

## DSC Capstone Project 2

### Donors Choose: Donation Recommendation System

#### Milestone Report 2

**Problem:** DonorsChoose.org funds classroom requests through the support of over 3 million donors. DonorsChoose.org is unique because it allows the donor to fund specific projects for specific schools, and specific teachers. Donors are able to browse projects, and fund ones that are important to them. Requests come in the form of all kinds of items, from books, to furniture, to electronics. In order to continue supporting public school classrooms, DonorsChoose.org must work to ensure that first time donors continue to make ongoing donations. Donors typically first visit DonorsChoose.org to make a one time donation, and for the sustainability of DonorsChoose.org it is important to convert one time donors into recurring donors.

**Dataset:** The dataset was obtained from <https://www.kaggle.com/donorschoose/io>. This dataset had been shared during a previously held contest. The data was initially broken up into 6 different .csv files (donations, donors, projects, teachers, resources and schools). In order to understand the data better, all data was initially examined separately.

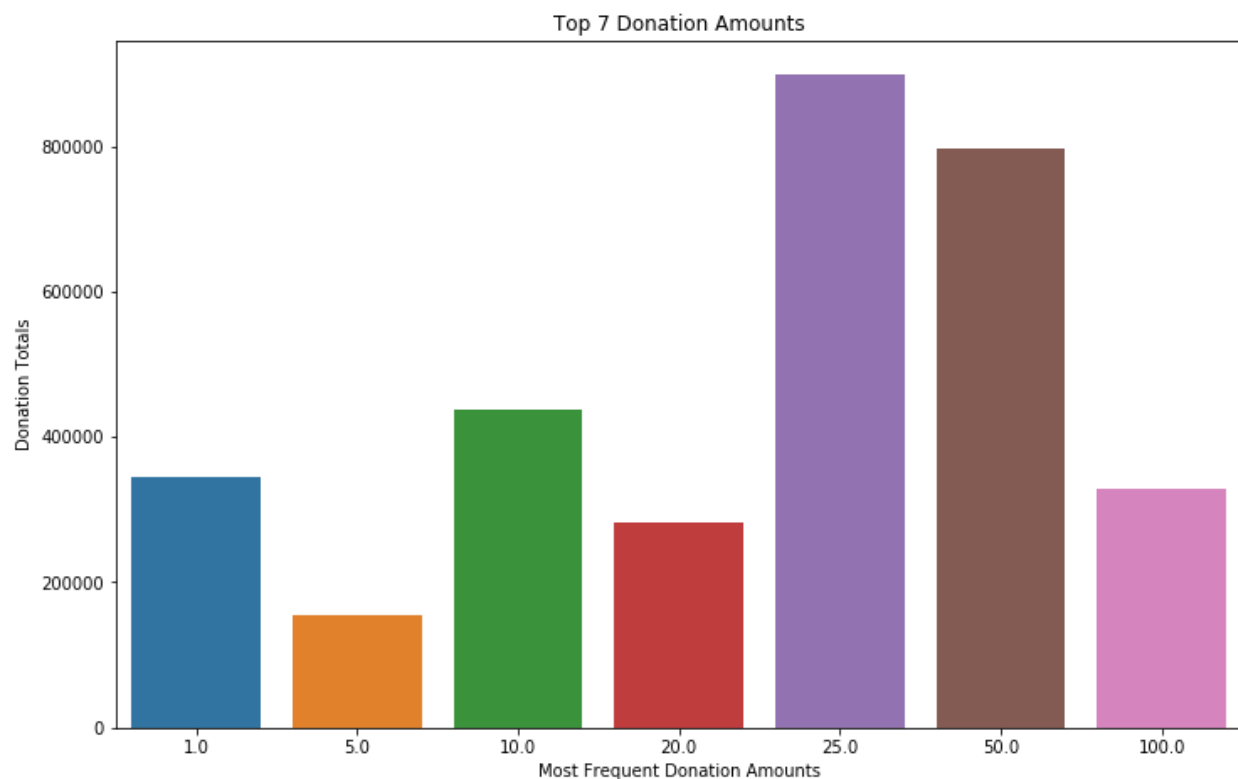
I initially cleaned and began exploring this dataset in Jupyter notebook, primarily using pandas. I quickly discovered that using my current toolset in Jupyter notebook was not going to be effective to process the large amount of data that I was working with. I moved my work to Databricks in order to begin using PySpark. This proved to be an exciting challenge.

**Data Cleaning:** In order to process the data in Databricks, the datasets were uploaded to an S3 bucket in AWS. I used EMR to create a cluster and notebook in AWS, due to the size of the dataset. In order to prepare the Projects .csv for import, a schema needed to be created that addressed each data format type per column. Additionally, there were initially issues with multiline text and the initial data import. Text columns, like project essay, were spilling over from their respective column into neighboring columns, and sometimes into the next row depending on the length. This was addressed by ensuring that the Projects .csv had the escape key parameter changed to ' " ', along with multiLine set to True, and the schema applied at import.

Once all of the data was read into their respective dataframes, I combined them into one dataset by joining on the ID columns (school\_id, project\_id, teacher\_id).

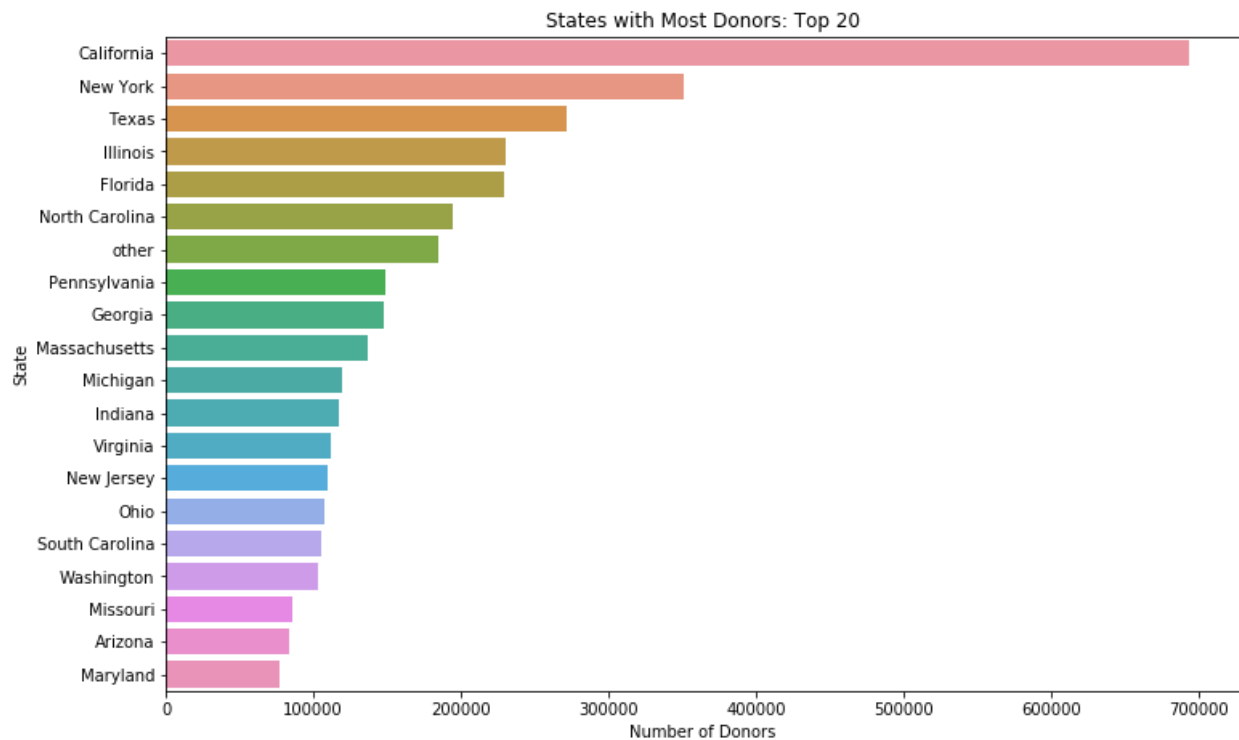
**Initial Findings:** In order to better understand the behavior of donors, much of the exploratory data analysis is focused on the donors, their donations, and the schools and projects that they donate to most frequently. In examining the data, I was curious who was making the highest donations, how frequently donors made recurring donations, did their donations occur at the same schools, and were the donations made for similar projects. I was also interested in determining whether there was a specific time of the year when donations increased or decreased.

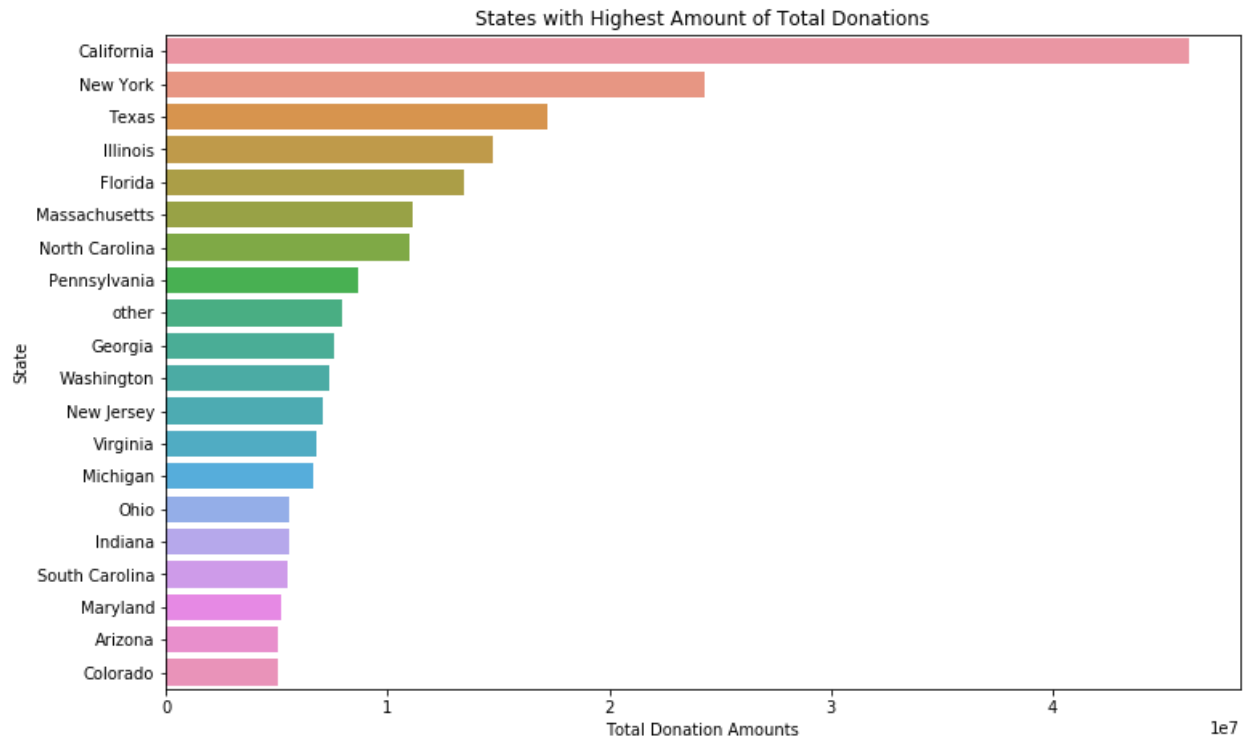
- There are 2,122,640 donors - 552,941 are repeating donors. This only accounts for 26% of donors! Additionally, I found that the number of recurring donors significantly decreases after donors who make more than two donations:
  - More than 2 donations - 278,039
  - More than 5 donations - 98,487
  - More than 10 donations - 40,299



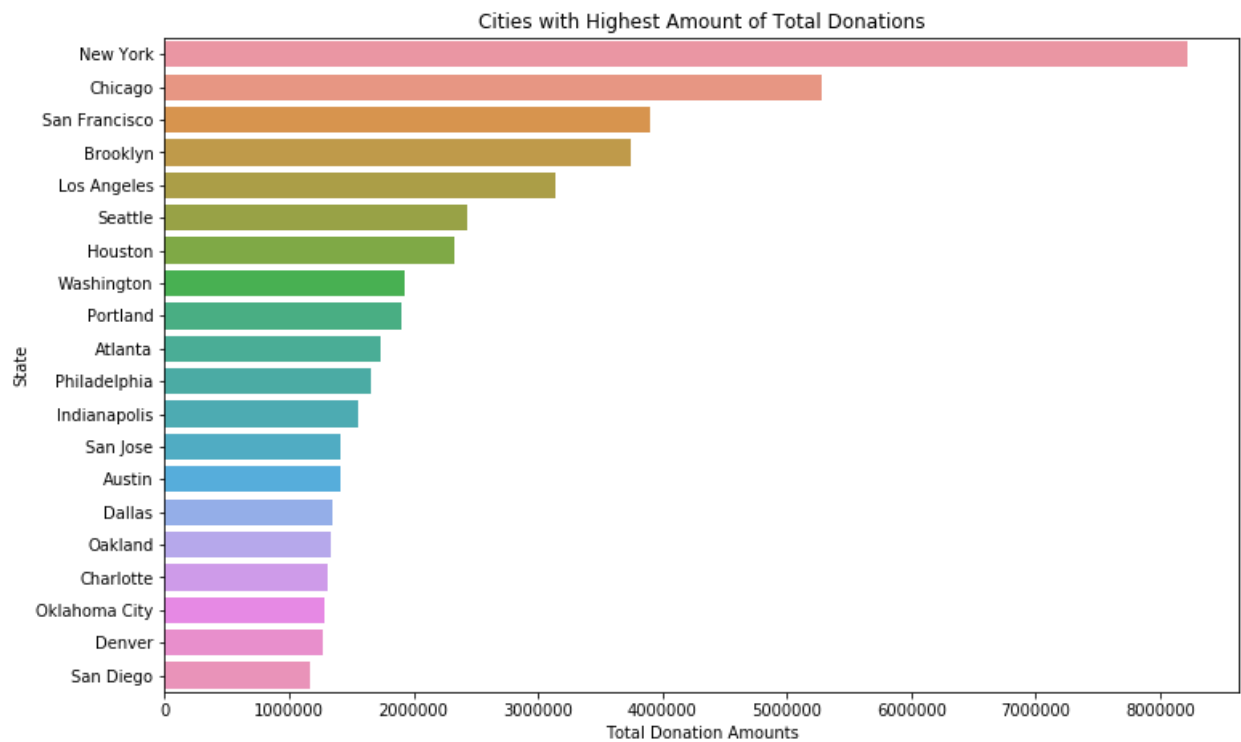
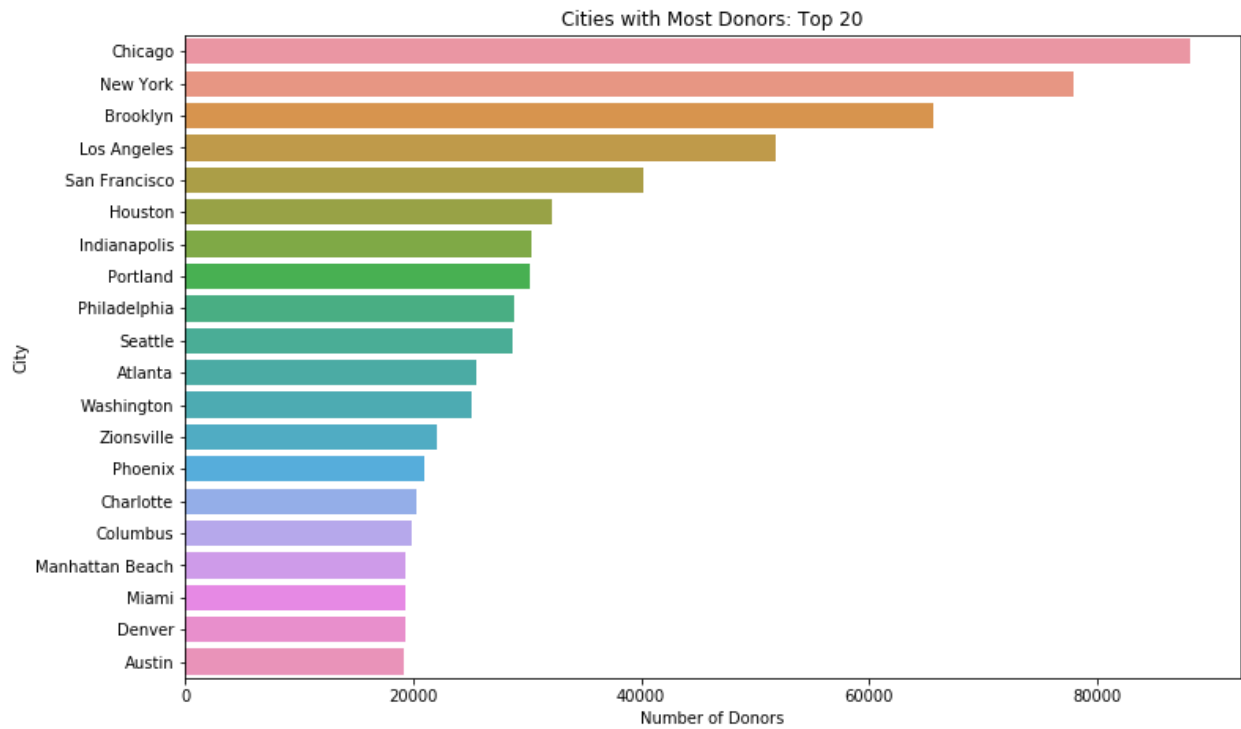
- \$25, \$50, \$10, \$1.00, \$100.00, \$20, and \$5 are the most frequently donated amounts - it's likely that these are pre-set donation options on a donation page.

- Donors are identifiable by a randomized donor ID - based off of this, we're able to determine that the top donor has made more 18035 total donations. This top donor has donated a total of \$37,121 and is from Manhattan Beach, California.
- Surprisingly, that donor ID is not the same as the donor ID who has made the most (total) donations to Donors Choose projects, totaling \$1,864,016. This is an interesting find, and leads me to believe that there are donors who donate small amounts, over a long period of time and there are donors who instead donate less frequently, but with higher donations.

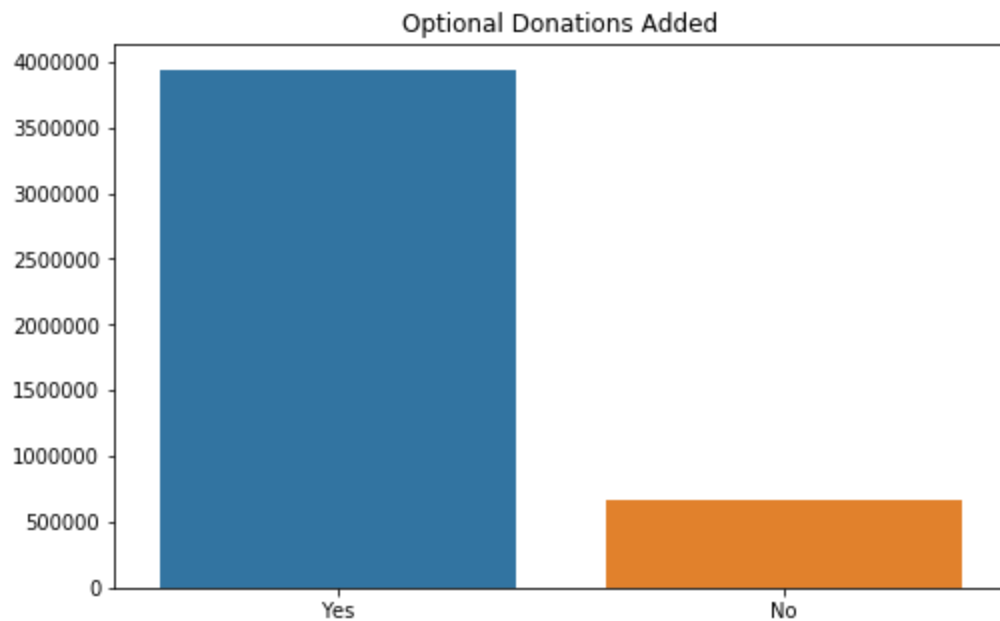




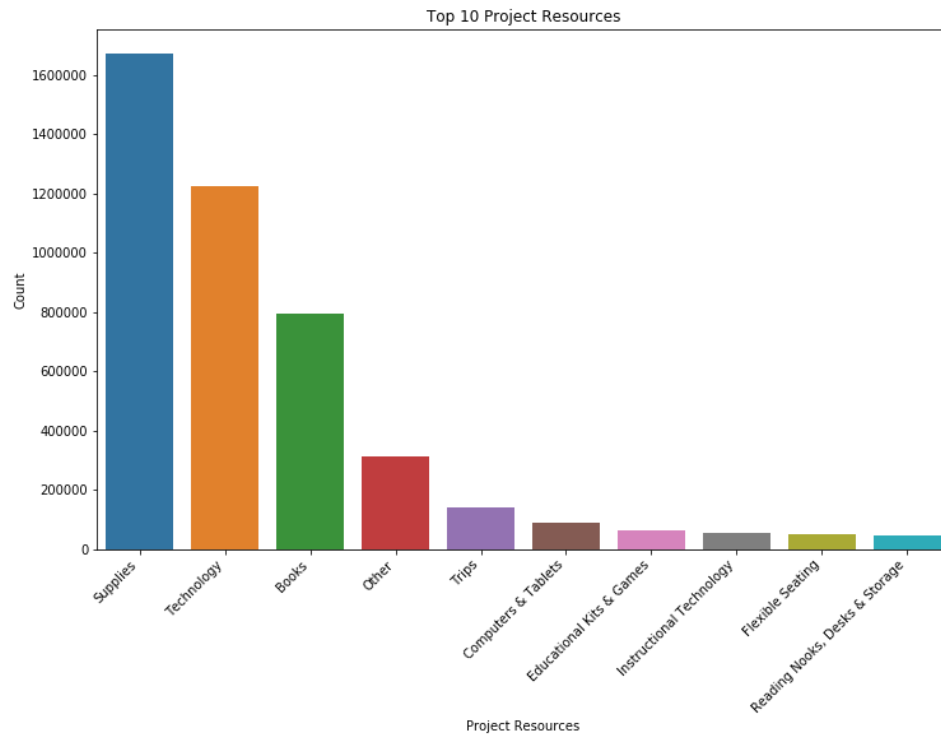
- The top three states with the most unique donors are also the same top three states that have the highest total of donations: California, New York, and Texas. Due to the size of these states, it makes sense that they would be the top three states for both the amount of donors and the highest total donations.



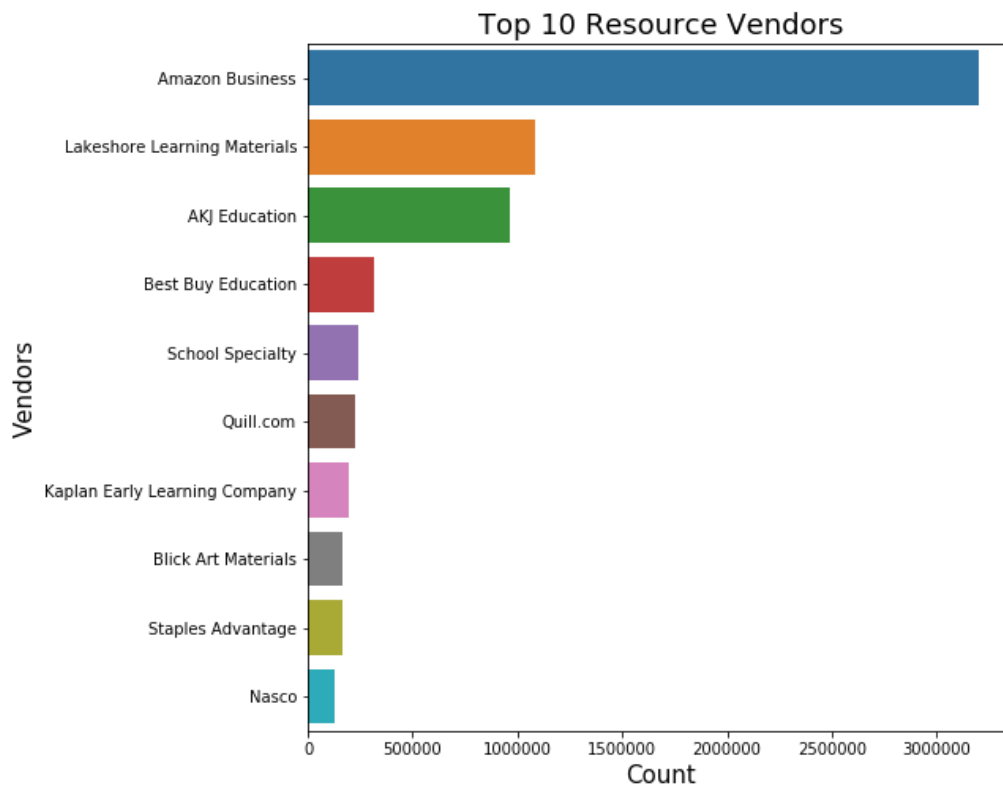
- Cities with the most donors and highest donations stray from the states slightly - the top three cities for unique donors are: California, New York and Brooklyn - while the top three cities for highest amount of donations are New York, Chicago and San Francisco.
- The majority of donors choose to add an optional donation when they are donating to a specific project. This means that they will donate 15% of their project donation to Donors Choose.

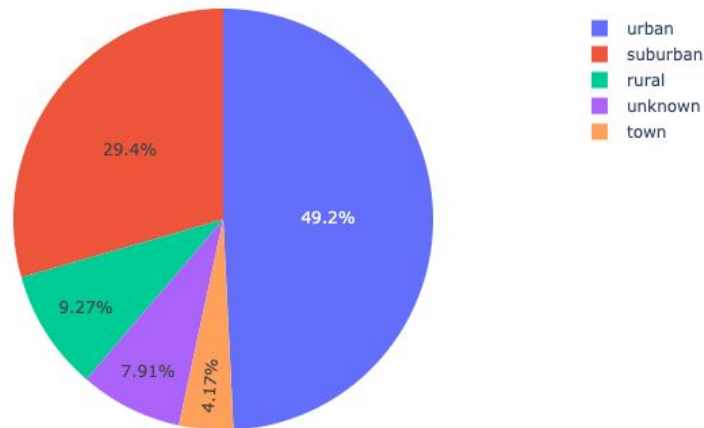


- The top three project main categories are literacy and language, math and science, and literacy and language/math & science, the subcategories match up with this - being: literacy, literacy/mathematics, literacy/literature and writing



- Amazon leads the top resource vendors, followed by Lakeshore Learning Materials and AKJ Education.





- School areas that are in the most need are in urban areas. Schools in urban areas also have the most projects funded by donations.

**Summary & Next Steps:** This exploratory analysis has helped to better understand donors who make one donation vs. donors who make recurring donations, and which areas one-time donors are most heavily populated in. This information, coupled with an understanding of the most requested supplies as well as areas and schools that are in the most need, provides important insight into potential future campaigns to convert one-time donors into recurring donors.

Based off of this initial exploratory data analysis, there is insight into where future donation campaigns may potentially focus in the future. The data provides a glimpse into understanding the areas of most need - urban area schools. The top cities for donors reflect this, as cities such as Chicago, New York, Brooklyn, Los Angeles, and San Francisco see the most donors. The projects most donated to, and resources most requested is an additional feature which can potentially help identify one time donors who could be converted to recurring donors. It is also interesting to note that Amazon Business is the top resource provider. Due to the amount of people who use Amazon.com regularly across the country, a future partnership with Donors Choose could potentially help retain donors as well.



**Modeling:** The objective of this project is to match donors with new projects that may interest them based on their prior donations. In doing so, this could create momentum for donors who have possibly only donated to one project in the past. There are three different types of recommendation systems that I explored using to solve this problem:

**Collaborative Filtering:** This method is used to make automatic predictions about project preferences that a donor may have. Based on projects that they've donated to in the past, this method will predict future behavior based on similar projects that other donors have donated to.

**Content Based Filtering:** Alternatively, content based filtering uses the information from the description, project essay and project need statement of previous projects that a donor has donated to. The content based filtering method then predicts other projects to these donors based upon their content similarities.

**Hybrid (mix between Content Based and Collaborative Filtering):** A hybrid approach combines collaborative filtering methods and content based filtering methods. By using a combination of project content and prior project donations, recommendations can be made for future projects.

*Pre-Process Data:* In order to prepare the data for the first collaborative filtering model, it needed to be prepared. Because our dataset does not have any sort of ranking system in which the donors can rate the projects they're donating to, I decided to use the donation amounts per project, per donor. In order to do this I originally created a table with the donor ID, project ID, and donation amount (sum):

donor_id ▲	project_id	donation_sum
17f3f75fe1d4c9058f9a36720b312fbc	0000c0bdc0f15bd239cffa884791a10	175
5d499cba7badd9eb0dce8b00b6a4370b	0000c0bdc0f15bd239cffa884791a10	70
72b04f710931d0699e0975d8d40a103e	0000c0bdc0f15bd239cffa884791a10	175
99488c9de95c1f7ae3ca4748bf211744	0000c0bdc0f15bd239cffa884791a10	175
9f845a4a86f2982c3f5f402de41377f5	0000c0bdc0f15bd239cffa884791a10	140
a3cfc28ca16f2927d02203c142279a16	0000c0bdc0f15bd239cffa884791a10	158.7600021
a7fa881a86cb4c7e70815b38168d1702	0000c0bdc0f15bd239cffa884791a10	21
e1a9ef3d18fad06f5ab2f7c1aab38e2e	0000c0bdc0f15bd239cffa884791a10	350
ef5f051d5008f36e81f14f0d38b0e69a	0000c0bdc0f15bd239cffa884791a10	140
fa843f037182711c3ea4483983c8af2b	0000c0bdc0f15bd239cffa884791a10	350
fcec0988e54a7607f172e13079839b4e	0000c0bdc0f15bd239cffa884791a10	175

Once this table was made, I created a new column in the table that is the log of each donation amount, as it is a convenient way to express large numbers:

donor_id	project_id	donation_sum	donation_log
17f3f75fe1d4c9058f9a36720b312fbc	0000c0bdc0f15bd239cffffa884791a10	175	5.164785974
e1a9ef3d18fad06f5ab2f7c1aab38e2e	0000c0bdc0f15bd239cffffa884791a10	350	5.857933154
ef5f051d5008f36e81f14f0d38b0e69a	0000c0bdc0f15bd239cffffa884791a10	140	4.941642423
fa843f037182711c3ea4483983c8af2b	0000c0bdc0f15bd239cffffa884791a10	350	5.857933154
99488c9de95c1f7ae3ca4748bf211744	0000c0bdc0f15bd239cffffa884791a10	175	5.164785974
a3cfc28ca16f2927d02203c142279a16	0000c0bdc0f15bd239cffffa884791a10	158.7600021	5.067393641
5d499cba7badd9eb0dce8b00b6a4370b	0000c0bdc0f15bd239cffffa884791a10	70	4.248495242
9f845a4a86f2982c3f5f402de41377f5	0000c0bdc0f15bd239cffffa884791a10	140	4.941642423
a7fa881a86cb4c7e70815b38168d1702	0000c0bdc0f15bd239cffffa884791a10	21	3.044522438
fcec0988e54a7607f172e13079839b4e	0000c0bdc0f15bd239cffffa884791a10	175	5.164785974
72b04f710931d0699e0975d8d40a103e	0000c0bdc0f15bd239cffffa884791a10	175	5.164785974
ab00b62ea86ec20c2eb4a69b40c6fa80	0000c0bdc0f15bd239cffffa884791a10	350	5.857933154

In order to finish prepping the data for the CFM model, the donor\_id column and project\_id column were selected into their own data frames by distinct value. Once in these dataframes, integer IDs were assigned to each before being re-joined with the original table. After renaming and dropping old columns, the table ready for modeling was completed:

donor_int_ID	project_int_id	donation_log
0	0	5.164785974
1	0	5.857933154
2	0	4.941642423
3	0	5.857933154
331792	0	5.164785974
331793	0	5.067393641
663060	0	4.248495242
663061	0	4.941642423
994514	0	3.044522438
994515	0	5.164785974

Lastly, all columns were formatted as integers, as that is the only format that ALS will accept. With the data frame cleaned up, we can begin creating our model!

*Model Preparation:* The algorithm I chose to use for this collaborative filtering based model is Alternating Least Squares (ALS). In order to prepare the data for the model, I split it into a training and testing set, using a .6/.4 split. Once the data was split, the training data was fit to the ALS model with default parameters. The testing data was then used to make predictions.

*Alternating Least Squares:* The ALS model with default parameters was run on the training and test data, and then evaluated using Root Mean Square Error (RMSE). The training data was first fit to the ALS model, and then the testing data was used to create predictions. When the ALS model was run on the testing data, using default parameters, the resulting RMSE was 1.9441559327012328.

*Hyperparameter Tuning:* In order to determine what the best parameters would be for this model, a parameter grid was created, along with a 5 fold cross validation. By fitting the training data to the cross validator, we were able to determine the best parameters for both rank and maxiter. With these parameters added to the ALS model, the new RMSE is RMSE: 1.9401127004230694. Tuning these hyperparameters continues to be a challenge due to memory constraints and the overall size of this data set.

The RMSE determines that the predicted recommendations would fall within the range of 1.94, below or above. The total range of the donation sum, converted to logs is -4 and 15. With that in mind, a 1.94 RMSE seems to be a fairly decent indicator of recommendation accuracy.