

Project: Capstone Project 1: Milestone Report

Three questions can be answered from this data that can be considered once machine learning is implemented. The first question was the question asked in the Kaggle competition in which this data was created. Based on all of the items purchased, what is a proper item one can predict will be bought in the next order. More specifically, what is the following product a user will buy based on all of the items purchased in the history of the users' orders? The second question that can be answered from the data deals with an instance of when a customer is exploring the items to be purchased. If a customer is exploring the things available, what items can be recommended that will allow them to 'discover' a new product? It is similar to the recommendations that Netflix and Amazon offer on their websites. Perhaps a customer is a frequent buyer of ice cream. I think it would be a safe bet to recommend a product to accompany this ice cream, maybe some sprinkles or chocolate syrup. This idea leads me to the third question that can be answered from the data, which includes, "What products can usually be purchased together?" Some products are purchased together all the time, such as taco shells and chop meat or chicken. One thing cannot be consumed without the other. Let's say a purchase of a Keurig Coffee Machine exists within the purchase history of a user, to aid in increasing one's sales, recommendations can be made for the user to discover new coffee cups or can make it easy for the user to purchase their favorite brand. A common practice with recommendation systems involves making an upsell. This would be recommending a higher end product that ultimately costs more. A sound recommendation system increases sales improves customer experience and provides insight for the customer, the company, and its employees.

I obtained this data from the Kaggle website smoothly and efficiently. The link for the data downloaded a BZip2 Compressed Tar file than contained six CSV files. This was a relational set of files that deals with more than 200,000 Instacart users. Each user had anywhere from 4 to one hundred of their orders. Patoolib is a library that makes extracting data from an archive such as, tar.gz simple. It consists of seven methods to be used and is easily controlled through the os library from python. I used Patoolib to deflate the files from its compression and create data frames. Initially, I created a function that produced the data frames for these six files, but due to assigning an index, I lost the order_id column in one of the data frames. This prevented me from combining all the files and linking user_id with order_id. After discovering this problem, after reviewing the data again on Kaggle's website, I decided not to use a function and create the six data frames individually. Next came inspecting the data. There were not any null or NaN values that required attention. NaN values existed for first purchases within the days_since_prior column. There was missing information regarding aisle and department for some products, but luckily, they have been filed as missing within the aisle and department data frames. To combine the data into one, more manageable, data frame, the

Project: Capstone Project 1: Milestone Report

first data frames combined was the prior and train files. This completes the overall sample and includes all of the orders sampled together in one data frame. The columns consist of order id, product id, add to cart order and reorder and were named full. To eliminate the product id and replace it with a recognizable name, first the aisles and departments were merged with the products data frame. The aisle id and department id were also replaced with the recognizable names of each aisle and department associated with the product name and the numeric identifiers were eliminated. Next, the products data frame was merged with the full data frame to remove the inconspicuous numbers and allow a person to read the names of the products, aisles and departments instead of some arbitrary numbers. The next merge performed allowed the user id to be matched with the order id. This allowed me to give context to each order and each user. Aside from the order id and user id, the new columns included in the full data frame consisted of the days since prior order, the time of day the order was purchased, the day of the week the order was purchased and identified which file this information came from, it could have come from the prior history data or the train data files. It was essential to keep the eval set column for later use within machine learning. After some more thought, I will include even more information to the train data set. This information will be saved for an explanation in a future report. Now that one data frame was constructed, with non-relevant id numbers eliminated and names included to identify instances, exploratory analysis can be completed, and a story can be told from the data.

There were many things one could infer from the data. For example, it was necessary to start simple with some simple questions. How many products do we have to choose from, aisles and departments? There are forty-nine thousand six hundred and eighty-eight available to customers from one hundred thirty-four aisles and twenty-one departments accessible. There is a total of two hundred six thousand two hundred nine customers and three million four hundred twenty-one thousand eighty-three orders. Included in the information, one thousand two hundred and fifty-eight products are not assigned to an aisle or department. There are seventy-seven thousand three hundred ninety-six instances of these products being purchased. The aisles and departments are not completely important, but they should be considered in future exploration. To compare the aisles and departments concerning the products available a bar chart (Fig 2) was included with color-coded instances to determine popular aisles or departments these items come from. This bar chart showed the popularity of each department and what aisle was involved. Missing data was the most populated, but that is a reasonable consideration. The top twenty items purchased (Fig 3) were computed and it seems that organic items are trendy, making up half of the top ten purchased items list. Bananas and organic bananas were by far the most purchased items and had the highest probability of being selected. I was interested in whether or not the most frequent items purchased were also the most frequent items reordered (Fig 4). One could deduce this from speculation, but it can be proved with the data. After all, with almost half

Project: Capstone Project 1: Milestone Report

of the items on both lists, it should be tested to see if there is a difference. It was also essential to look at the items that were chosen first. (Fig 5). This information was included in the add to cart order column. Were the items reordered and purchased the most also the first pick for a customer? Again, there were similar items, but I will hold off on more exploration until the statistics portion of this analysis. Fifty-four percent of the items available have been reordered before in this sample, and the rest have not been reordered. (Fig 6) I think the reorder column and the order purchased is critical in informing us about what should be recommended to the customer.

Continuing exploration, there was a mean of sixteen orders per users (Fig 7) and a mean of ten items per purchase. (Fig 8). This can give us a basket size to think of for later use, depending on how we decide the recommendation system will work. An additional column had not been analyzed in terms of distribution. Days since prior column had a mean of approximately sixteen and had a sharp increase at the end of its distribution. This was due to the customers who no longer purchased after thirty days. (Fig 9). Next, a scatter plot was used to determine if there was a relationship between the number of days since the last order and the basket size of the order, in other words, how many items were bought in relation to the number of days it had been since the last order.(Fig 10). This plot distribution showed a fan-shaped distribution and a noticeable non-linear relationship. An appealing aspect of the data to measure is the relationship of time with orders and the items purchased. A heat map was created to visualize the relationship of the time of day and the day of the week these orders took place. (Fig 11). This showed that most of the orders occur within the hours of seven in the morning to about eight at night. I would consider these the peak hours meaning that the ten percent threshold of orders were crossed.(Fig 13). Most of these orders occur early in the week.(Fig 12). One noticeable difference is the time of day on the weekend in which these orders take place. Perhaps because people like to sleep in on the weekends, most orders made occur in the afternoon rather than the morning. The frequency of such orders does increase on the weekends. Saturday has the most, Sunday has second most and Monday has the third most orders during the week. Perhaps this is when the customers like to stock their refrigerator for the week.

After some exploratory data analysis, one type of product was showing up frequently. Organic food seemed to be trendy and considering produce was the most popular department, this idea had to be explored further (Fig 14) to test if there was a difference in reordering rate based on a product being organic or non-organic, a Bayesian approach was chosen to be the most efficient approach. Probability within this sample would allow us to quantify our beliefs on how different organic and non-organic products are. In the initial step to using a Bayesian method, a beta distribution was created. Looking at the situation as binary, alpha was used as the products that were organic and beta as the non-organic products. Samples were generated from this data

Project: Capstone Project 1: Milestone Report

one thousand times using Numpy's Beta method. Once the distributions were plotted from the samples it was easy to see that the graphs of the two products never converged. (Fig 15). This was evidence enough to conclude that these two products had a significant likelihood of having different reorder rates. The distribution of each graph was input into a histogram to endure its normal distribution.(Fig 16, 17). A binomial distribution was not needed to multiply with the beta distribution. To get a percentage of likelihood, the two samples were subtracted from each other and the mean was found of the difference. Because of the large difference in the distributions, every value was over zero and therefore a one hundred percent confidence interval could be established.

Another hypothesis was tested using a Bayesian method. In this case, the data was not binary, and a different distribution needed to be used. To trade off the distributions used, instead of a beta distribution, a Dirichlet distribution was used. To trade the binomial distribution for the distribution of the likelihood, a multinomial distribution was used keeping in mind the conjugate distribution. The hypothesis tested was if the top thirty items purchased were also the top thirty items reordered. Again, dealing with likelihood, a Bayesian method was more appropriate. First, a Dirichlet distribution provided samples. (Fig 18). This could have been done multiple times, but there was enough evidence to use just one instance. Also, because the graphs of the data converged so closely, it was evident that this was very likely. The multinomial distribution samples were calculated, and its values counted to establish the likelihood of the equation. (Fig 19). The Dirichlet prior was then multiplied with the multinomial distribution, likelihood, and confidence was set. It was found that only forty percent of the values were over zero maintaining the null hypothesis that there was no difference in the top items purchased and top items reordered. There was little difference between the two. (Fig 20).

Two other hypotheses test were performed using Frequentists statistics. These two hypotheses tested if the order in which the product was picked affected its reorder rate. Items that were selected first were to be tested first, and then any order number was examined. The means of each sample were calculated, and its standard deviation computed by pooling together the variance. Once the standard error could be computed, a T-value was chosen from a T- table and used in the two-tailed T-test. The confidence interval for both hypotheses tested did not contain zero, and therefore, the null hypotheses were rejected. It was found that item's chosen order does affect whether they are reordered.

Obtaining the data for this particular project was easy compared to my project regarding NYC parking tickets. There is so much to continue to explore but I am now looking forward to using the information gathered during the EDA to create an efficient model. A sound recommendation system increases sales improves customer experience and provides insight for the customer, the company, and its employees. Given a large set of sample data, it will be interesting to continue to explore and improve on the machine learning model to come.

Project: Capstone Project 1: Milestone Report

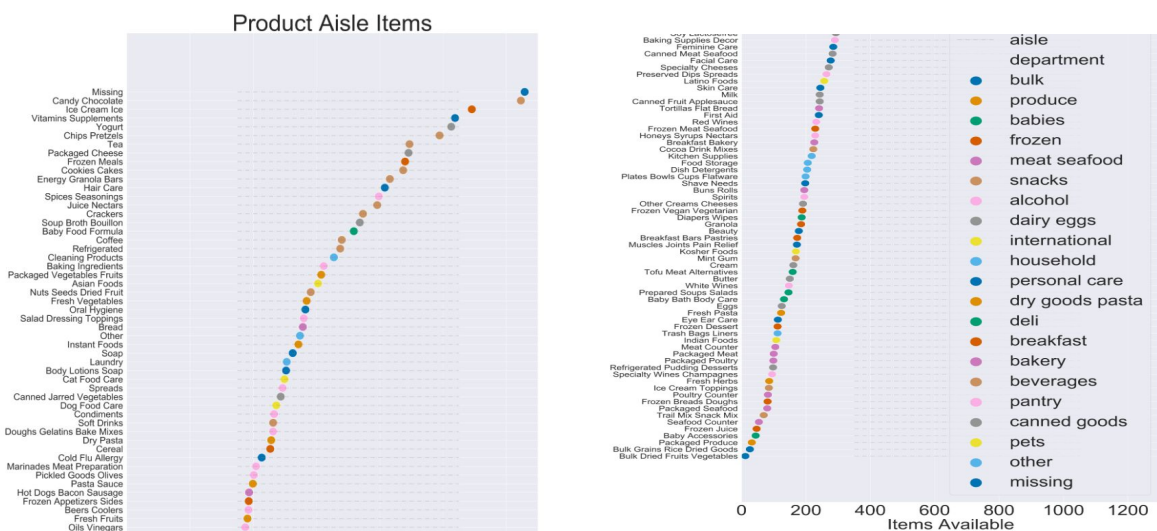


Fig 2

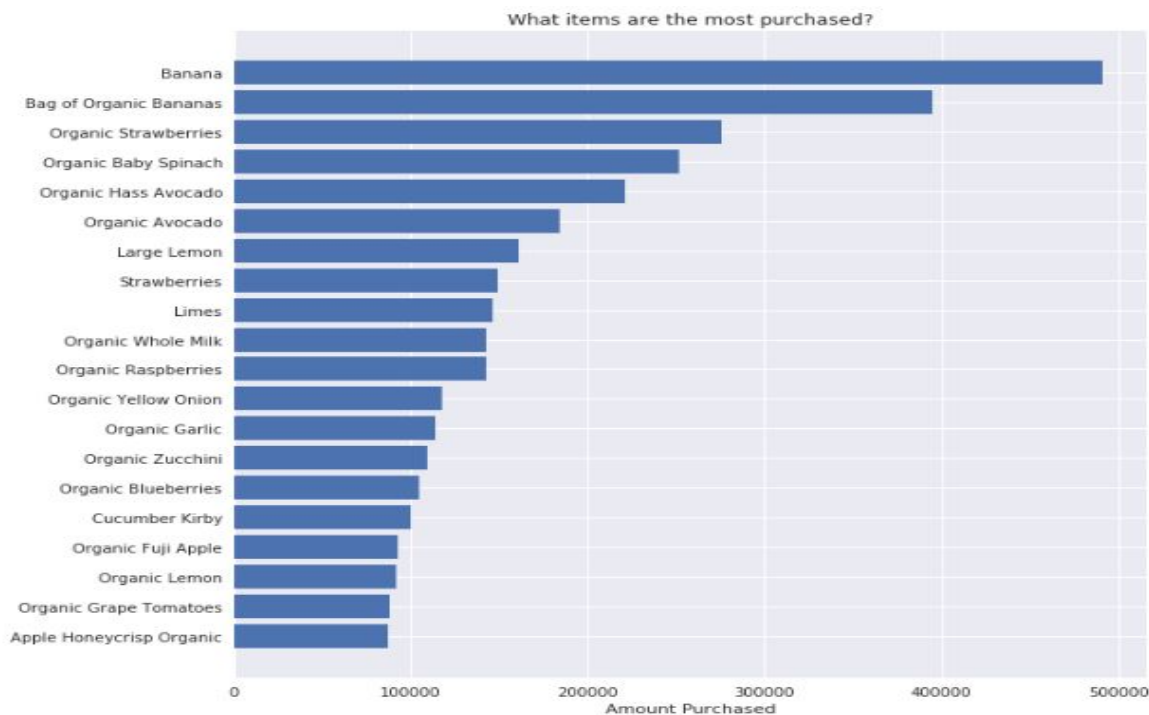


Fig 3

Project: Capstone Project 1: Milestone Report

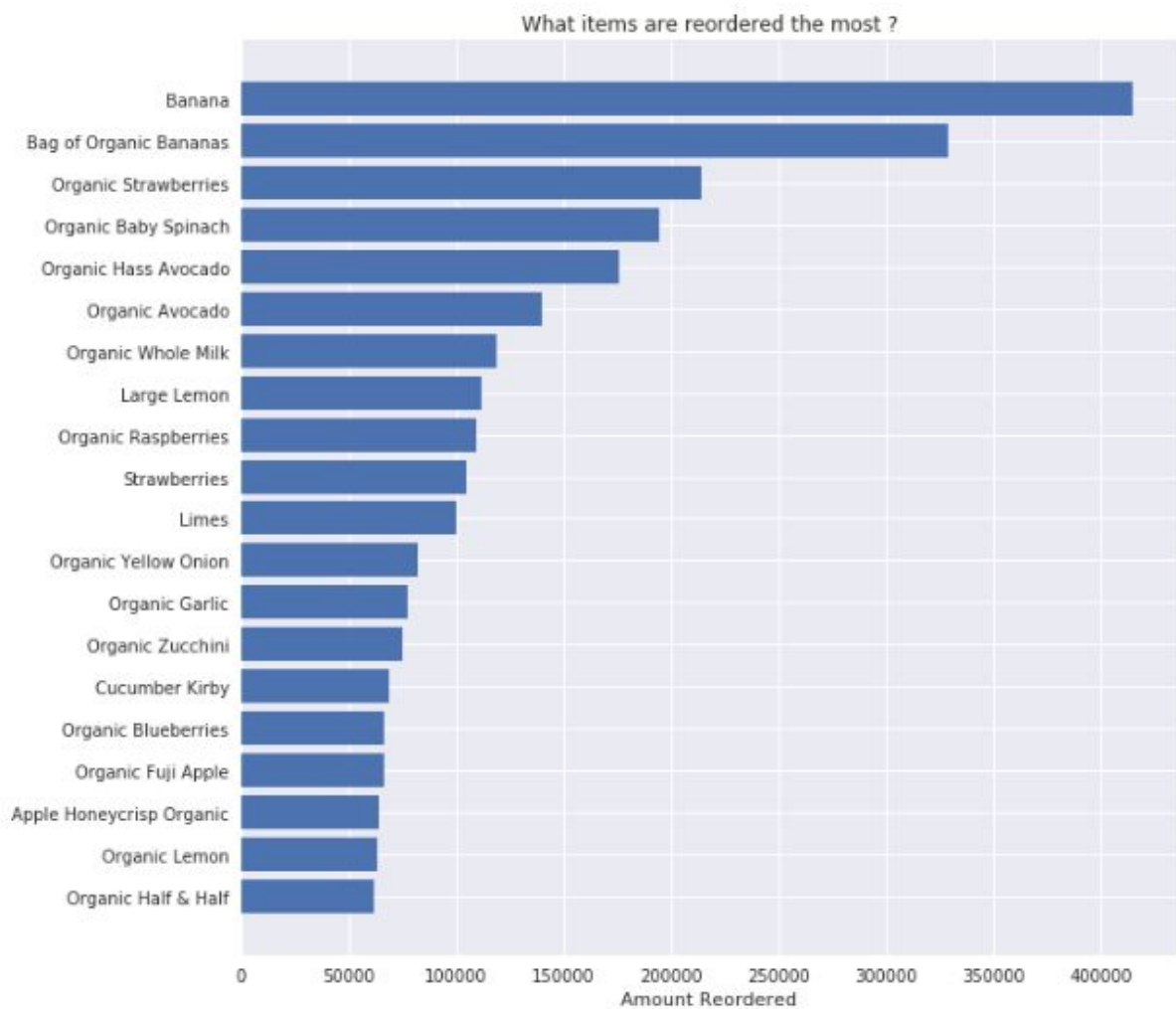


Fig 4

Project: Capstone Project 1: Milestone Report

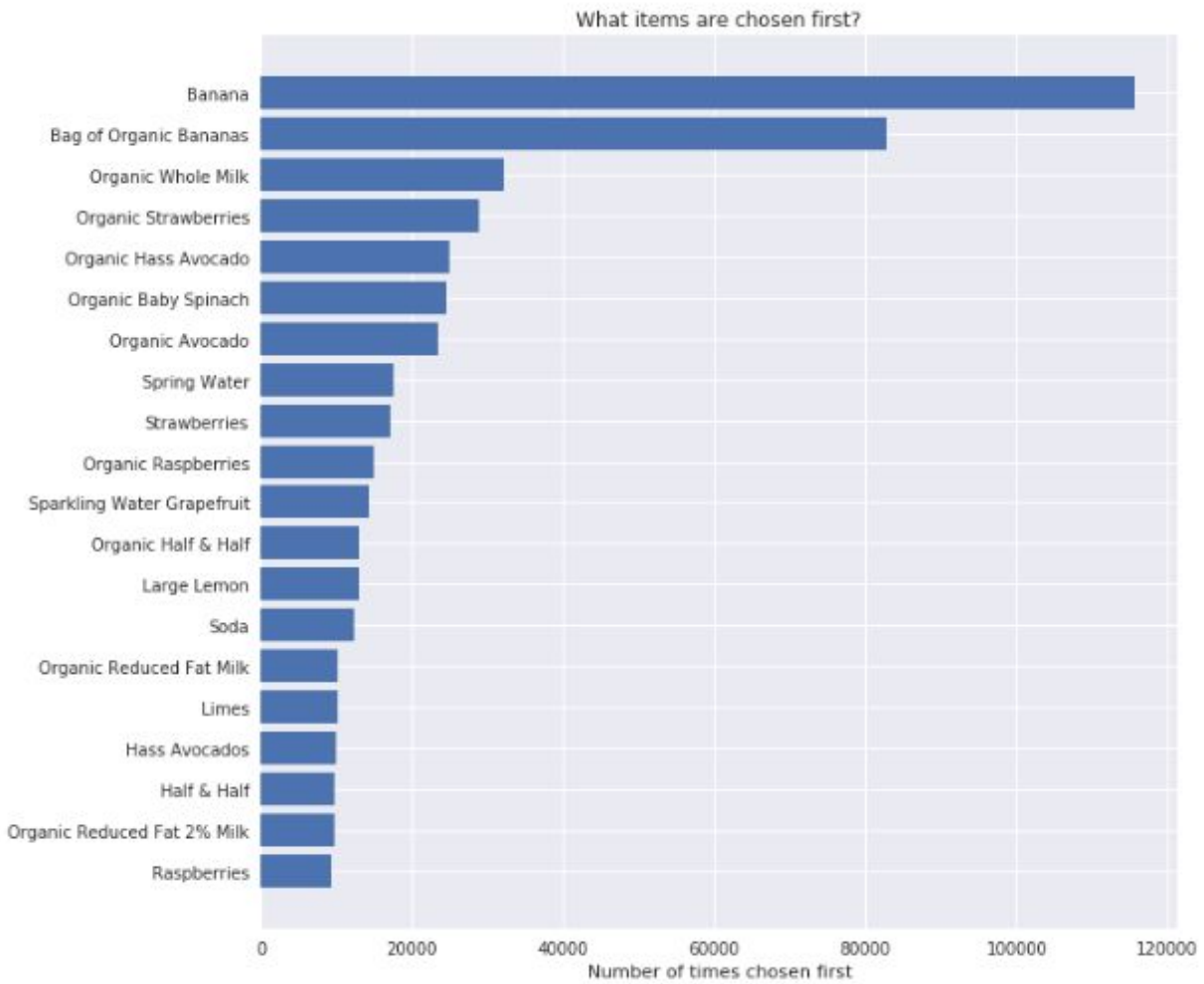


Fig 5

Project: Capstone Project 1: Milestone Report

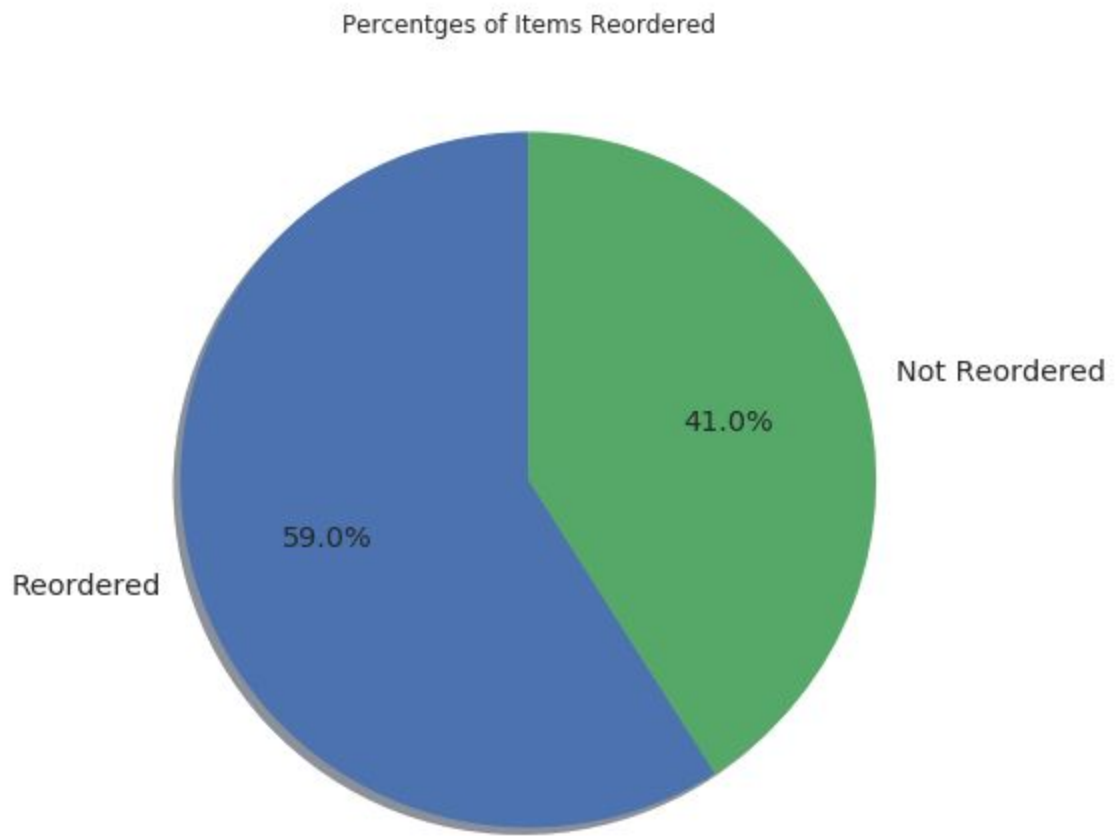


Fig 6

Project: Capstone Project 1: Milestone Report

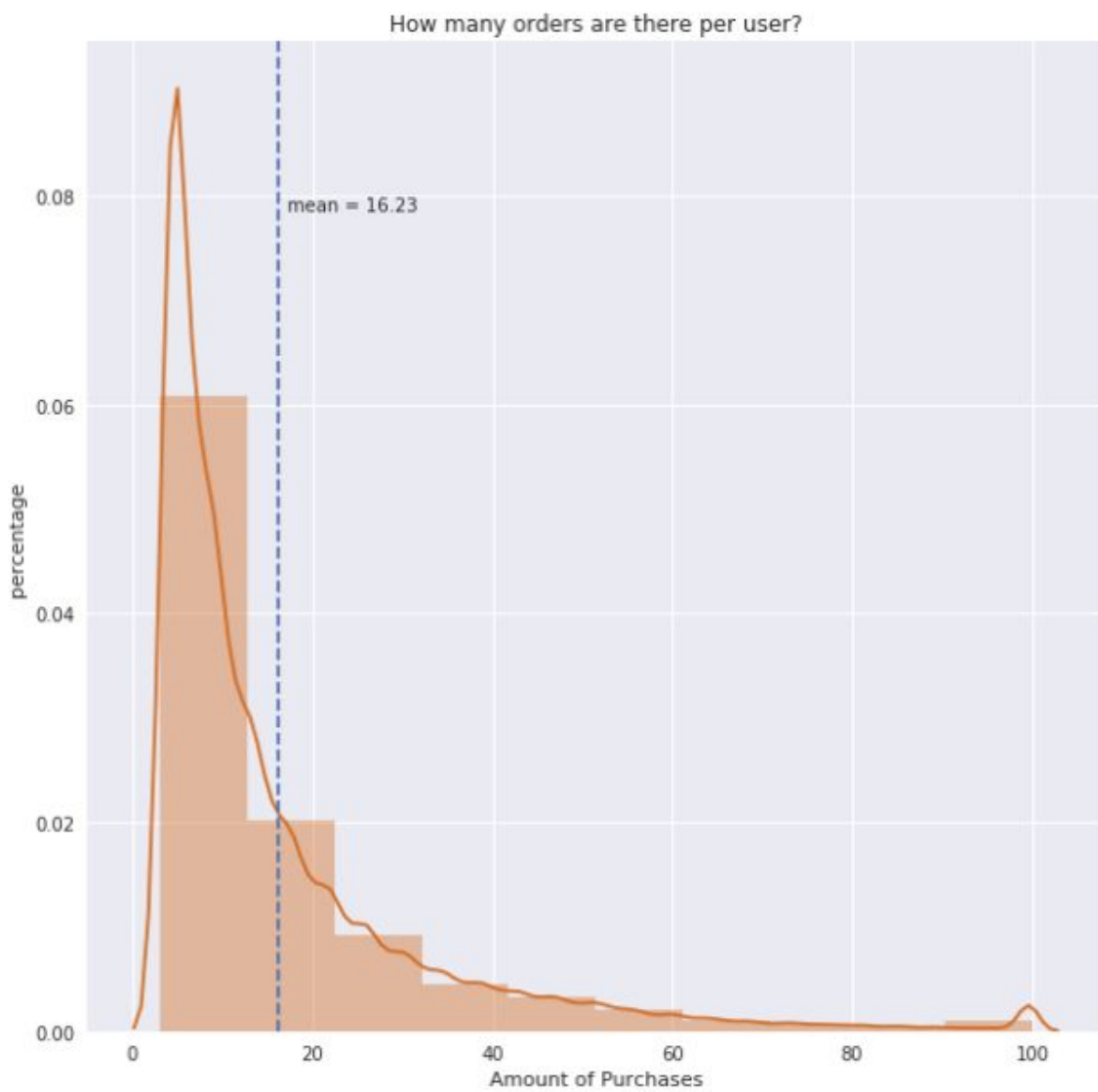


Fig 7

Project: Capstone Project 1: Milestone Report

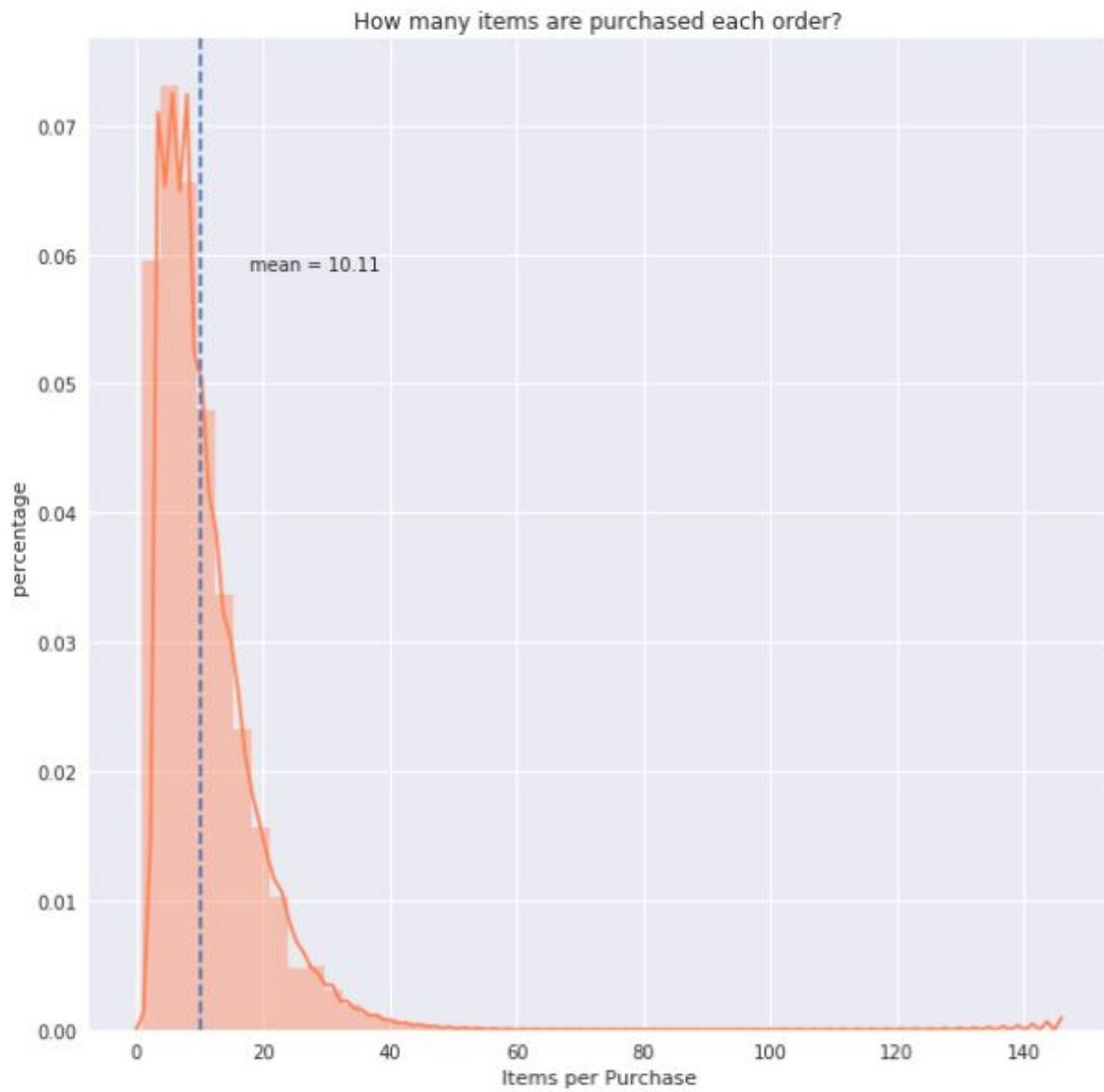


Fig 8

Project: Capstone Project 1: Milestone Report

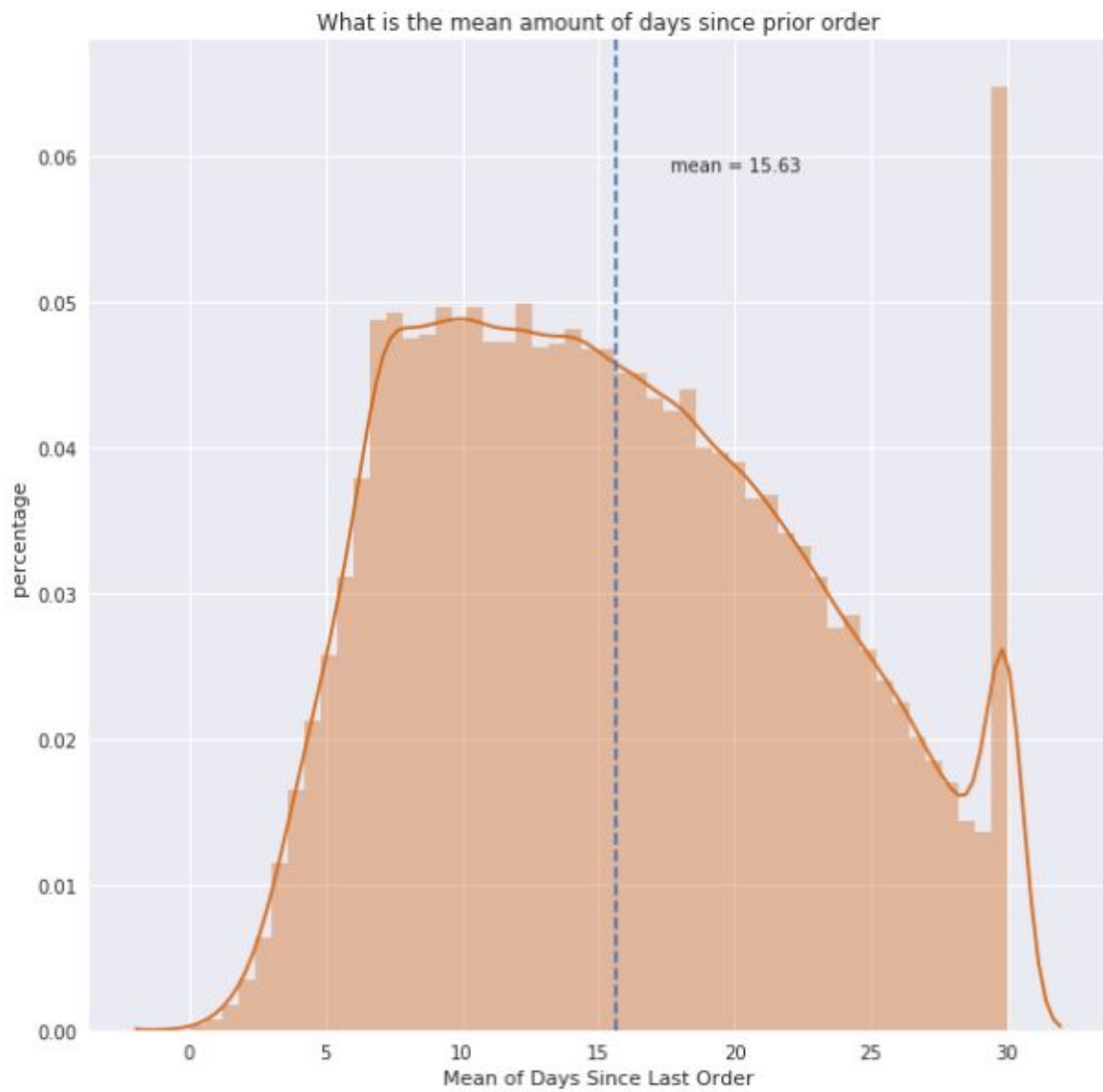


Fig 9

Project: Capstone Project 1: Milestone Report

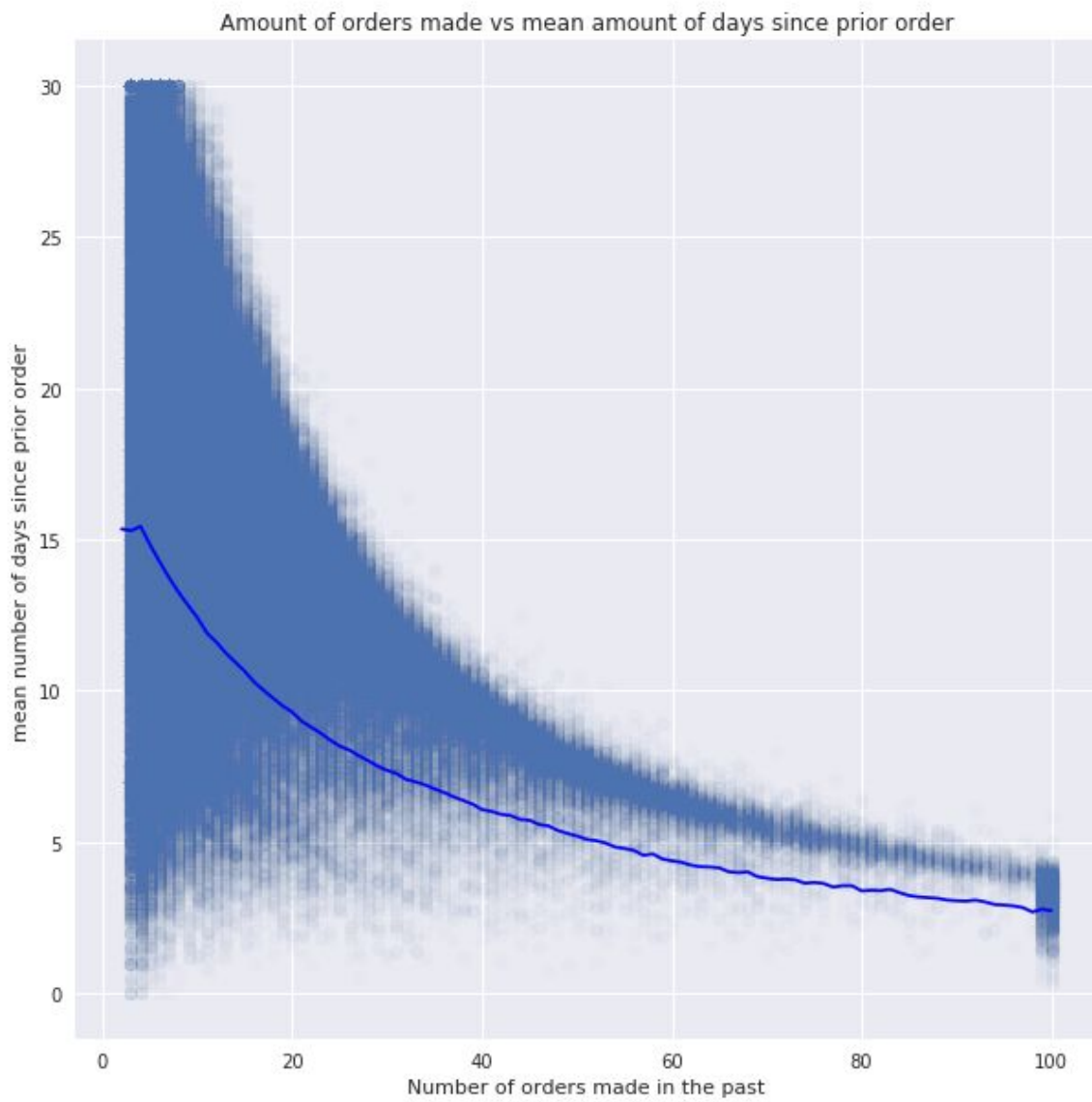


Fig 10

Project: Capstone Project 1: Milestone Report

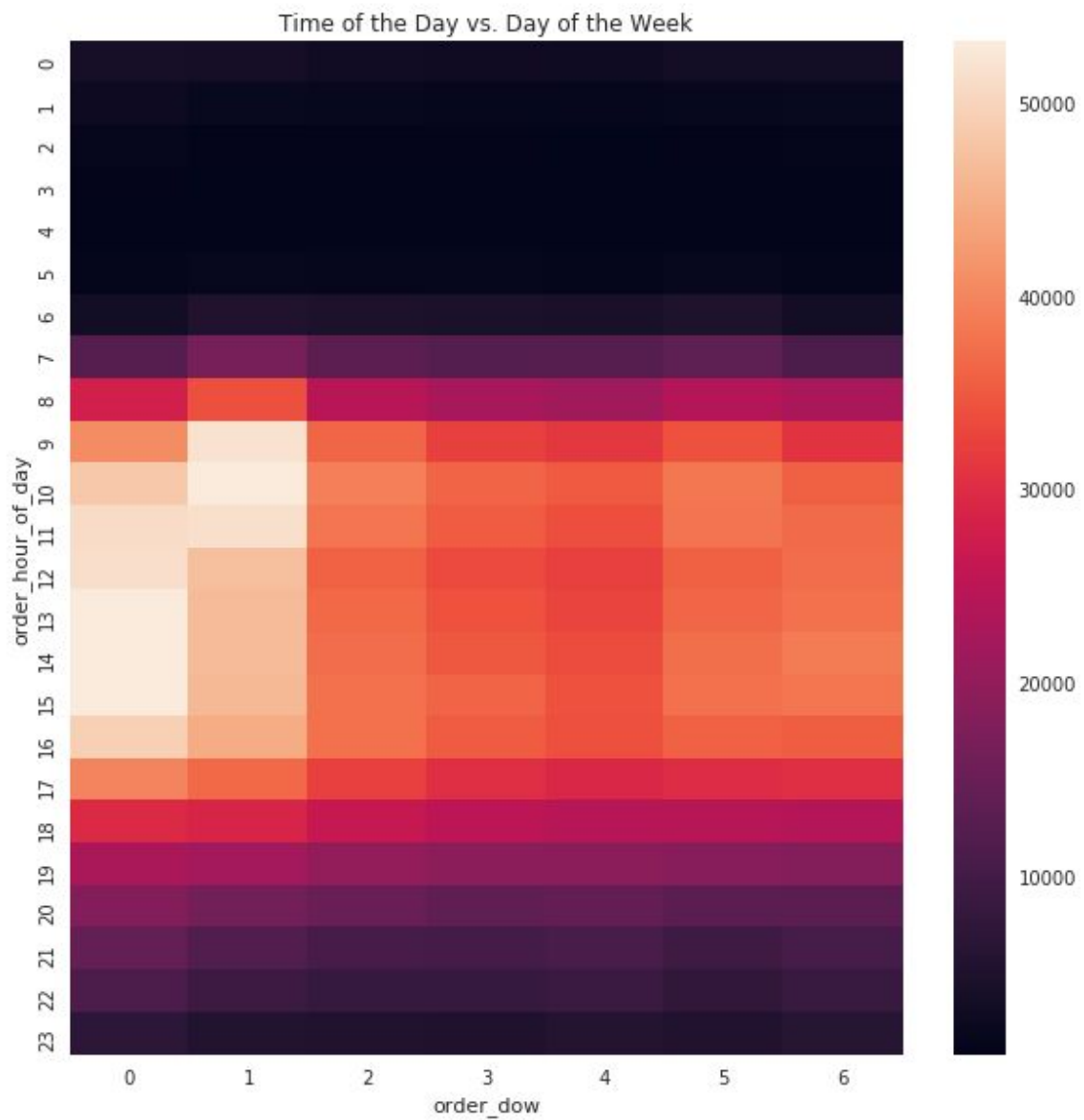


Fig 11

Project: Capstone Project 1: Milestone Report

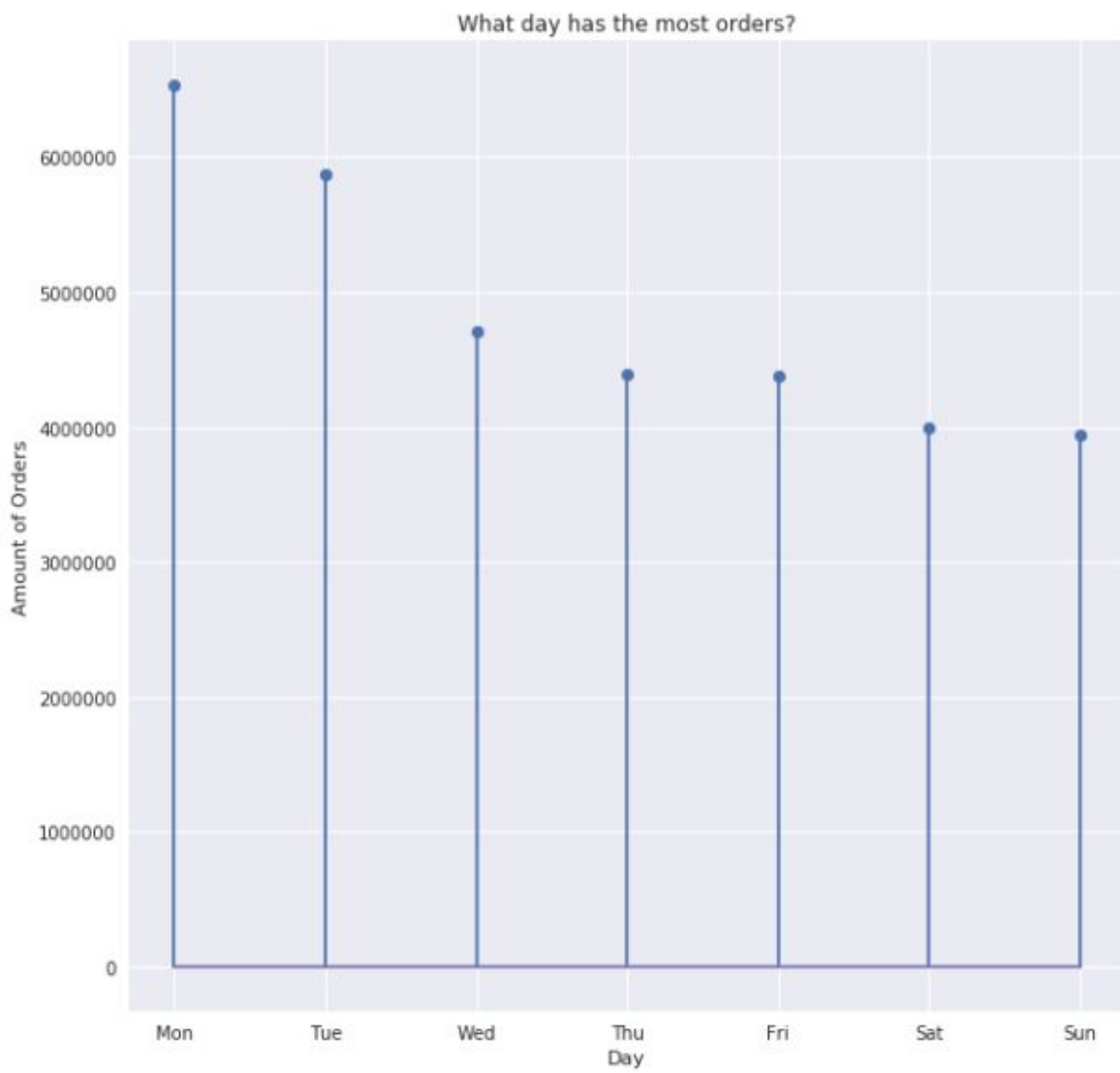


Fig 12

Project: Capstone Project 1: Milestone Report

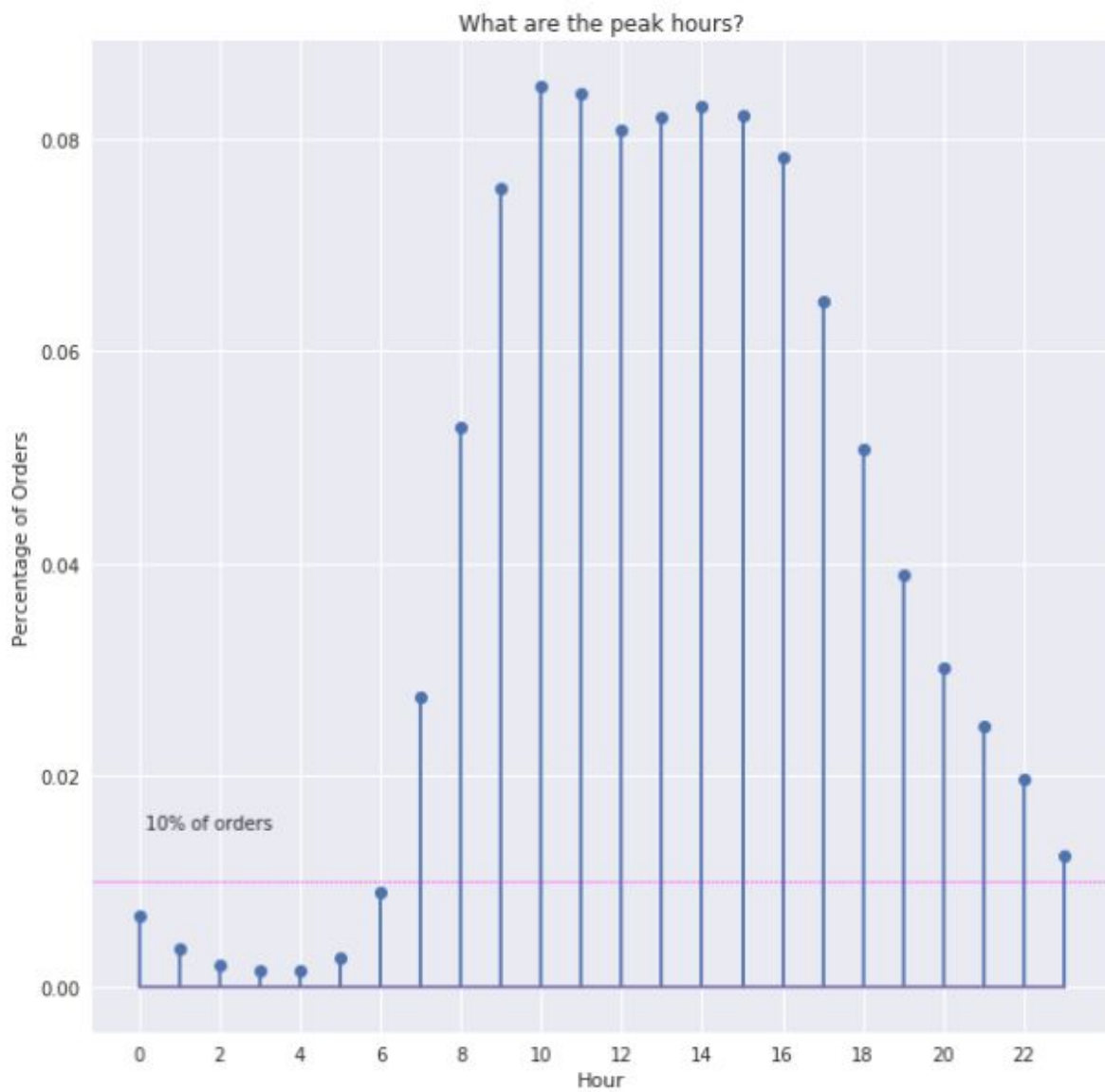


Fig 13

Project: Capstone Project 1: Milestone Report

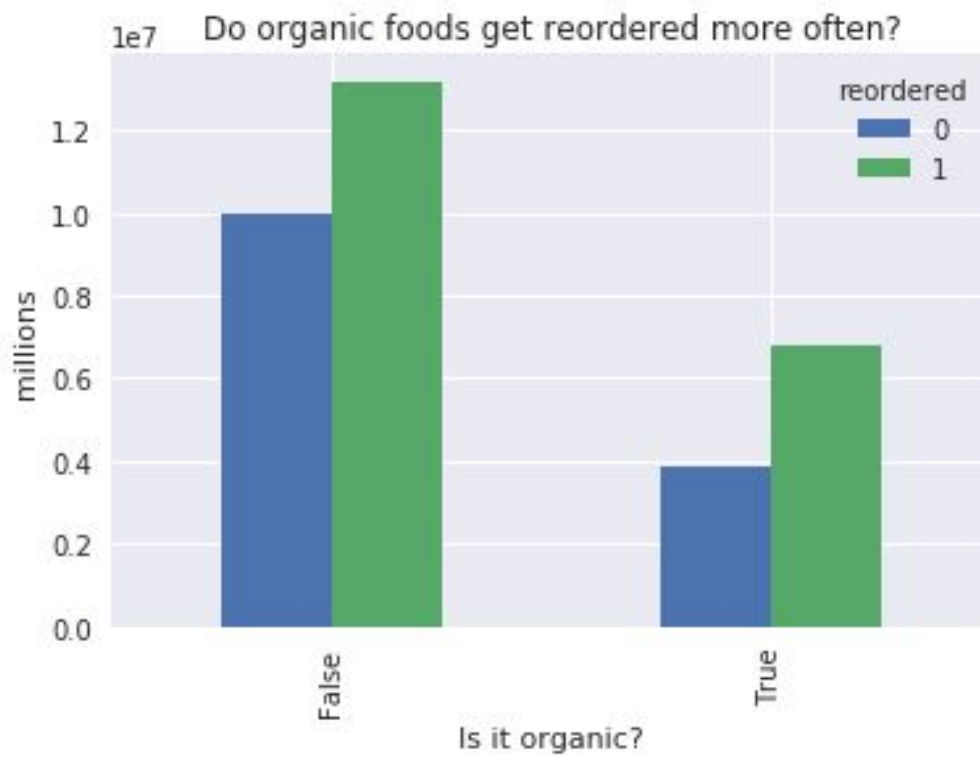


Fig 14

Project: Capstone Project 1: Milestone Report

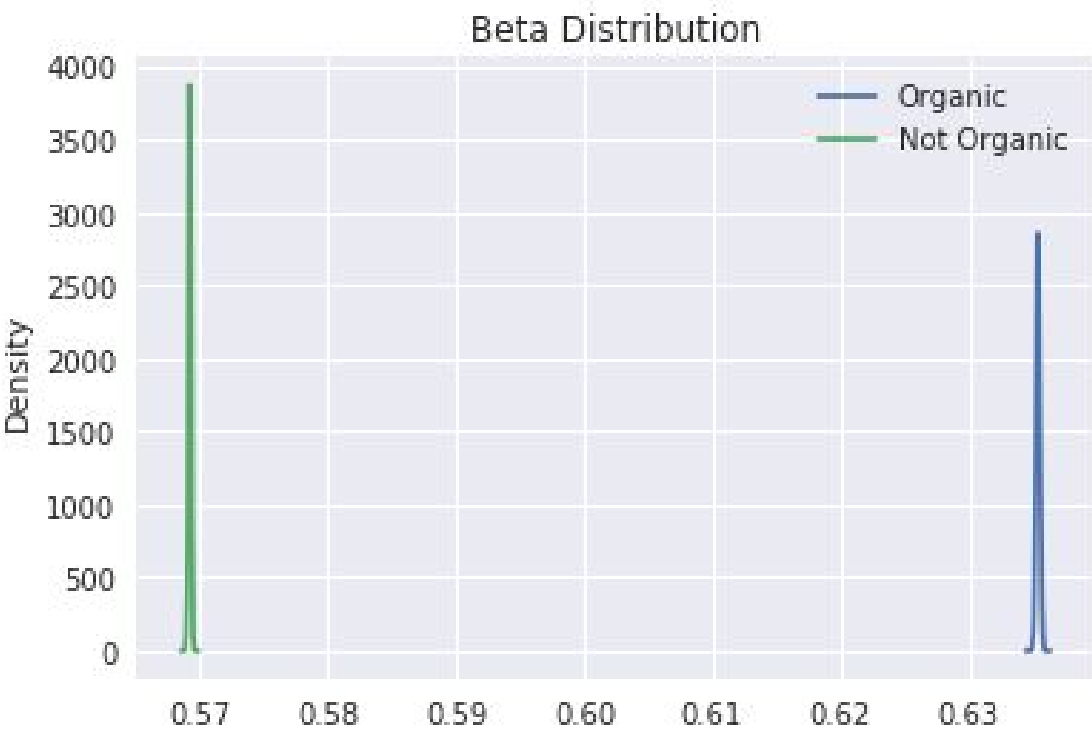


Fig 15

Project: Capstone Project 1: Milestone Report

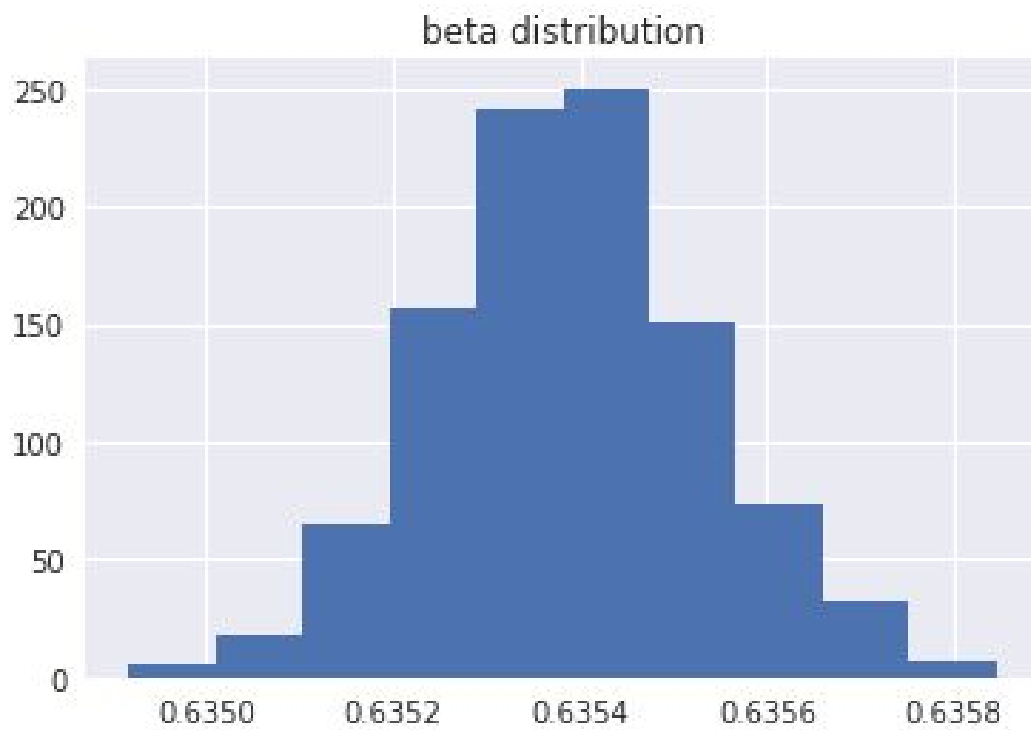


Fig 16

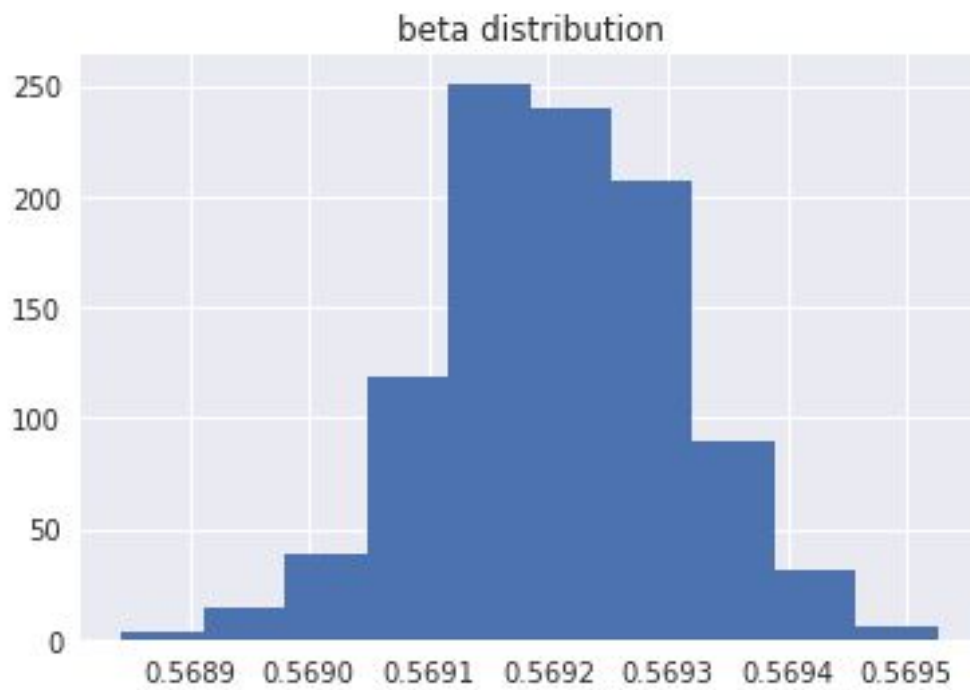


Fig 17

Project: Capstone Project 1: Milestone Report

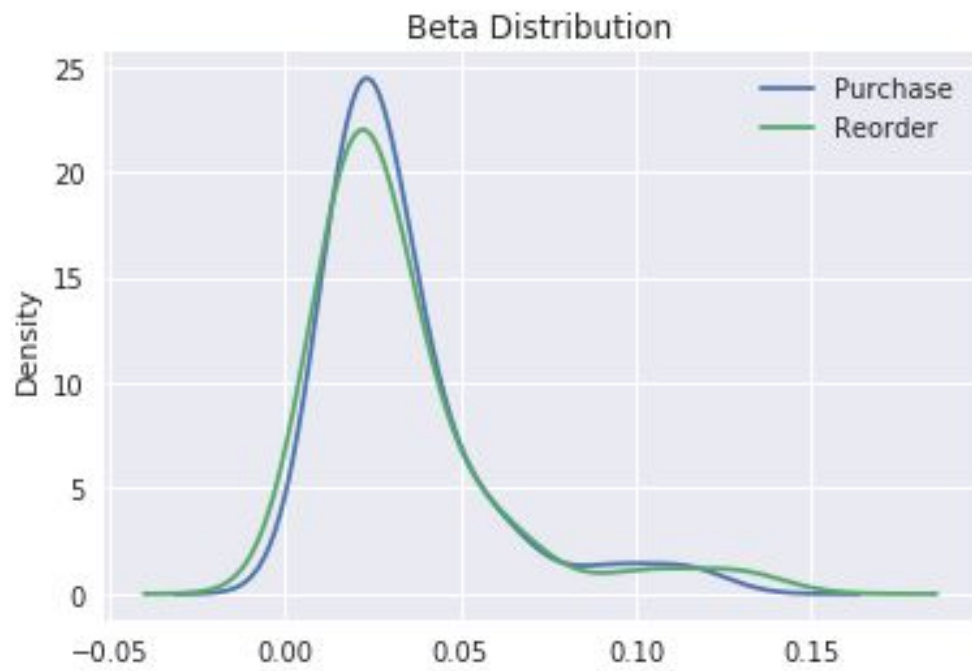


Fig 18

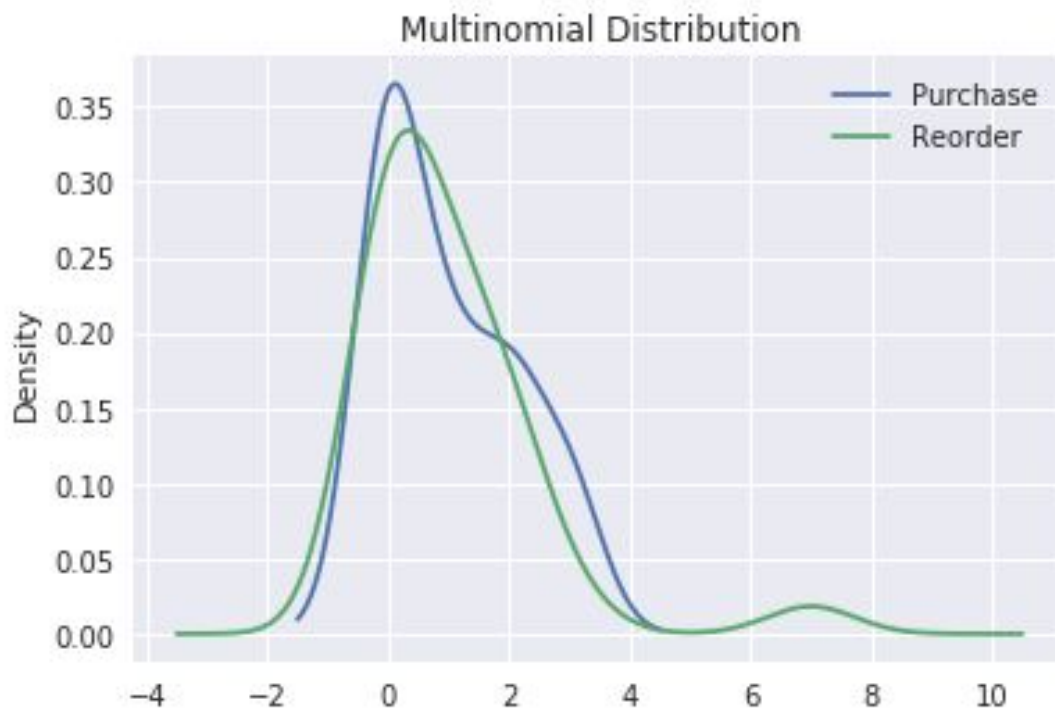


Fig 19

Project: Capstone Project 1: Milestone Report

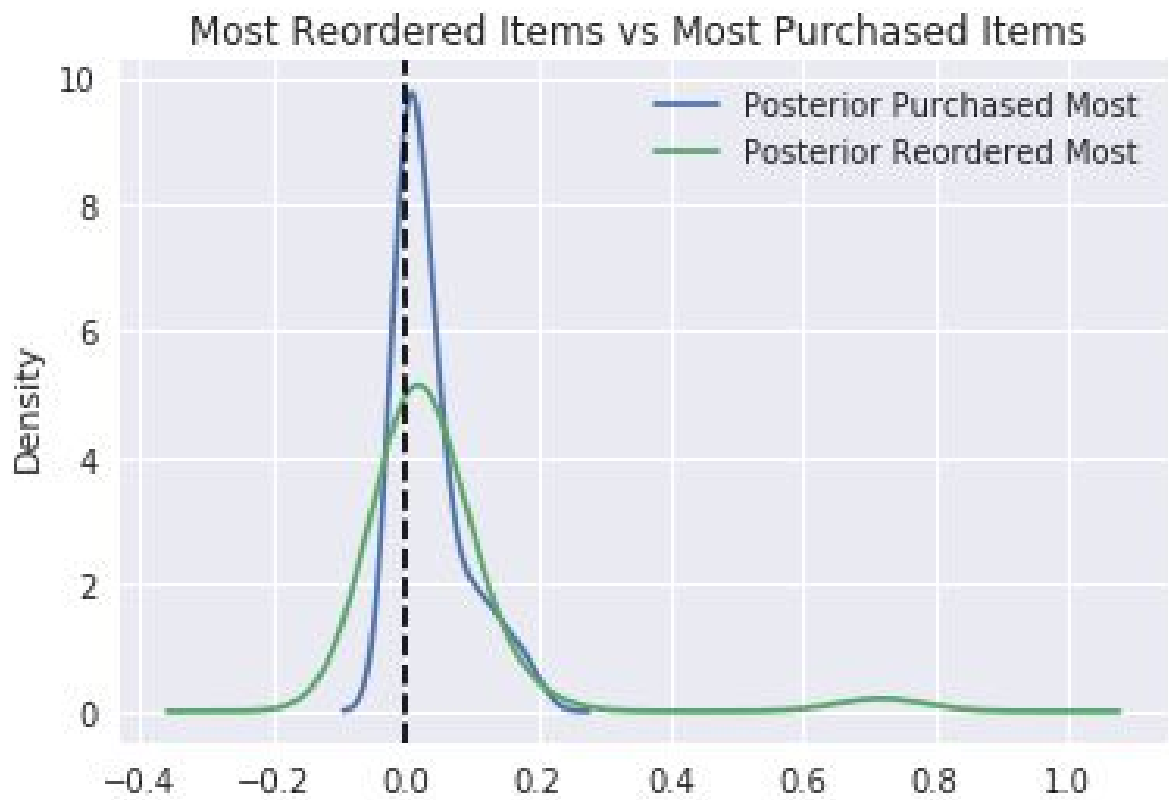


Fig 20