

## Data-Mining Project Write-Up

### Executive Summary

This report explores the data provided by Instacart in 2017 as part of a public competition to predict which products will be in a user's next order. Our analytical approach includes insights such as time-series comparisons of customer, product, and department traffic across working hours and days of the week, feature correlation, and relationship mapping for aspects such as popularity and association, to name a few. The goal is to extract meaningful insights, patterns, and conclusions by exploring customer purchasing behaviors from various facets. Utilizing ggplot2, SQL, and various other standard R libraries, we provide visualizations and analysis, breaking down the data and patterns from a descriptive analytics perspective. All technical work is in the R code uploaded to the group repository, while the data was left out due to file size constraints set by GitHub but can be downloaded here: <https://www.kaggle.com/c/instacart/data>.

### Our Data

Working with the dataset provided by Instacart for their Instacart Market Basket Analysis competition, we started with a relational set of files describing three million grocery orders from more than 200,000 customers over time. These seven CSV files included 25 columns total and were initialized as follows:

- **Departments.csv**
- **Aisles.csv**
- **Order\_products\_\*.csv (train & prior)**
- **Orders.csv**
- **Products.csv**
- **Sample\_submission.csv**

From here, we dropped repetitive and/or unuseful files to our descriptive task. Specifically, we dropped 'sample\_submission.csv' (file specific to the competition hosted by Instacart) and 'order\_products\_prior.csv' since it was a larger version of order\_products\_train.csv thus increasing the computational cost and heavily slowing execution time.

We identified and handled missing values and irrelevant columns to ensure the data was suitable for analysis. Of our remaining five files stored as data frames, missing values were only found in the orders table, specifically in the 'days\_since\_prior\_order' column. Since it could not be deduced whether these missing values in this column represented first-time orders or simply information that was not provided, the result would have caused a reduction of just ~ 6% of the total records out of over 200,000, so we decided to omit them.

Next, we handled some column modifications. We dropped the 'eval\_set' column as it was irrelevant to our analysis and created time-related columns 'order\_day' and 'order\_hour' from the ambiguous descriptions provided (0-6 for day and 0-23 for hour) to add context for understanding customer behavior across different periods.

We merged tables to reduce redundancies. Since the products table was the only one referencing the information provided by the aisles and departments tables, we merged them into products to reduce the number of tables to reference and the frequency of joins.

The resulting data we were to work with now includes 3 data frames with 18 columns total:

- **Orders\_cleaned**
- **Products\_cleaned**
- **Order\_products\_train**

## Analysis

### Visualization 1: Table

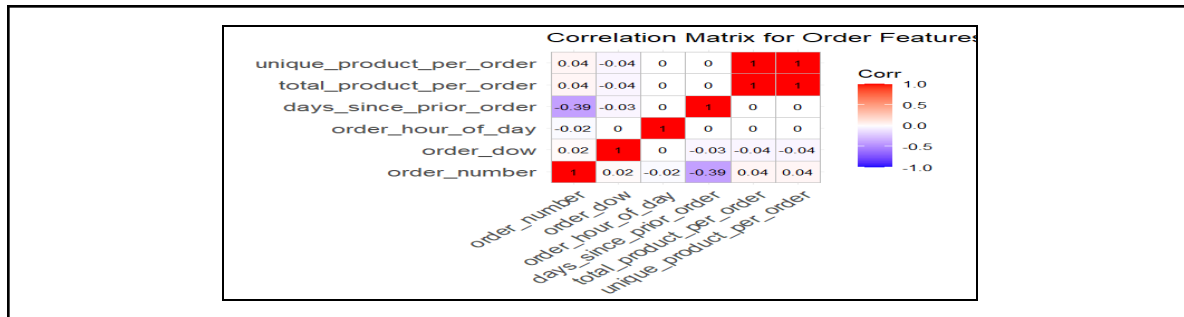
Some initial observations to take note of when observing the summary statistics of our data are listed in the table below. By tracking statistics on each order like the aisle and department food is kept, we can identify trends to enhance the predictive analysis abilities of the app. The average time between orders for all users was just over eleven days but ranged anywhere from later that same day to over 30 days later while the average order size was around 10 with a cart high of over 80 products. This wide spread of values in a number of the impactful variables causes complications and complexity when working with and attempting to gain insights out of their relationships thus requiring specific case handling to workaround.

<pre># Combine important metrics/aggregations from the summary statistics of the different tables into one structure summary_table &lt;- tibble(   Metric = c('Orders', 'Days Between Orders (Avg)', 'Customers', 'Order Size (Avg)',             'Products', 'Departments', 'Aisles'),   Value = c(     summary_orders\$total_orders,     round(summary_orders\$avg_days_between_orders, 2),     summary_orders\$total_customers,     round(avg_order_size, 2),     summary_products\$total_products,     summary_products\$total_departments,     summary_products\$total_aisles   ) )  # Pretty print a table visualization summary_table %&gt;%   kable() %&gt;%   kable_styling(bootstrap_options = 'striped', full_width = FALSE)</pre>	<table><tr><th>Metric</th><th>Value</th></tr><tr><td>Orders</td><td>3214874.00</td></tr><tr><td>Days Between Orders (Avg)</td><td>11.11</td></tr><tr><td>Customers</td><td>206209.00</td></tr><tr><td>Order Size (Avg)</td><td>10.55</td></tr><tr><td>Products</td><td>49688.00</td></tr><tr><td>Departments</td><td>21.00</td></tr><tr><td>Aisles</td><td>134.00</td></tr></table>	Metric	Value	Orders	3214874.00	Days Between Orders (Avg)	11.11	Customers	206209.00	Order Size (Avg)	10.55	Products	49688.00	Departments	21.00	Aisles	134.00
Metric	Value																
Orders	3214874.00																
Days Between Orders (Avg)	11.11																
Customers	206209.00																
Order Size (Avg)	10.55																
Products	49688.00																
Departments	21.00																
Aisles	134.00																

### Visualization 2: Correlation Matrix

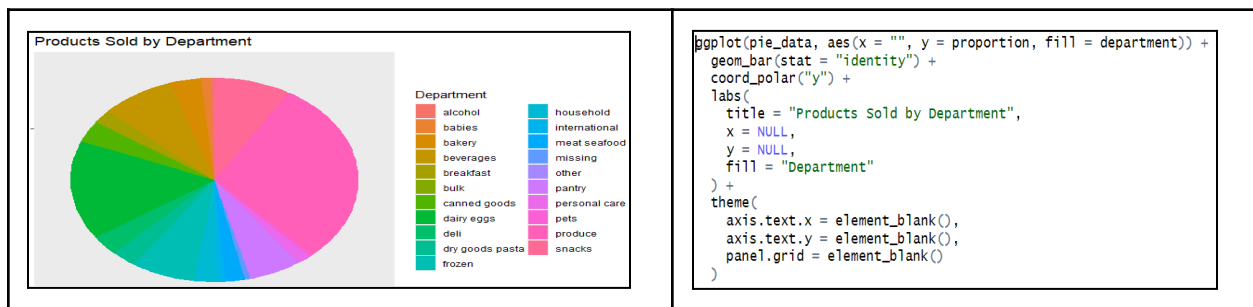
A correlation matrix is intended to visualize the strength and direction of relationships between variables of interest, specifically linear. However, as you can see with our visualization, the numeric values we extracted through the SQL query do not exhibit linear relationships. This can be explained by the nature of the variables we are working with. For example, order\_dow encodes a categorical variable (day of the week), but this as a number simply does not provide meaning when drawing relationships. In other words, through this visualization, we were able to pick out this abnormality, the relationships are more complex and thus may benefit from other visualizations that can map and present these relationships, potentially nonlinear, better than a correlation matrix can.

<pre>corr_matrix &lt;- cor(corr_numeric_data)  ggcorrplot(corr_matrix, lab = TRUE, lab_size = 3, title = 'Correlation Matrix for Order Features',   colors = c('blue', 'white', 'red'))</pre>
---



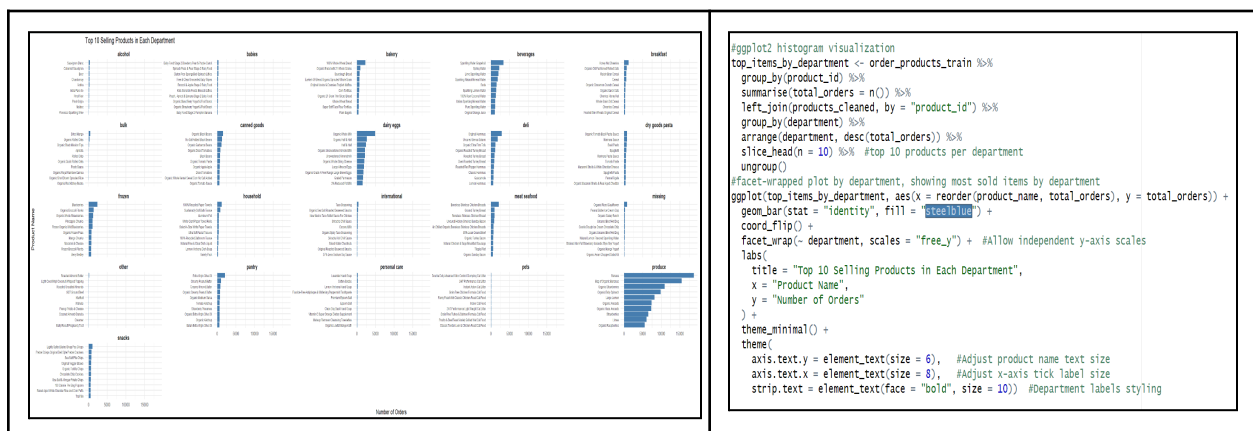
### Visualization 3: Pie Chart

To get an accurate representation of the amount of products being sold in each grocery store department, we decided to make a pie chart. Grouping by department and calculating proportions with respect to other departments, we were able to view with this query and plot that the most popular departments are dairy and eggs, produce, and snacks. Those can be explained clearly by human behaviors. Dairy and eggs are essential for everyday cooking and baking, making them frequent purchases, while snacks are often bought on impulse. It looks like refrigerated goods are the most commonly purchased type of product on Instacart.



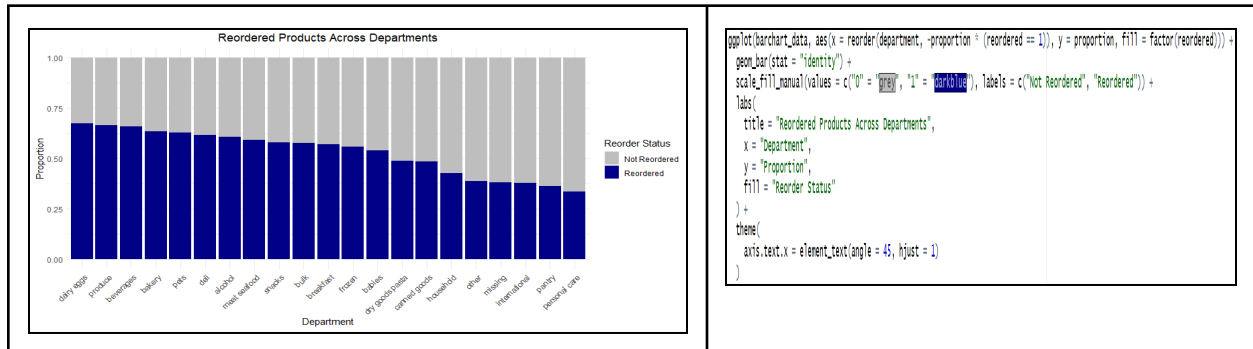
### Visualization 4: Histogram

We created histograms providing us with the top ten most commonly purchased products for each department found in a grocery store. Using `facet_wrap`, we could look at each department's top-performing items and compare them to other departments. You can see that most departments have a favored good, which came as a surprise as we imagined that there would be more competition within departments rather than a stand-out product.



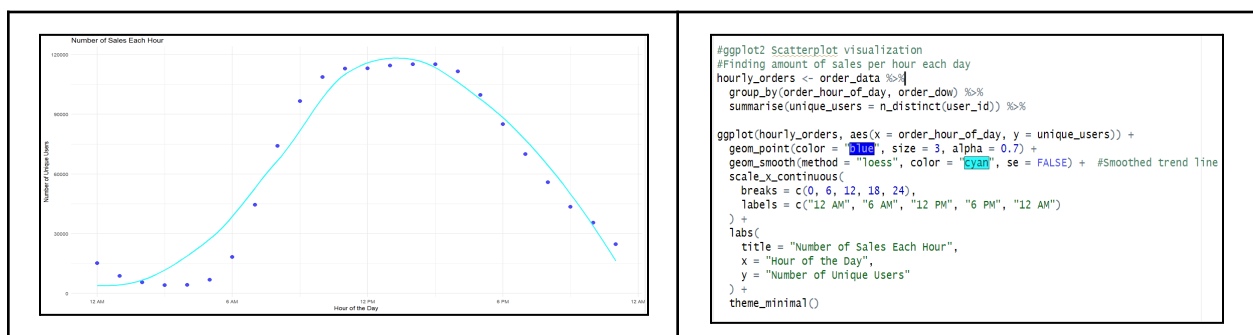
## Visualization 5: Stacked Barchart

By reordering the contents of our SQL query we are able to draw a few relationships out of similar departments. There is a general decreasing trend in the proportion of reordered products as you move from left to right across the departments. Similar to what we mentioned previously with customer purchasing behaviors, customers are frequenting departments centralized on produce. While on the lower half, we see lower reordering rates which can be attributed to lower consumption/turnover like personal care items and miscellaneous things that. The usability of a product along with the rate at which it needs to be replenished appears to have a direct correlation with its reordering potential which is to be expected.



## Visualization 6: Scatterplot

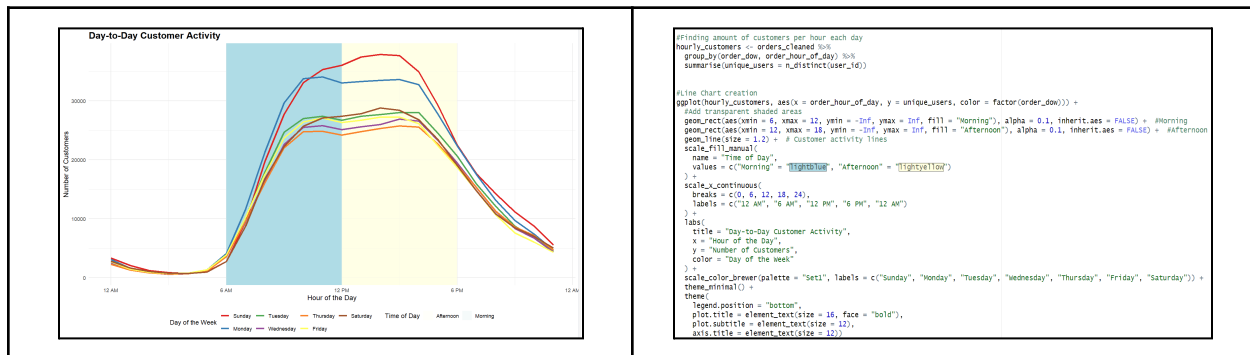
By leveraging ggplot(), geom\_point(), geom\_smooth(), and other operators, we were able to form a scatterplot that compares the number of orders at each hour of the day or in other words, the customer frequency. With geom\_smooth(), we placed a smooth trend line over the plot to identify linear or non-linear relationships. Without a trend line, we would also be unable to see any patterns not shown by the variability in the data. Most orders coming late morning into the afternoon make sense, as people may be shopping for items they will cook for dinner that night.



## Visualization 7: Line Chart

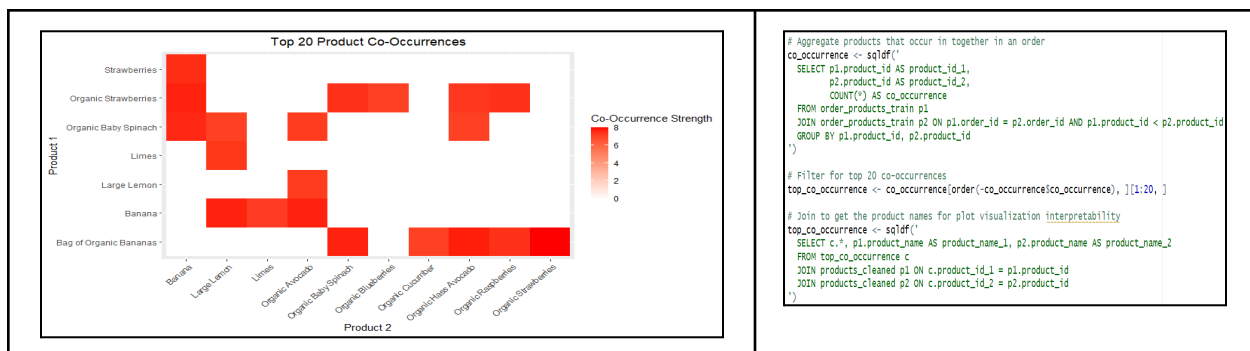
By using ggplot(), geom\_rect(), and geom\_line(), we created a line chart to compare the amount of periodic grocery store activity on each day of the week. Using geom\_rect(), we could distinguish even further the peak times, which are the morning and afternoon. We also distinguished each day of the week by color and used scale\_color\_brewer() to label the numeric variable with the true days of the week. It was expected that Sunday would be the most popular day of the week to grocery shop due to most people not working, and this assumption was

correct. What was shocking, however, was the similarities between Monday and Sunday's activity. Monday's activity proves to be much larger than every day except Sunday.



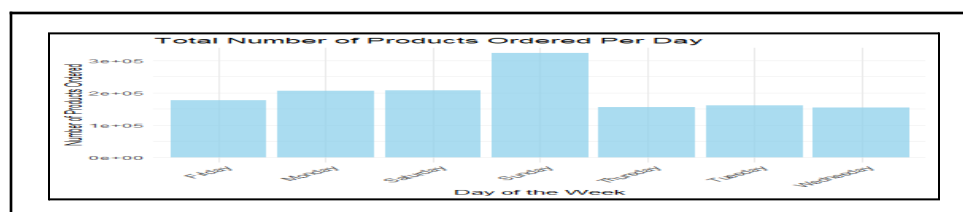
## Visualization 8: Heatmap

The heatmap highlights patterns in customer purchasing behavior, helping to identify which products are often bought together. It helps track customer habits and preferences like common meal combinations. Think of fruit for example, organic strawberries and regular strawberries are often purchased together, showing one of the strongest co-occurrences. This suggests customers balance preferences for organic options with price or availability. It is also offers consideration for arranging items and aisles in such a way that products complement each other and drive positive trends in sales because of customer interest and attention.



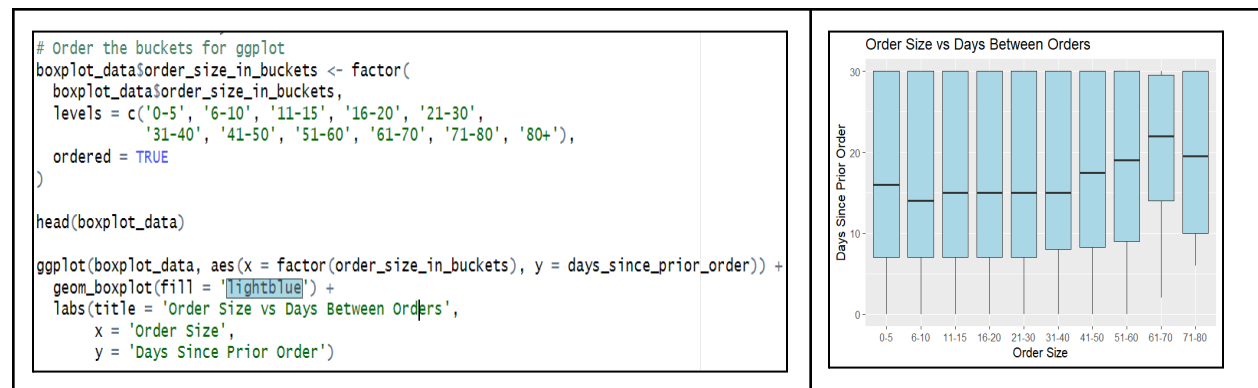
## Visualization 9: Bar Chart

Is there potential for providing deals and recommendations on specific days to help further drive sales? These insights can help Instacart predict a user's next purchase by leveraging specific shopping patterns tied to the days of the week. For Tuesdays and Wednesdays, when many users are busy and non-active, Instacart can suggest smaller purchases, like snacks or quick meal ingredients, and offer reminders or midweek discounts as encouragement. There are obvious advantages of knowing the high-traffic days like the weekend.



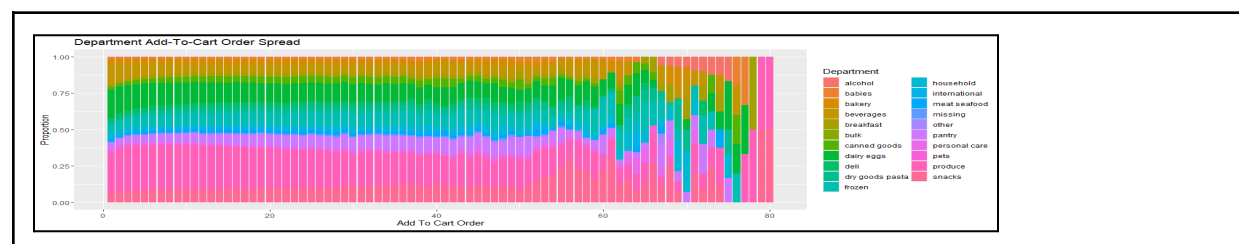
## Visualization 10: Boxplot

Instacart can predict that users who place smaller orders will shop more frequently, often within 1-2 weeks, and can suggest replenishment items. For users who place larger orders, Instacart can expect less frequent purchases, about every 3-4 weeks, and recommend bulk or pantry staples. The visual highlights the median, showing that larger orders generally have longer intervals between purchases.



## Visualization 11: Stacked Histogram

This stacked histogram illustrates the proportion of items added to the cart across various grocery store departments as order size increases with each color representing a specific department. Larger orders (to the right) include a broader mix of items like snacks (pink), pantry (red), and beverages (blue), indicating more diverse shopping habits. Lower cart order counts tend to follow similar trends in terms of item positions with departments while higher counts tend to be more random and sporadic possibly due to the intent of purchase (for an event, you might get many random items that aren't normally coupled but still required) and range of products selected.



## Visualization 12: Treemap

To create the Treemap visualization, we first needed to install the Treemapify package. Then, we could left-join `order_product_train` to `products_clean` to analyze the sales per department in stores. This plot did not bring many surprises, as it was expected that produce and dairy/eggs would be frontrunners in sales as those items are extremely popular and abundant in United States grocery stores. We were surprised at how small of a share of the plot that personal care got however, because this is an essential department that almost every person uses. This could be because people may turn to alternate stores, such as enterprises which specialize in personal care items, rather than Instacart.



## Queries

### Query 1: Summary statistics of the cleaned dataset

We ran queries to analyze the patterns in our datasets once they were cleaned. In this query, we summarized all of the variables in the dataset.

<pre>summary_orders &lt;- sqldf(' SELECT COUNT(*) AS total_orders, AVG(days_since_prior_order) AS avg_days_between_orders, COUNT(DISTINCT user_id) AS total_customers FROM orders_cleaned ')  avg_order_size &lt;- sqldf(' SELECT AVG(order_size) AS avg_order_size FROM (   SELECT order_id, COUNT(*) AS order_size   FROM order_products_train   GROUP BY order_id ) ')  avg_order_size_val &lt;- avg_order_size*avg_order_size  summary_products &lt;- sqldf('SELECT COUNT(*) AS total_products, COUNT(DISTINCT department) AS total_departments, COUNT(DISTINCT aisle) AS total_aisles FROM products_cleaned')</pre>	<pre>summary_products total_products total_departments total_aisles 49688          21          134 summary_orders total_orders avg_days_between_orders total_customers 3214874      11.11484      206209 avg_order_size avg_order_size 10.55276</pre>
--	---

### Query 2: Correlation coefficients between numeric variables of interest

For this query, we extract all relevant numeric variables to be put into a correlation matrix in order to assess the coefficients for potential linear relationships along with the strength and direction of them if applicable.

<pre>customer_order_sum &lt;- sqldf(' SELECT order_id, COUNT(*) AS total_product_per_order, COUNT(DISTINCT product_id) AS unique_product_per_order FROM order_products_train GROUP BY order_id ')  # Combine numeric columns from orders_cleaned and order_products_train to visualize # any linear relationships via correlation corr_numeric_data &lt;- sqldf(' SELECT o.order_number, o.order_dow, o.order_hour_of_day, o.days_since_prior_order, s.total_product_per_order, s.unique_product_per_order FROM orders_cleaned o INNER JOIN customer_order_sum s ON o.order_id = s.order_id ')</pre>	<pre>&gt; head(corr_numeric_data) order_number order_dow order_hour_of_day days_since_prior_order total_product_per_order unique_product_per_order 1      11      4      8      14      11      11 2      15      1      11      30      31      31 3      5      0      11      6      9      9 4      21      2      11      6      9      9 5      4      1      14      10      18      18 6      4      6      10      30      22      22</pre>
--	--

### Query 3: Comparison of department popularity

In this query, we want to look at the proportion of department popularity in percentage terms and products sold to better grasp department share. Our findings make sense as produce and dairy/eggs are large departments with many customers. It is surprising how little of a percentage of orders that the smaller departments get, such as pets and bulk.

<pre>pie_data   department_id  department  products_sold  proportion 1      12      produce      49688      0.385433775 2      1      dairy eggs      118882      0.085446776 3      14      frozen      118046      0.085256495 4      13      bakery      100426      0.075298043 5      12      pantry      82242      0.0635747093 6      12      canned goods      48394      0.0349511814 7      12      meat seafood      48799      0.037392383 8      12      breakfast      38713      0.0279593562 9      9      dry goods pasta      4791      0.003698071 10     12      household      30307      0.021883634 11     14      breakfast      29520      0.021035107 12     12      personal care      23170      0.015783152 13     12      baby      24943      0.0107907096 14     12      pet care      21701      0.016819675 15     21      misc      8253      0.0058590486 16     5      other      1598      0.0040439053 17     8      bulk      4497      0.003478295 18     5      other      3755      0.002863674 19     10     bulk      1115      0.000834988</pre>	<pre>pie_data &lt;- sqldf(' SELECT p.department_id, p.department, COUNT(o.product_id) AS products_sold FROM order_products_train o JOIN products_cleaned p ON o.product_id = p.product_id GROUP BY p.department_id ORDER BY products_sold DESC ')  # Calculate department popularity proportion against all departments pie_data\$proportion &lt;- pie_data\$products_sold / sum(pie_data\$products_sold)  pie_data</pre>
---	---



#### Query 4: Top Performing product in each Department

For the fourth query, we wanted to get a snapshot of what the top products were in each department. Finding out exactly what item is being purchased the most posed many insights, as it was shocking that Sparkling Water Grapefruit was the most popular beverage, as we thought it would be a soda or water.

<pre># SQL query to find the top product per department (Histogram) top_item_per_department &lt;- "SELECT department, product_name FROM C SELECT   p.department,   p.product_name,   COUNT(o.product_id) AS total_orders,   ROW_NUMBER() OVER (PARTITION BY p.department ORDER BY COUNT(o.product_id) DESC) AS rank FROM   order_products_train o JOIN   products_cleaned p ON   o.product_id = p.product_id GROUP BY   p.department, p.product_name ) AS subquery WHERE rank = 1" sqldf(top_item_per_department)</pre>	<table><thead><tr><th>department</th><th>product_name</th></tr></thead><tbody><tr><td>1</td><td>alcohol</td></tr><tr><td>2</td><td>babies</td></tr><tr><td>3</td><td>bakery</td></tr><tr><td>4</td><td>beverages</td></tr><tr><td>5</td><td>breakfast</td></tr><tr><td>6</td><td>bulk</td></tr><tr><td>7</td><td>canned goods</td></tr><tr><td>8</td><td>dairy</td></tr><tr><td>9</td><td>deli</td></tr><tr><td>10</td><td>dry goods</td></tr><tr><td>11</td><td>frozen</td></tr><tr><td>12</td><td>household</td></tr><tr><td>13</td><td>international</td></tr><tr><td>14</td><td>meat seafood</td></tr><tr><td>15</td><td>missing</td></tr><tr><td>16</td><td>other</td></tr><tr><td>17</td><td>pantry</td></tr><tr><td>18</td><td>personal care</td></tr><tr><td>19</td><td>pets</td></tr><tr><td>20</td><td>produce</td></tr><tr><td>21</td><td>snacks</td></tr></tbody></table>	department	product_name	1	alcohol	2	babies	3	bakery	4	beverages	5	breakfast	6	bulk	7	canned goods	8	dairy	9	deli	10	dry goods	11	frozen	12	household	13	international	14	meat seafood	15	missing	16	other	17	pantry	18	personal care	19	pets	20	produce	21	snacks
department	product_name																																												
1	alcohol																																												
2	babies																																												
3	bakery																																												
4	beverages																																												
5	breakfast																																												
6	bulk																																												
7	canned goods																																												
8	dairy																																												
9	deli																																												
10	dry goods																																												
11	frozen																																												
12	household																																												
13	international																																												
14	meat seafood																																												
15	missing																																												
16	other																																												
17	pantry																																												
18	personal care																																												
19	pets																																												
20	produce																																												
21	snacks																																												

#### Query 5: Likelihood of reordering based on department

In this query, we look at how often departments can expect to get their products reordered versus not reordered. We can see larger disparities for departments where there is high consumption and turnover of its products such as with baking goods and alcohol as there is a much larger reorder rate associated. This is also attributed to the items being ordered having a higher product count in general.

```
# Aggregate by department and reorder status to calculate proportion
barchart_data <- sqldf('
  SELECT p.department,
         o.reordered,
         COUNT(*) AS product_count,
         COUNT(*) * 1.0 / SUM(COUNT(*)) OVER (PARTITION BY p.department) AS proportion
  FROM order_products_train o
  JOIN products_cleaned p ON o.product_id = p.product_id
  GROUP BY p.department, o.reordered
')
```

```
head(barchart_data)
```

```
> head(barchart_data)
```

	department	reordered	product_count	proportion
1	alcohol	0	2201	0.3931761
2	alcohol	1	3397	0.6068239
3	babies	0	6857	0.4589385
4	babies	1	8084	0.5410615
5	bakery	0	17702	0.3657891
6	bakery	1	30692	0.6342109

#### Query 6: Customers by the Hour

This query shows us the hours of the day people order food the most. The findings were not shocking, as it makes sense that the late morning and early afternoon have the most orders.

<pre>#sqldf query to see how many sales by the hour hourly_sales &lt;- "SELECT   order_hour_of_day,   COUNT(DISTINCT user_id) AS unique_users FROM   orders_cleaned GROUP BY   order_hour_of_day ORDER BY   order_hour_of_day" sqldf(hourly_sales)</pre>	<pre>&gt; sqldf(hourly_sales)</pre> <table><thead><tr><th>order_hour_of_day</th><th>unique_users</th></tr></thead><tbody><tr><td>1</td><td>15281</td></tr><tr><td>2</td><td>8854</td></tr><tr><td>3</td><td>5549</td></tr><tr><td>4</td><td>4225</td></tr><tr><td>5</td><td>4280</td></tr><tr><td>6</td><td>6866</td></tr><tr><td>7</td><td>18352</td></tr><tr><td>8</td><td>44579</td></tr><tr><td>9</td><td>74119</td></tr><tr><td>10</td><td>96666</td></tr><tr><td>11</td><td>108667</td></tr><tr><td>12</td><td>112986</td></tr><tr><td>13</td><td>113025</td></tr><tr><td>14</td><td>114484</td></tr><tr><td>15</td><td>115060</td></tr><tr><td>16</td><td>115105</td></tr><tr><td>17</td><td>111551</td></tr><tr><td>18</td><td>99676</td></tr><tr><td>19</td><td>85006</td></tr><tr><td>20</td><td>70025</td></tr><tr><td>21</td><td>55804</td></tr><tr><td>22</td><td>43537</td></tr><tr><td>23</td><td>35528</td></tr><tr><td>24</td><td>24778</td></tr></tbody></table>	order_hour_of_day	unique_users	1	15281	2	8854	3	5549	4	4225	5	4280	6	6866	7	18352	8	44579	9	74119	10	96666	11	108667	12	112986	13	113025	14	114484	15	115060	16	115105	17	111551	18	99676	19	85006	20	70025	21	55804	22	43537	23	35528	24	24778
order_hour_of_day	unique_users																																																		
1	15281																																																		
2	8854																																																		
3	5549																																																		
4	4225																																																		
5	4280																																																		
6	6866																																																		
7	18352																																																		
8	44579																																																		
9	74119																																																		
10	96666																																																		
11	108667																																																		
12	112986																																																		
13	113025																																																		
14	114484																																																		
15	115060																																																		
16	115105																																																		
17	111551																																																		
18	99676																																																		
19	85006																																																		
20	70025																																																		
21	55804																																																		
22	43537																																																		
23	35528																																																		
24	24778																																																		

#### Query 7: Peak Hours of Each Day of the Week

This query shows us the top two hours of every day of the week, along with their respective customer amount. What was slightly surprising is that Monday, having the second most customers per day, has different peak hours than every other day of the week.



<pre># Query to find the top 2 hours(for sake of output length) with the most unique customers for each day of the week top_hours_by_day &lt;- 'SELECT day_of_week, order_hour_of_day, unique_customers FROM (   SELECT     CASE order_dow       WHEN 0 THEN 'Sunday'       WHEN 1 THEN 'Monday'       WHEN 2 THEN 'Tuesday'       WHEN 3 THEN 'Wednesday'       WHEN 4 THEN 'Thursday'       WHEN 5 THEN 'Friday'       WHEN 6 THEN 'Saturday'     END AS day_of_week,     order_hour_of_day,     COUNT(DISTINCT user_id) AS unique_customers,     ROW_NUMBER() OVER (PARTITION BY order_dow ORDER BY COUNT(DISTINCT user_id) DESC) AS rank   FROM     orders_cleaned   GROUP BY     order_dow, order_hour_of_day ) AS ranked_hours WHERE rank &lt;= 2' sqlDF(top_hours_by_day)</pre>	<pre>&gt; sqlDF(top_hours_by_day)   day_of_week order_hour_of_day unique_customers 1      Sunday                14          37888 2      Sunday                15          37708 3      Monday                 11          34082 4      Monday                 10          33767 5      Tuesday                16          28037 6      Tuesday                15          28006 7      Wednesday             15          26913 8      Wednesday             16          26565 9      Thursday               15          25744 10     Thursday               16          25541 11     Friday                 14          27269 12     Friday                 15          27267 13     Saturday               14          28778 14     Saturday               15          28426</pre>
--	--

## Query 8: Items frequently bought together

This query displays common pairings of products or, in other words, which items are frequently bought together. Ranking at the top of this list is the pairings of bananas with the likes of strawberries, avocado, and spinach. This can serve as an indicator of healthy choices where a customer's purchasing pattern can be predicted via nutritional value and produce association.

<pre># Aggregate products that occur in together in an order co_occurrence &lt;- sqlDF('   SELECT p1.product_id AS product_id_1,          p2.product_id AS product_id_2,          COUNT(*) AS co_occurrence   FROM order_products_train p1   JOIN order_products_train p2 ON p1.order_id = p2.order_id AND p1.product_id &lt; p2.product_id   GROUP BY p1.product_id, p2.product_id ')  # Filter for top 20 co-occurrences top_co_occurrence &lt;- co_occurrence[order(-co_occurrence\$co_occurrence), ][1:20, ]  # Join to get the product names for plot visualization interpretability top_co_occurrence &lt;- sqlDF('   SELECT c.*, p1.product_name AS product_name_1, p2.product_name AS product_name_2   FROM top_co_occurrence c   JOIN products_cleaned p1 ON c.product_id_1 = p1.product_id   JOIN products_cleaned p2 ON c.product_id_2 = p2.product_id ')</pre>	<pre>&gt; top_co_occurrence   product_id_1 product_id_2 co_occurrence product_name_1 product_name_2 1      13176      21137          3074 Bag of organic Bananas Organic Strawberries 2      13176      47209          2420 Bag of organic Bananas Organic Hass Avocado 3      13176      21903          2236 Bag of organic Bananas Organic Baby Spinach 4      24852      47766          2216 Banana Organic Avocado 5      21137      24852          2174 Organic Strawberries Banana 6      24852      47626          2158 Banana Large Lemon 7      21903      24852          2000 Organic Baby Spinach Banana 8      16797      24852          1948 Organic Strawberries Banana 9      13176      27966          1780 Bag of organic Bananas Organic Raspberries 10     21137      27966          1670 Organic Strawberries Organic Raspberries 11     21137      21903          1639 Organic Strawberries Organic Baby Spinach 12     26209      47626          1595 Lines Large Lemon 13     21137      47209          1539 Organic Strawberries Organic Hass Avocado 14     21903      47766          1402 Organic Baby Spinach Organic Avocado 15     47626      47766          1349 Banana Organic Avocado 16     24852      26209          1331 Banana Lines 17     21137      39275          1269 Organic Strawberries Organic Blueberries 18     13176      30391          1268 Bag of organic Bananas Organic Cucumber 19     21903      47209          1252 Organic Baby Spinach Organic Hass Avocado 20     21903      47626          1238 Organic Baby Spinach Large Lemon</pre>
--	--

## Query 9: Product popularity by day of the week

In this query, we are looking at the number of products purchased on each separate day of the week. Looking at our findings, it is not surprising that the end of the week has fewer purchases, as most customer may be completing their grocery shopping at the beginning of the week.

<pre>bar_data &lt;- sqlDF('   SELECT o.order_day, COUNT(p.product_id) AS num_products   FROM orders_cleaned o   JOIN order_products_train p ON o.order_id = p.order_id   GROUP BY o.order_day ')</pre>	<pre>&gt; bar_data   order_day num_products 1      Friday      176910 2      Monday      205978 3      Saturday      207279 4      Sunday       324026 5      Thursday      155481 6      Tuesday      160562 7      Wednesday     154381</pre>
--	---

## Query 10: Order size versus days since the prior order

For this query, we look at the amount of items in a single order in comparison to the days since the user last ordered. Through this, we can analyze if there is a correlation between the size of the order and the frequency of their order history. The findings do not translate to much correlation, which may be because people make specific orders based on personal decisions and ways in which they cook, varying in amount and people whom they are buying for.

<pre> boxplot_data &lt;- sqldf(" SELECT o.user_id, o.order_number, o.days_since_prior_order, CASE   WHEN s.order_size BETWEEN 0 AND 5 THEN '0-5'   WHEN s.order_size BETWEEN 6 AND 10 THEN '6-10'   WHEN s.order_size BETWEEN 11 AND 15 THEN '11-15'   WHEN s.order_size BETWEEN 16 AND 20 THEN '16-20'   WHEN s.order_size BETWEEN 21 AND 30 THEN '21-30'   WHEN s.order_size BETWEEN 31 AND 40 THEN '31-40'   WHEN s.order_size BETWEEN 41 AND 50 THEN '41-50'   WHEN s.order_size BETWEEN 51 AND 60 THEN '51-60'   WHEN s.order_size BETWEEN 61 AND 70 THEN '61-70'   WHEN s.order_size BETWEEN 71 AND 80 THEN '71-80'   ELSE '80+' END AS order_size_in_buckets FROM orders_cleaned o JOIN   SELECT order_id, COUNT(product_id) AS order_size FROM order_products_train GROUP BY order_id ) o ON o.order_id = s.order_id" </pre>	<pre> &gt; head(boxplot_data) user_id order_number days_since_prior_order order_size_in_buckets 1      1           11                14                11-15 2      2           15                30                31-40 3      5            5                 6                 6-10 4      7           21                 6                 6-10 5      8            4                10                16-20 6      9            4                30                21-30 </pre>
--	--

## Query 11: Positions that products get added to the cart at

This query joins the necessary tables then calculates the proportion that departments commonly get added to the cart at a specific position. As one example, this can be interpreted that beverages are added to a customer's cart first roughly 14% of the time.

<pre> &gt; head(stacked_hist_data) add_to_cart_order department count proportion 1      1      alcohol  1268 0.009663971 2      1      babies   857 0.006531564 3      1      bakery  4625 0.035249106 4      1  beverages 18073 0.137742076 5      1  breakfast  2144 0.016340343 6      1      bulk    138 0.001051757 </pre>	<pre> combined_df &lt;- sqldf(" SELECT o.add_to_cart_order, p.department FROM order_products_train o LEFT JOIN products_cleaned p ON o.product_id = p.product_id" )  # Calculate the proportion departments are added to a customer's cart at specific spots stacked_hist_data &lt;- sqldf(" SELECT add_to_cart_order, department, COUNT(*) AS count, COUNT(*) * 1.0 / SUM(COUNT(*)) OVER (PARTITION BY add_to_cart_order) AS proportion FROM combined_df GROUP BY add_to_cart_order, department ORDER BY add_to_cart_order, department" ) </pre>
---	---

## Query 12: Total purchases in Each Department

In this SQL query, we look at department popularity based on the number of orders from each department. The findings were not shocking, as the stand-out departments are produce and dairy/eggs, which is what we predicted.

<pre> # create new table to compare sales in each department department_sales &lt;- products_cleaned %&gt;%   left_join(order_products_train, by = "product_id") %&gt;%   count(department) sales_by_department &lt;- "SELECT department, n FROM department_sales ORDER BY n DESC" sqldf(sales_by_department) </pre>	<pre> &gt; sqldf(sales_by_department) department n 1      produce 409238 2      dairy eggs 217467 3      snacks 120106 4      beverages 114885 5      frozen 101019 6      pantry 82506 7      bakery 48597 8      canned goods 47176 9      deli 44494 10     dry goods pasta 39045 11     household 36744 12     meat seafood 30437 13     breakfast 29683 14     personal care 23819 15     babies 15153 16     international 12158 17     missing 8604 18     alcohol 6005 19     pets 4714 20     other 1970 21     bulk 1362 </pre>
--	---

## Conclusion

By using data analysis techniques, we explored aspects of customer behavior related to the Instacart experience. We learned working with data of this scale can quickly accelerate to large computational costs and slow execution times, but insights derived from this analysis are highly actionable. Understanding customer patterns and product trends which are made available via querying and visualizations can directly drive sales with marketing strategies and inventory management. We saw potential continuation of this exploration with the incorporation of aspects like predictive modeling which then could provide personalized recommendations to customers, ultimately leading to improved business outcome and revenue.