within Baltimore, Montgomery and Prince George's county. Most regions have outliers, i.e. youth moved outside of the area that was initially attached to one of the DJS advocacies, however this team continues to provide services to that particular youth.

## 3.4 Choice DJS Weekly Phone Call Heat Map

The second view requested by the client needed to address the best time to contact youth by phone during the week. This visualization requires a time scale and some means to show frequency of successful calls. Our solution of using a heat map comes from Van Wijk and Van Selow. Their approach to visualizing large amounts of data across the year is through cluster analysis paired with a calendar based visualization. Our problem is similar as it spans across displaying the events across the year, but our visualization ultimately should project the patterns and trends of a week to the client. However, the usage of a heat map is still viable. Our heat map is designed to show which day of the week and what time of day the phone call was made. The time of day is listed along the top and the day of the week listed vertically on the left to create a grid of squares. The colored squares represent the total number of calls in
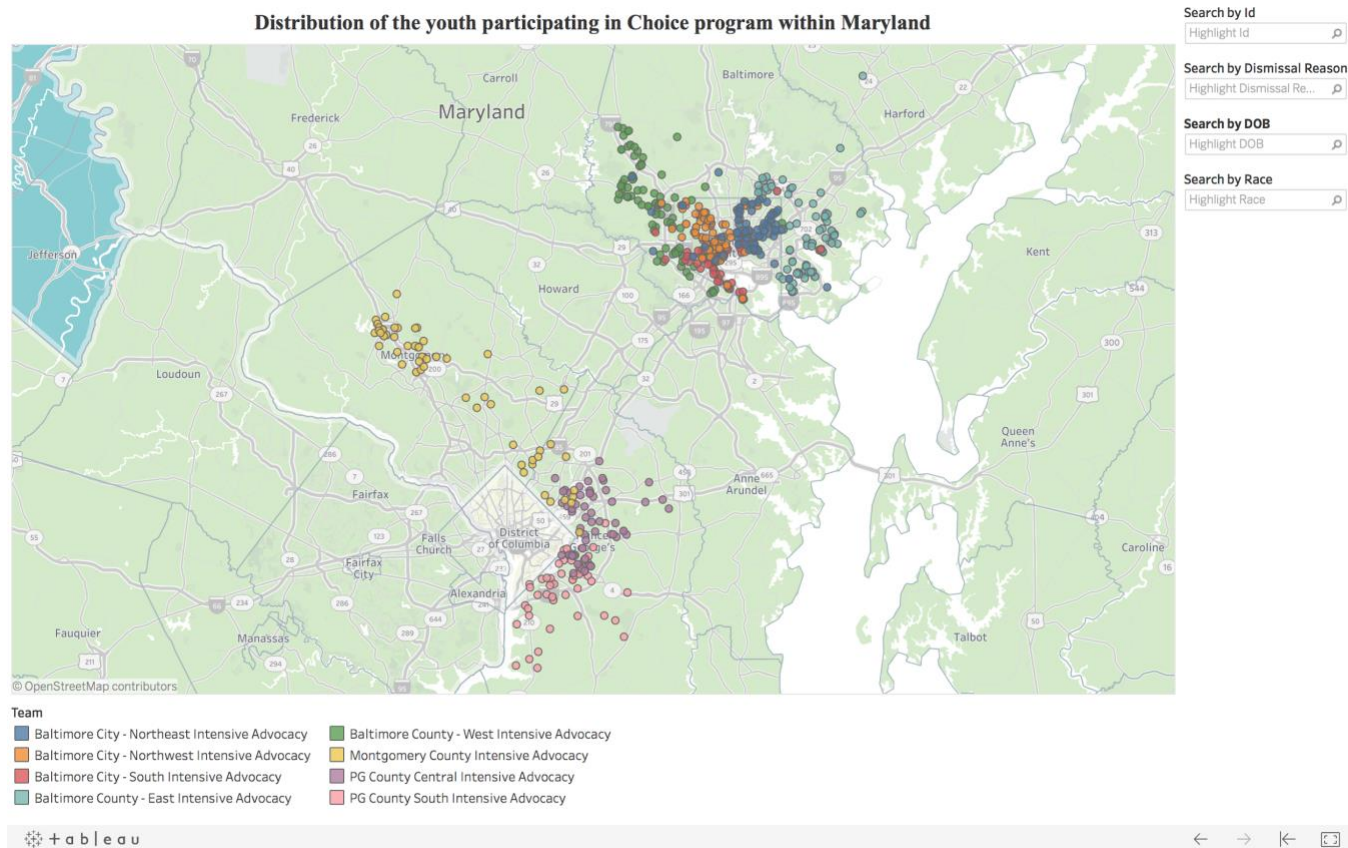


Fig. 1. Distribution of the youth participating in Choice program within Maryland
Interactive Link

which youth were successfully contacted across the entire year. Lighter colors represent fewer successful calls and darker blue squares represent higher numbers of successful calls. A tooltip was added so that the specific number of calls made within a time slot could be seen by hovering the mouse or a cursor over the colored square. The heat map is shown in Figure 2 and contains a link to the interactive visualization.

### 3.5  Choice DJS Team's Effectiveness

Another request from our client was a way of determining the effectiveness of the eight teams that assist youths who participated in this program. Initially, and as mentioned above, effectiveness was simply based on total home visits. However, after reviewing the Choice dataset, we felt the initial approach resulted in many unanswered questions. Such questions included examining if a team can have a high number of visits, but have a negative outcome. Hence, we propose a visualization which not only satisfies our client's initial goals of seeing effectiveness based on visit, but would also be used for exploring further trends in the data. The resulting product was a parallel set. As with any visualization, there will always be drawbacks and advantages, and parallel sets is no exception. The major advantage of parallel sets is its ability to show the flow of data across many categories and to visualize different trends. The downside is that having many categories can make your visualization extremely congested, which makes it more difficult to glean any pattern. The dataset for Choice DJS contains many categories, however, since we are interested in exploring additional trends and outcomes, that better characterize effectiveness for each team, we decided to use parallel sets nonetheless. Since we were aware of the disadvantages of parallel sets, we ensured that the visualization produced would have additional functionality that will allow our users to limit the number of categories, and further drill down into any specific subset of the data.

The visualization we produced (see Figure 3 & 4) allows for the exploration of various dataset points (categories) that support exploring the effectiveness of the eight teams. Home visits are included to resolve our client's need. Additional categories, such as outcomes from the program, and race and gender are also included. Furthermore, the application that we used to create our visualizations has additional functionality that allow importing large datasets. Hence our client can continue to use this application in the foreseeable future to analyze Choice team's effectiveness. Figure 5 shows sample data as .csv being imported into our application.

#### 3.5.1 Overview of Figure 3 and 4

Figure 3 shows one of the many possible trends that can be explored in this visualization. It shows team effectiveness based on whether

Table 1. Initial Outcomes Conversion Table

| Modified Outcomes | Initial Outcomes |
|---|---|
| **Successful** | DJS Terminated - Successful - living in community, Services Completed - living in community, Services Completed - not living in community, Successful Completion - living in community |
| **Unsuccessful** | Case closed - Unsuccessful - living in community, DJS Terminated - Unsuccessful - living in community, Youth detained |
| **Ongoing** | Ongoing |
| **Other** | AWOL/Runaway - living in community, Referred elsewhere, Out of catchment area - living in community, Other: Warrant/Writ, At other advocacy program - living in community, Disruptive Behavior - living in community, Youth placed outside of home - living in community, Youth refused treatment - living in community, Deceased |

the youths they served successfully completed the program or their required community service. The visualization also shows the number of home visits made by the AmeriCorps teams from July 2016 to July 2017. Each team is represented by a different colored bar along the top with the width of the bar corresponding to the number of youths that the team served during the year. Hovering over the team names along the top will show the number of youths served by that individual team compared to all participants in the Choice DJS program. The row through the middle divides the data according to the reasons why the participants were dismissed from the program: ongoing, successful, unsuccessful, or other. Along the bottom, the number of home visits that each youth participant received was divided into three categories: 1-74, 75-149, and 150 and above. Since each team is represented by a single color, it is easy to follow the percentage of youths that completed the program and to see the corresponding number of home visits they received. Hovering over each line in the visualization will highlight all connected components, as well as give its percentage of the dimensions above and below.

Figure 4 shows a possible trend that can be explored with this tool. It presents an overview of Parallel Sets V2.1. The left column displays the various categories that a user can choose to explore. The image that is displayed analyzes the effectiveness of the
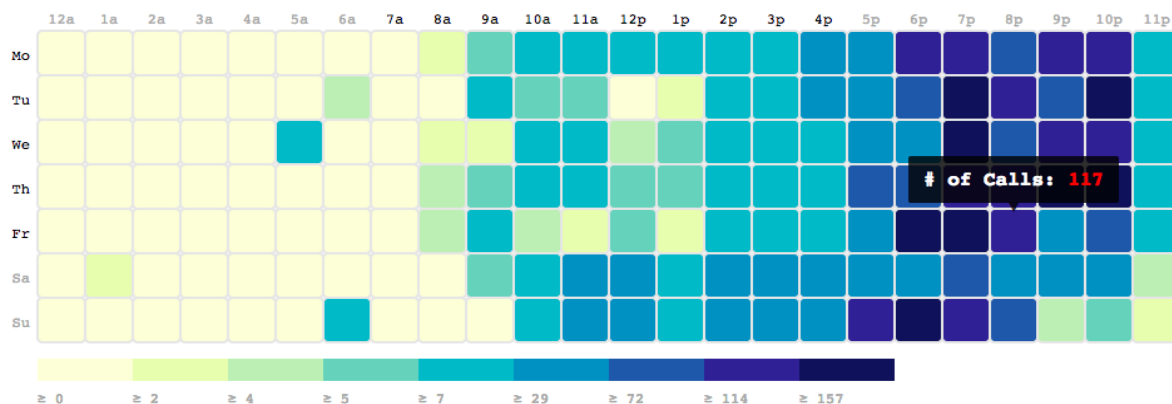


Fig. 2. Heat map of successful phone calls made July 2016 - July 2017 by AmeriCorps members to youths participating in the Choice DJS program
Interactive Link

Montgomery County team based on gender, outcomes, and home visits. Hovering over any of the lines gives the percentages of youths within this team that correspond to the selected category.

## 3.6 Implementation

Each of the visualizations were created using separate tools. This section describes each of the aforementioned visualizations and the process of their implementation.

### 3.6.1. Dot Distribution Map

The map that shows how youth are distributed geographically within different regions of Maryland was created using Java and Tableau. Java was used to parse and analyze the data, and to generate the youths' locations as latitudes and longitudes. In order to show correlations and discover patterns pertaining to the youths' locations, the information from the dataset was converted into coordinates as latitudes and longitudes based on the addresses provided in a dataset. The next step was to upload the modified dataset into Tableau to generate a dot distribution map that highlights visual clusters of the data, which in this case, are the locations of the youth within several distinct regions in the state of Maryland. Geographic data was plotted using custom longitude and latitude values to show an accurate location for each youth. Several layers were applied to the map such as land covers, State/Province borders, county borders, streets and highways. These features help the viewer to obtain the most accurate information about the region, what the surroundings are and where exactly the youth of interest is located. A low-saturated green color was used as a neutral color for background. A reference for choosing a color palette was the work of Cynthia Brewer on ColorBrewer [3]. Thus, the colors were picked to make

sure that they are extremely accessible and easy to distinguish.

### 3.6.2. Heat Map

The heat map was created using Python and D3. Python was used to parse and format the data. The data given to us by the client resembled logs containing information on each phone call. This data included the parties involved, the duration of the call, and the date and time at which the phone call was made. For the purposes of the heat map, we reduced the data to account for date and time and if the phone call was successful. Having reduced these logs, the next step required formatting the dates and times. Each date of the July 2016 - July 2017 year was then categorized into its day of the week (Monday, Tuesday, Wednesday, etc.). Similarly, each successful call was made at a certain time of day and is categorized into its hour of day. Using these two parameters, iterating through each log, each day-hour combination would increment in count to format into a .tsv file in a day-hour-value form. This data was then uploaded with the D3 day-hour heat map code. This application of D3's heat map also creates ranges for each color to scale across the map. The heat map proposes a sequential color scale for easy reading and mapping of the values while keeping a similar color scheme across the range of the data. Lastly, creating interaction with the heat map required using the D3-tip library. Each colored square represents a value, the summed counts of calls, and is shown in the tooltip. The tooltip is made and styled with the D3-tip library, and its integration required adding a listener to each square for a 'mouseover' and 'mouseout' event. By doing so, as each square is created and colored, a tooltip element is attached along with its basic functions. Therefore, hovering over a certain square will give the exact number of successful calls during
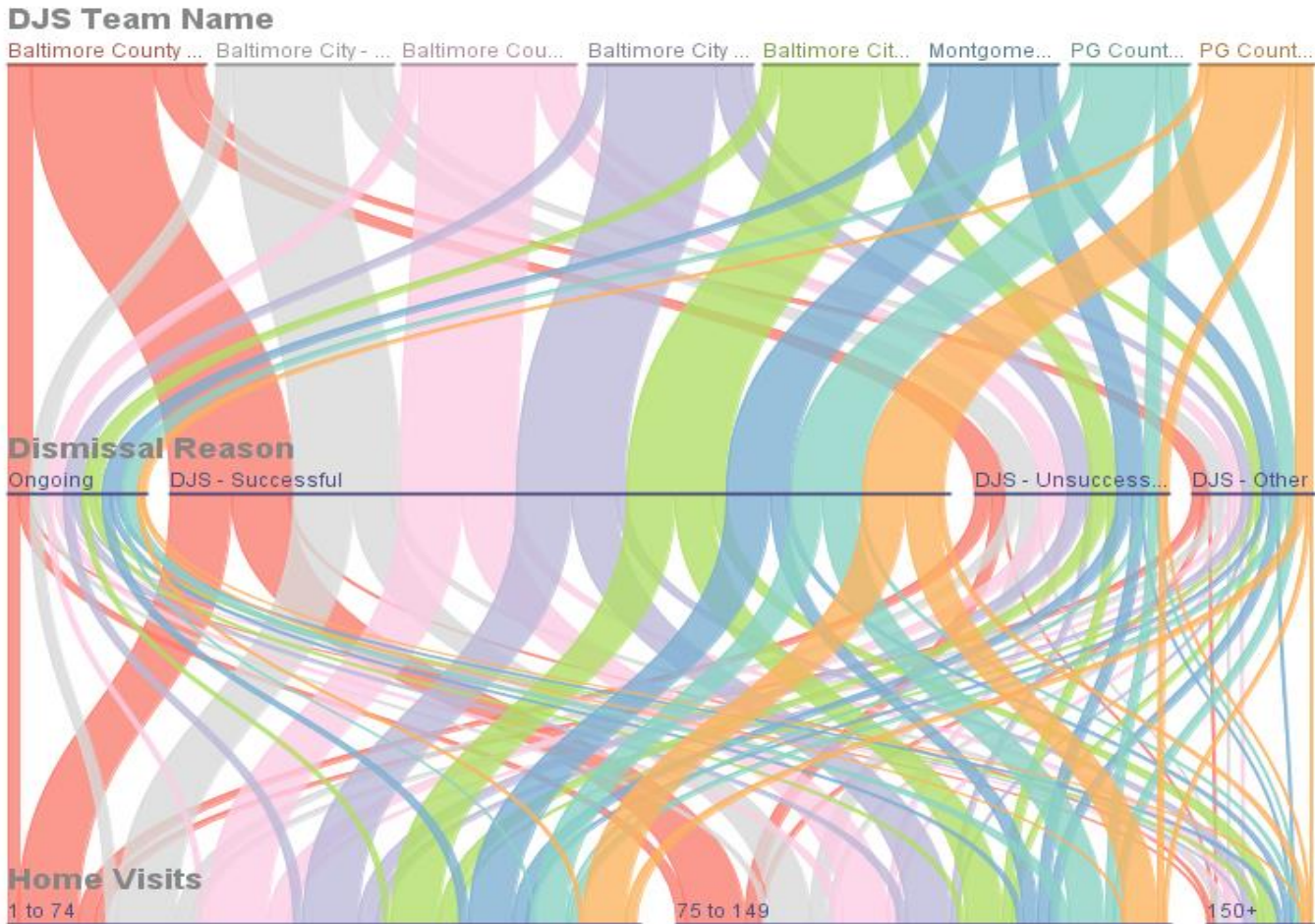


Fig. 3. Choice DJS Team's Effectiveness showing each of the eight teams and the relationship between the outcomes of the youths in each team and the number of home visits they received

that day and time from the year.

### 3.6.3. Parallel Sets

The parallel set was created using Parallel Sets V2.1, an open source visualization application for categorical data. Parallel Sets V2.1 is created with Java, openGL (via JOGL) and sqlite. JOGL is used for supporting 3D graphics, while sqlite serves as the backend for storing and querying the dataset. Parallel Sets V2.1 is limited in the sense that it only accepts categorical data in .csv files and it only works with Java jre 1.6. We have modified it such that it now works for jre 1.7 and above, and have included additional functionality, particularly with the colors to allow other trends to have a more shadowy and grayscale look when something is selected. Additionally, since some of our data is continuous, we used Java for parsing, removing errors, and converting all continuous data into categorical data as a way of limiting our data points for our visualization to be more effective. Hence, continuous data for home visits were converted into three ranges: 1 - 74, 75 - 149, and 150+ visits. Also, categorical data such as youths' outcomes, which initially had fourteen outcomes were reduced to four (see Table 1) to ensure that this visualization is less cluttered. This visualization can only be run locally, which is necessary, since Choice Data contains personal identifiable information, and our project has been designed to only be used by our client and other members of Choice.

## 4  RESULTS

Here we present the results of our findings from the three visualizations created. First, we present the visual findings and how they answered our goal questions. Then we present our client's opinion of each of the visualizations to satisfy the larger goal of meeting our client's needs through creating visualizations that helped her present and explore the dataset she provided.

### 4.1  Visualization Findings

The goal of the first visualization presented (Figure 1) was to show the distribution of youths participating in the Choice DJS program. This visualization allows the user to explore demographics of the youths and shows their location on a map. In this way, the user can look for clusters of youths that have additional characteristics in common, such as DJS team, race, or dismissal reason. Since the teams are each represented by a different color, it is easy to see which teams serve more youths. With each dot representing an individual participant in the program, there is clearly a larger distribution of youths within Baltimore City, around its edges, and clustered closer to the edge of Washington, D.C. than the outer areas of the different counties. This indicates that a lot of the serviced youth live in or near large cities.

The Phone Call Heat Map visualization (Figure 2) shows the number of calls made for which a Choice DJS member was able to contact a participating youth during a time of day, each day of the week, summed over the year. This answers the client's question of when is the most effective time to contact the youth. From the visualization, it is easy to see that the youths most frequently answered the phone in the evenings between 5 and 10 pm. Very few phone calls were answered in the early mornings between midnight and 9 am, with a few exceptions. The day and time labels that are shown in black (Mo-Fr and 7a-4p) indicate the time frame in which the youths are typically in school, and as expected, the colors are lighter within that block indicating that less youth contact was made during those times. Similarly, early mornings show fewer calls, as youth or AmeriCorps members are expected to be asleep.

The goal of the last visualization, using parallel sets (Figure 3 & 4), was to show any differences across the Choice DJS teams in terms of serving youths that had successful outcomes from the program. The visualization mostly focused on the number of home visits a youth received from his or her AmeriCorps team member as this could be a contributing factor for whether a youth is successful or unsuccessful. Each team is represented by a different colored bar along the top with the width of the bar corresponding to the number of youths that the team served during the year. From the
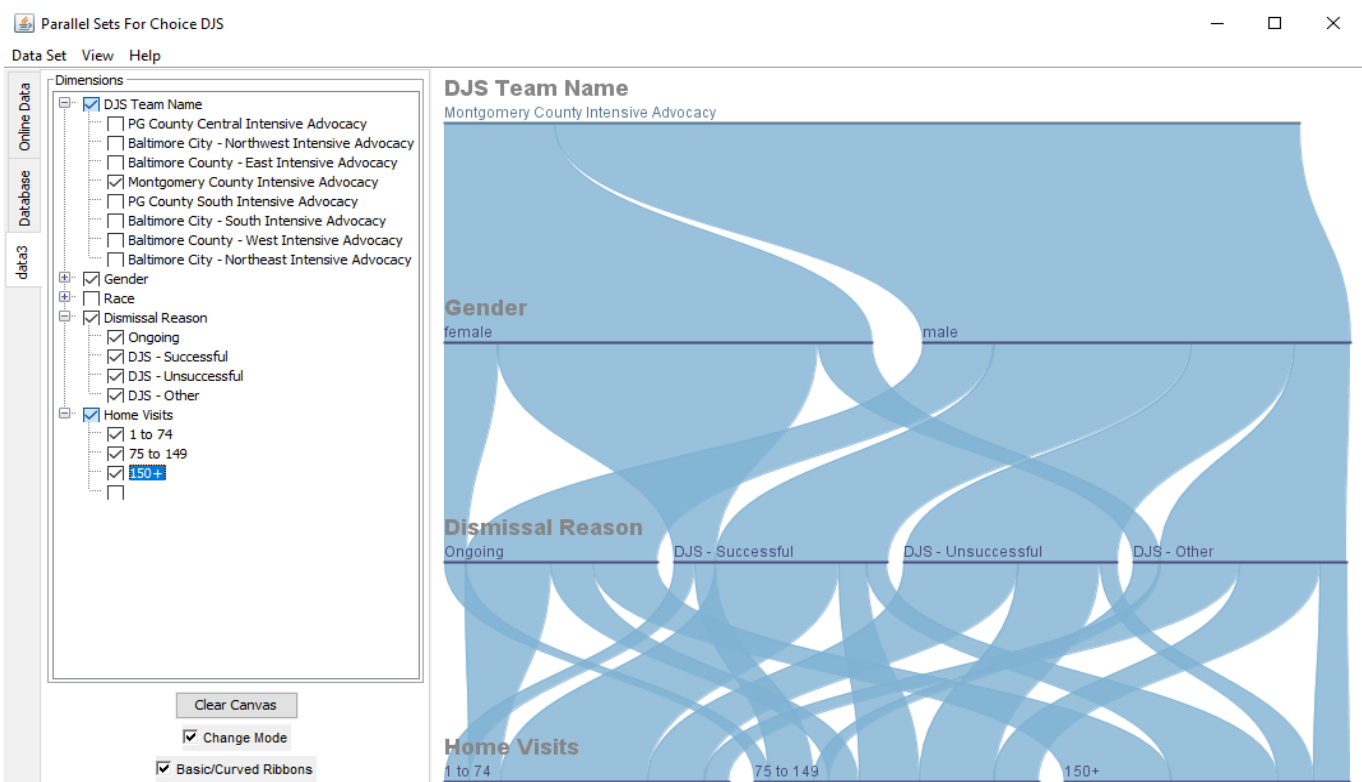


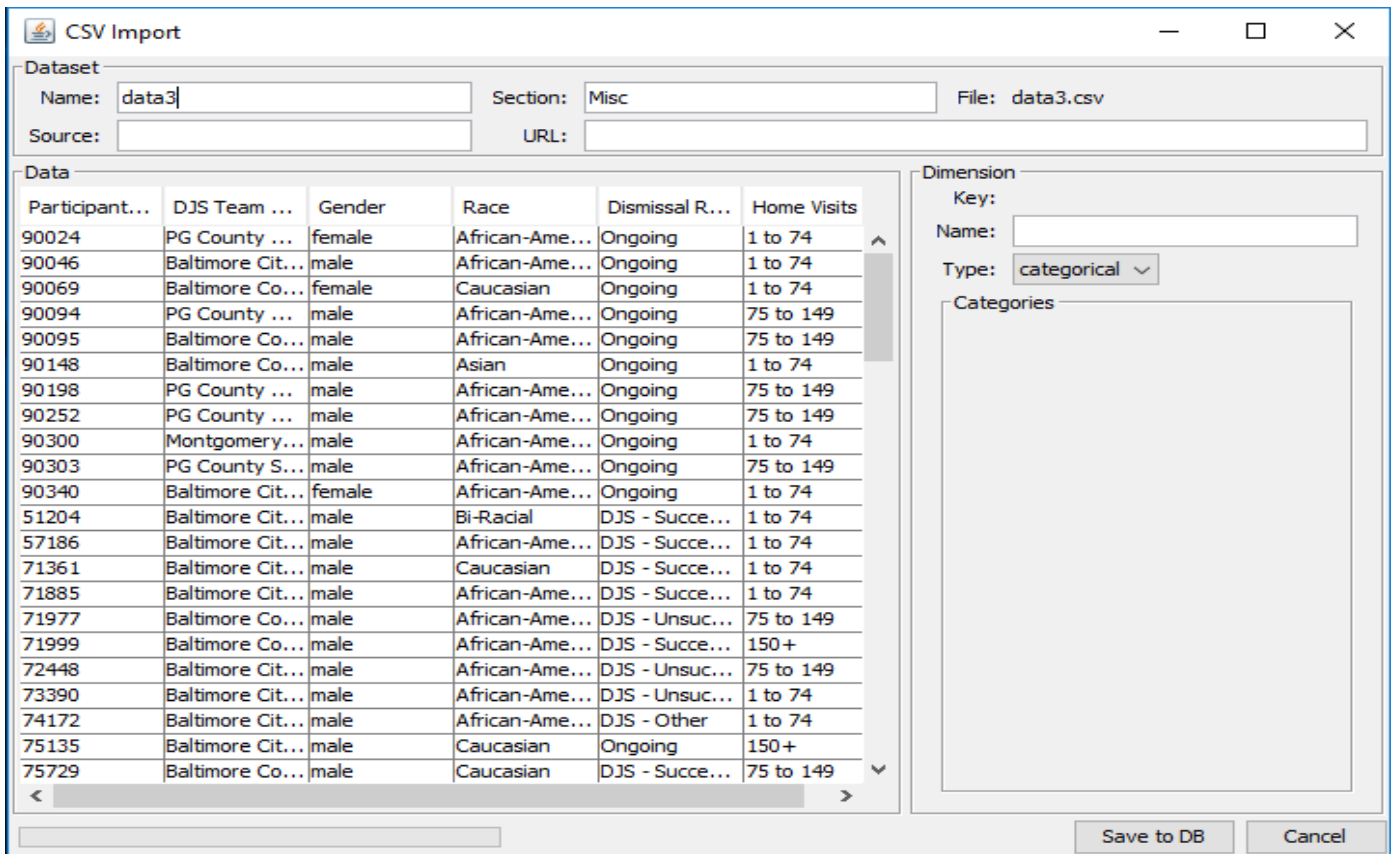Fig. 4. Montgomery County Team Effectiveness and Overview of Parallel Sets Tool

Fig. 5 shows a .csv file that has been imported to be saved in our database

visualization, it is clear that about 60% of youths successfully completed the program during the year of the given data. Since each team is represented by a single color, it is easy to follow the percentage of youths that completed the program and to see the corresponding number of home visits they received. About half of the youths receiving 75 or more home visits had successful outcomes. Hovering over each line in the visualization will highlight its overall connection and give its percentage of the dimension above and below. Through interaction with this visualization we can see that the Montgomery County Intensive Advocacy team had a higher percentage of unsuccessful youths and made fewer home visits than the more successful PG County South team. Further analysis of Montgomery County (Figure 4) shows an interesting trend which can be explored from this visualization. Twelve percent of the youth served by this team are females, and despite the number of visits they received, none were unsuccessful. On the other hand, 20% of the male participants were unsuccessful. Of the 20% that were unsuccessful, 54% received 1 - 74 visits, 38% received 75 to 149 visits, and the remaining 8% received 150+ visits. Overall, although only 46% of the Montgomery County team was successful, 70% of those that were successful received 1 - 74 visits. Hence, one can conclude that having a lower number of visits does evaluate to a decrease in effectiveness.

### 4.2 Client Response

The Choice DJS Location Distribution Map was the first visualization presented to our client. This visualization met our client's goal of easily identifying the distribution of the teams and the youths they served. She mentioned that her team frequently consults a map to view the distribution of the youth and that the colored map with filterable settings would be really useful.

The heat map answered our client's main goal of finding out the times when it was most effective to contact the youth by phone. Our client thought that our visualization was really valuable in showing

effective times to call. Additionally, she would have liked to add the ability to further split the phone call data by month to show a clearer picture, since the youths schedules can change depending on whether they have school. That addition will be addressed in future directions.

Our client expressed interest in learning how to use the parallel sets visualization tool. She did not get the opportunity to dive into exploring the data with our tool prior to the writing of this paper but she was enthusiastic about playing with the tool to further gain insight into what makes different teams effective. Our hope is that this tool will help her gain further knowledge to answer her questions about team effectiveness.

Overall, all three of our visualizations fulfilled our client's request and delivered upon her expectations.

### 5 CONCLUSION AND FUTURE DIRECTIONS

In order to fully realize the goals of our client, we created three distinct visualizations that all offer interactions with different aspects of the data. The first visualization allows for exploring trends in the data by location as all the participants of the program are placed on a map corresponding to their home address. The second visualization allows our client to answer her very specific goal of identifying the best time for AmeriCorps members to contact the youth they serve so that the teams can be the most effective in helping their clients. The final visualization can be used for exploring the data to find which teams have the most youth successfully completing the program.

The three visualizations presented answer many questions and allow for information to be discovered that was not apparent from looking at the raw text data. In addition to looking at how number of visits impact youth success, other components could be added into the parallel sets visualization to look for trends in other factors such as race or location. The heat map could also be expanded to add

additional breakdowns by month or team and allow the user to toggle through those views and look for potential differences.

Adding the ability to easily import new data sets into the first and second visualizations for the user could be an additional feature that could be added to the design. In this way, the user is already familiar with how to use the tools and they could continue to add data for each year for ongoing trend analysis.

## REFERENCES

[1] Aigner, W., Miksch, S., Müller, W., Schumann, H., and Tominski, C. "Visual Methods for Analyzing Time-Oriented Data", *Visualization and Computer Graphics IEEE Transactions on*, vol. 14, pp. 47-60, 2008, ISSN 1077-2626.

[2] Bendix, F., Kosara, R., and Hauser, H. "Parallel Sets: Visual Analysis of Categorical Data", Proceedings of IEEE Information Visualization 2005, pp. 133-140, 2005

[3] Brewer, C and Harrower, M. "COLORBREWER 2.0," ColorBrewer: Color Advice for Maps. [Online]. Available: http://colorbrewer2.org/#type=sequential&scheme=BuGn&n=3.

[4] Chi, Ed H. "A Taxonomy of Visualization Techniques using the Data State Reference Model" Proceedings of the IEEE Symposium on Information Visualization, 2000

[5] Dang, T.N., Anand, A., Wilkinson, L. "TimeSeer: Scagnostics for High-Dimensional Time Series", *Visualization and Computer Graphics IEEE Transactions on*, vol. 19, pp. 470-483, 2013, ISSN 1077-2626

[6] Few, St. "Introduction to Geographical Data Visualization", *Perceptual Edge Visual Business Intelligence Newsletter* , March/April 2009

[7] Lammarsch, T. "A compound approach for interactive visualization of time-oriented data", *Visual Analytics Science and Technology 2008. VAST '08. IEEE Symposium on*, pp. 177-178, 2008.

[8] Roth, S. and Mattis, J. "Data characterization for intelligent graphics presentation", C*HI '90 Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, pp. 193-200, March 1990

[9] Shneiderman, B. "The eyes have it: A task by data type taxonomy for information visualizations", *Proc. 1996 IEEE, Visual Languages*.

[10] Van Wijk, J. and van Selow, E, "Cluster and Calendar based Visualization of Time Series Data". Proceedings of Information Visualization 99, pp., 1999