



## User Recommendation Engine & Market Basket Analysis

March 22, 2020  
Xuechun Wang

Link: [https://github.com/XuechunWang/Cap\\_Recommendation](https://github.com/XuechunWang/Cap_Recommendation)

**01** Description of the Problem and Data Sets

**02** Insights and Exploratory Analysis

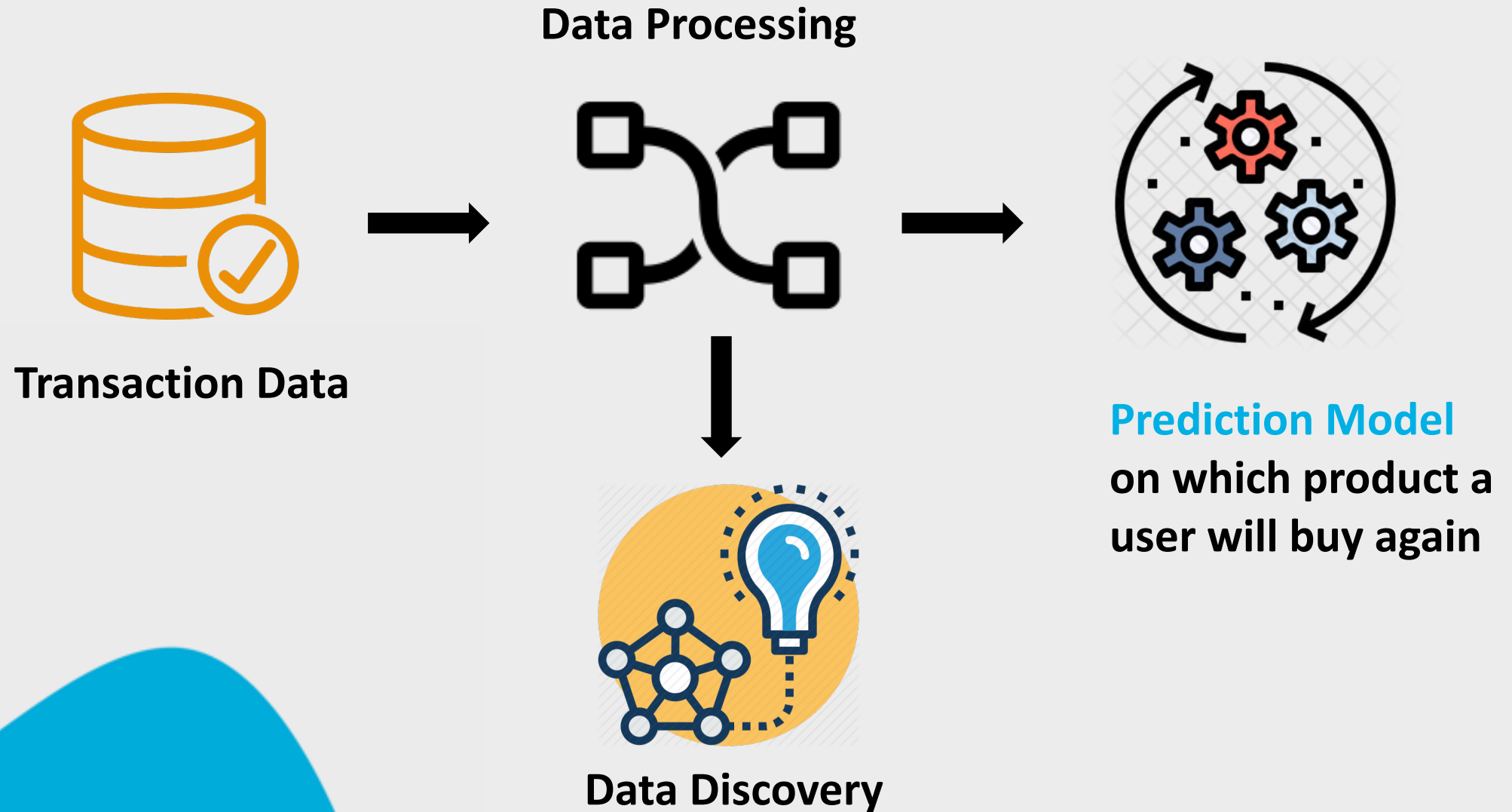
**03** Feature Engineering and Clustering Analysis

**04** Recommendation Engine Modeling and Diagnostics

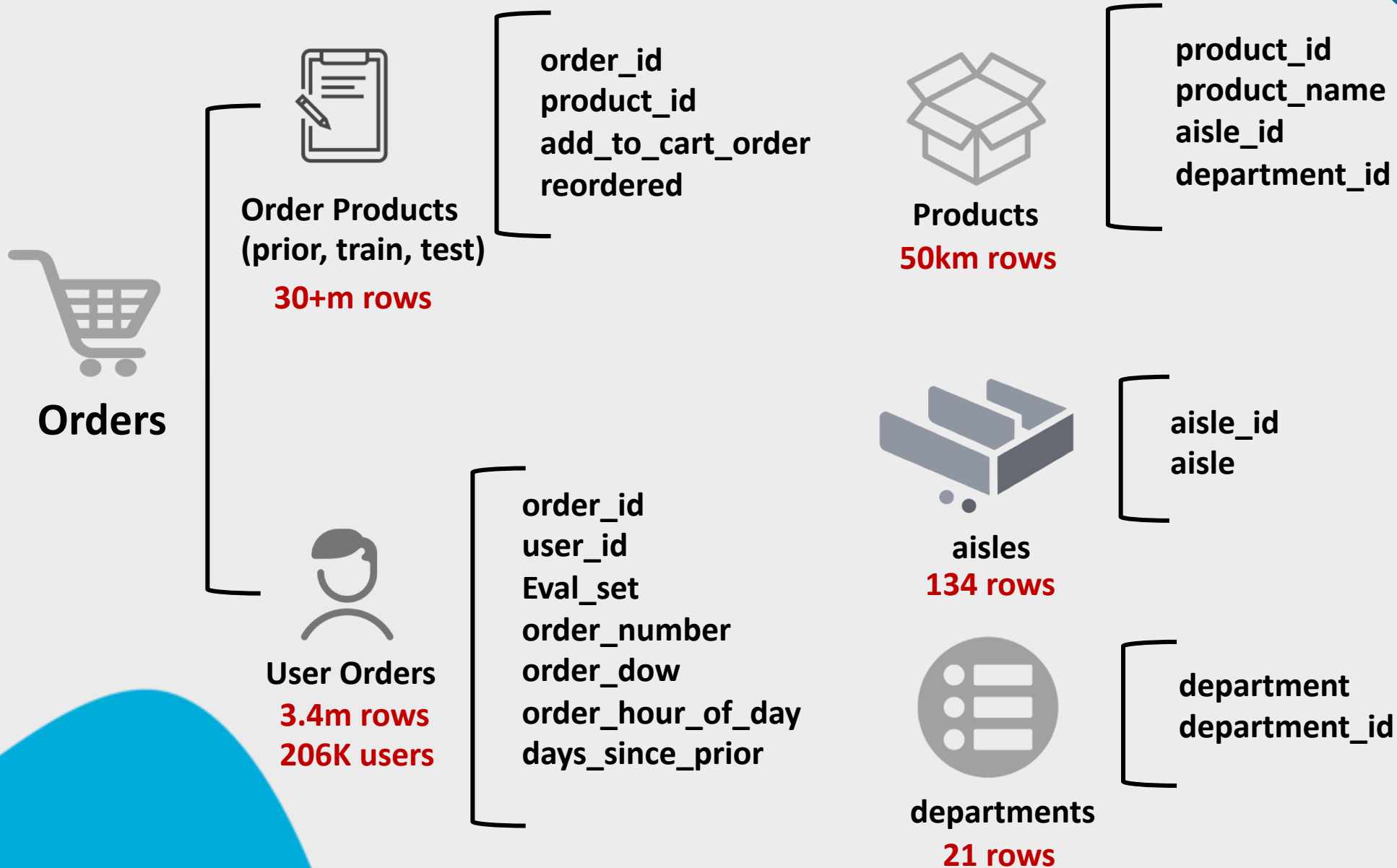
**05** Technology Used and Future Work



## Description of the **Problem** and Data Sets



# Description of the Problem and Data Sets



# Insights

and Exploratory  
Analysis



# Insights and Exploratory Analysis - Understanding the Products

## Top 3 Departments

How many unique products:

Personal care

Snacks

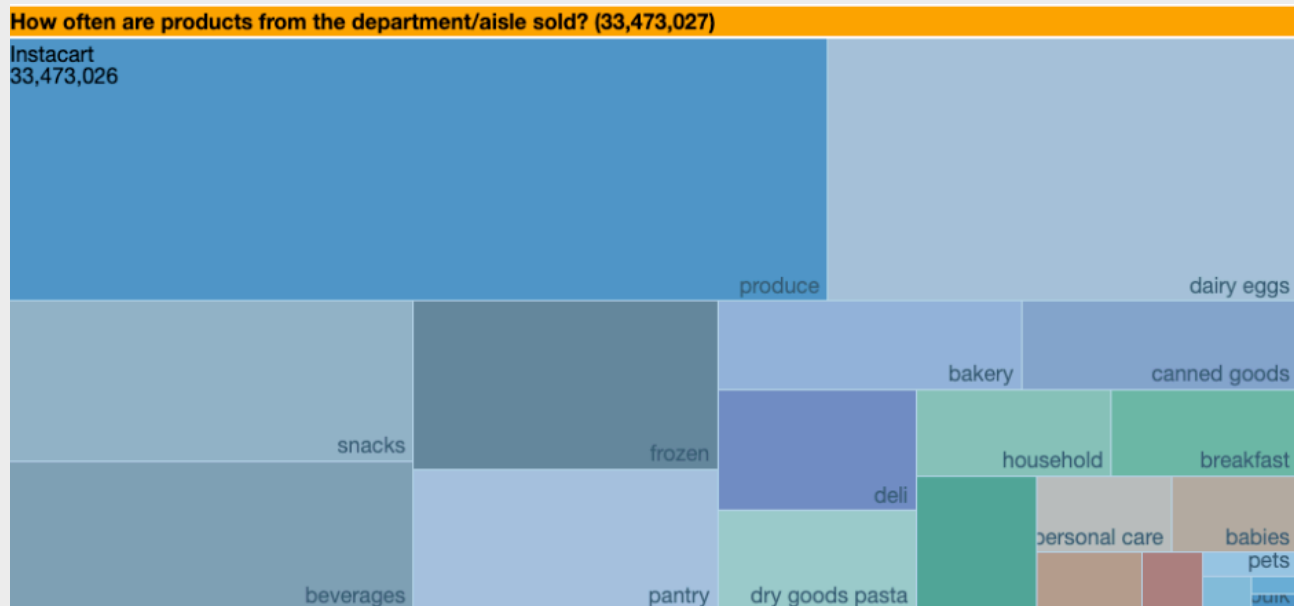
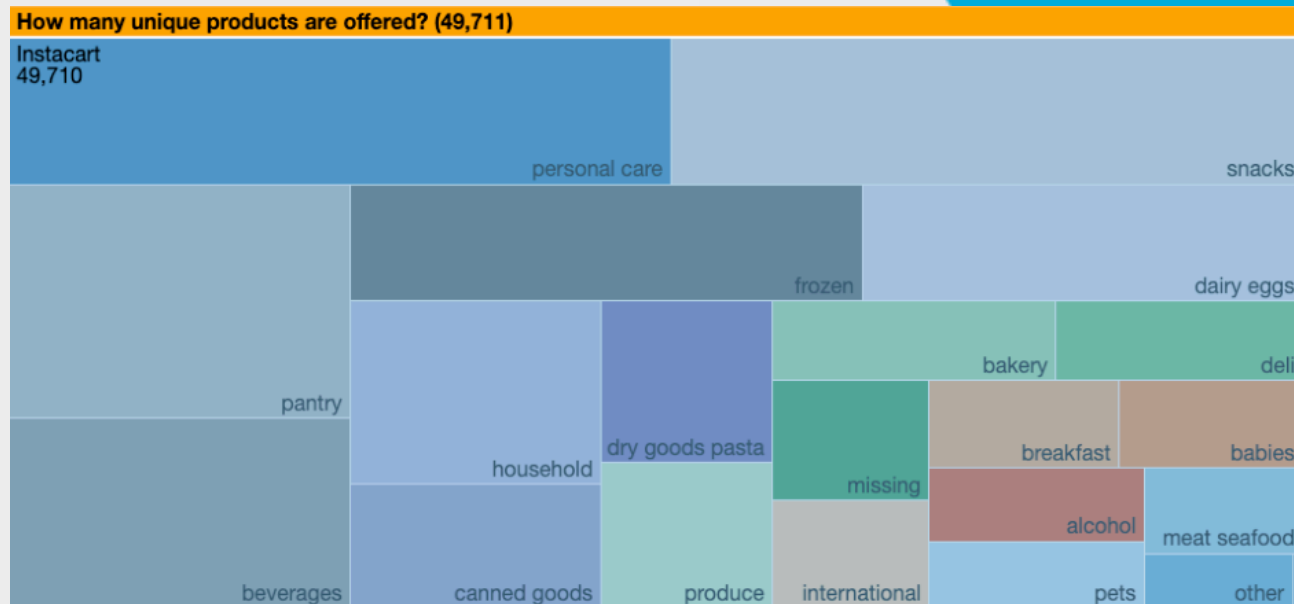
Pantry

How often purchased:

Produce

Dairy eggs

Snacks



# Insights and Exploratory Analysis – Never buy items

## Products that are never sold:

product_id		product_name	aisle_id	department_id
3629	3630	Protein Granola Apple Crisp	57	14
3717	3718	Wasabi Cheddar Spreadable Cheese	21	16
7044	7045	Unpeeled Apricot Halves in Heavy Syrup	88	13
37702	37703	Ultra Sun Blossom Liquid 90 loads Fabric Enhanc...	75	17
43724	43725	Sweetart Jelly Beans	100	21
45970	45971	12 Inch Taper Candle White	101	17
46624	46625	Single Barrel Kentucky Straight Bourbon Whiskey	31	



## Insights and Exploratory Analysis - Most frequently buy items

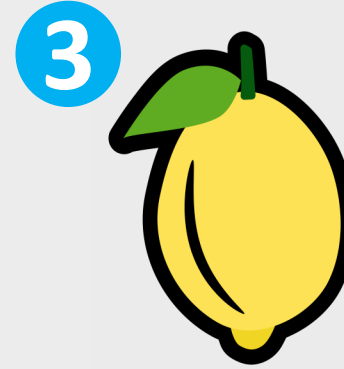
### Top List:



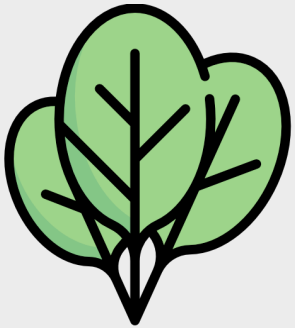
Banana



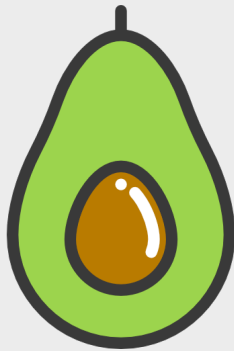
Strawberry



Lemon



spinach



avocado



milk



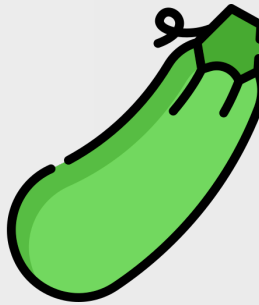
raspberry



onion



garlic



zucchini

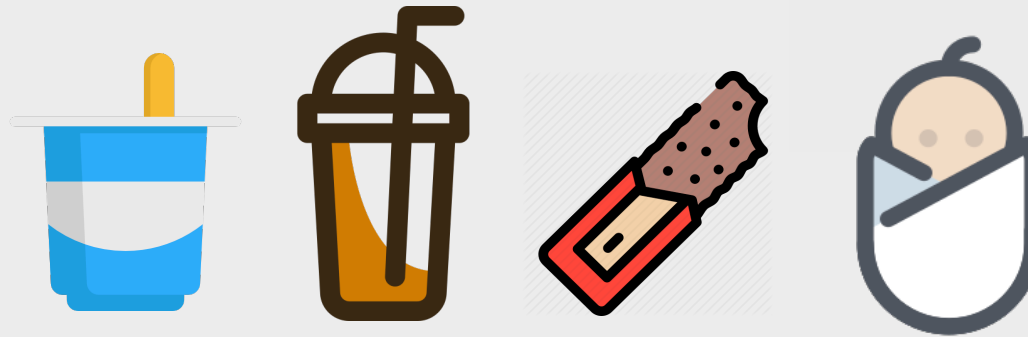




# Insights and Exploratory Analysis - Most frequently buy items

## From Association rule Analysis - **Apriori Algorithm**

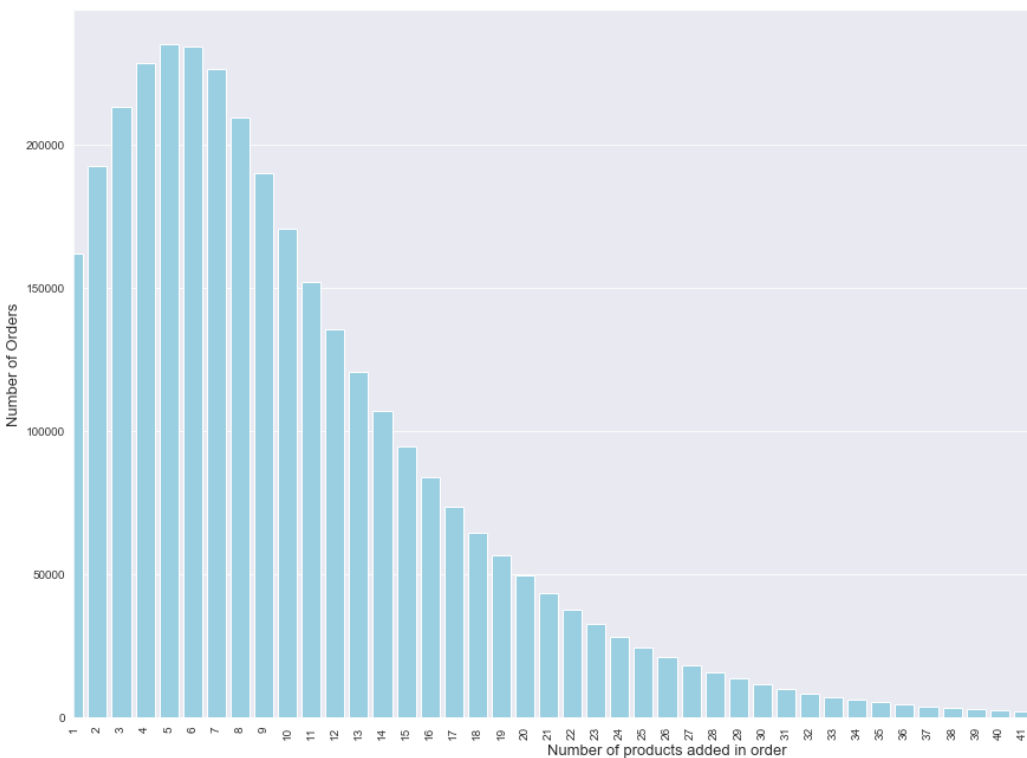
- One flavor of an item being purchased with another flavor from the same family being added to the order
  - Especially on **Yogurt, Beverage, and Food Bar, Baby Food**



- Recognized some common patterns:
  - Ex. **Moisture Conditioner - Hair Shampoos**

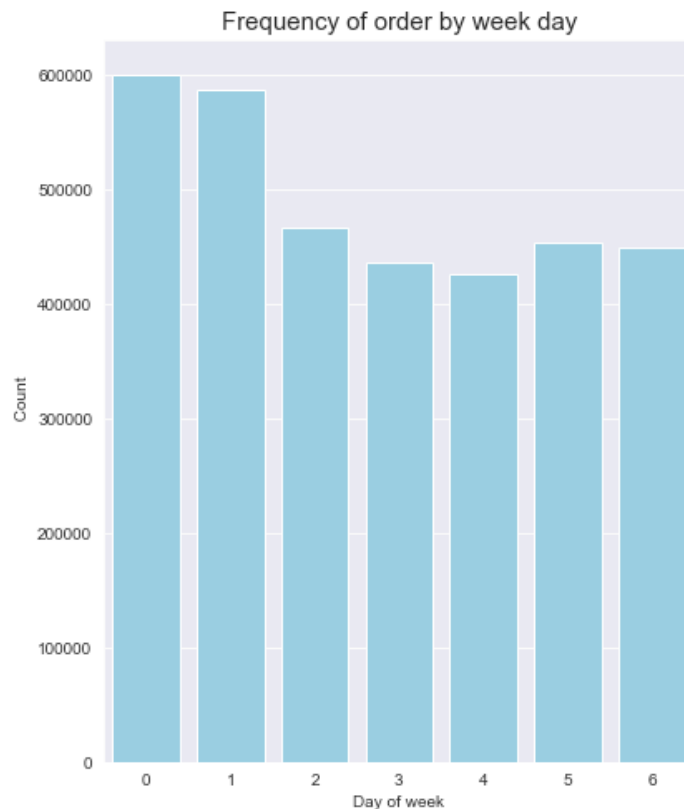


# Insights and Exploratory Analysis - Understanding the orders



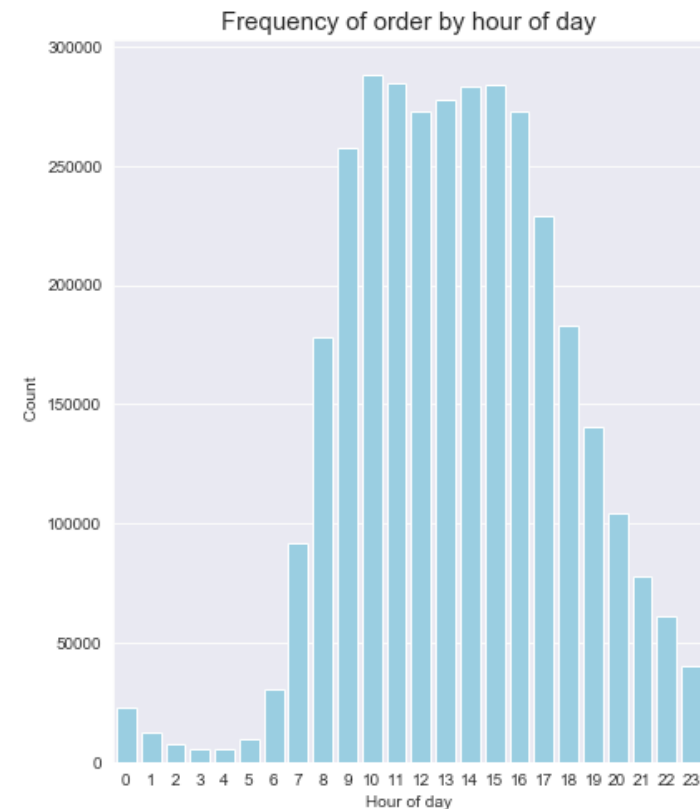
**Distribution of number of products per order**

Similar to Poisson distribution

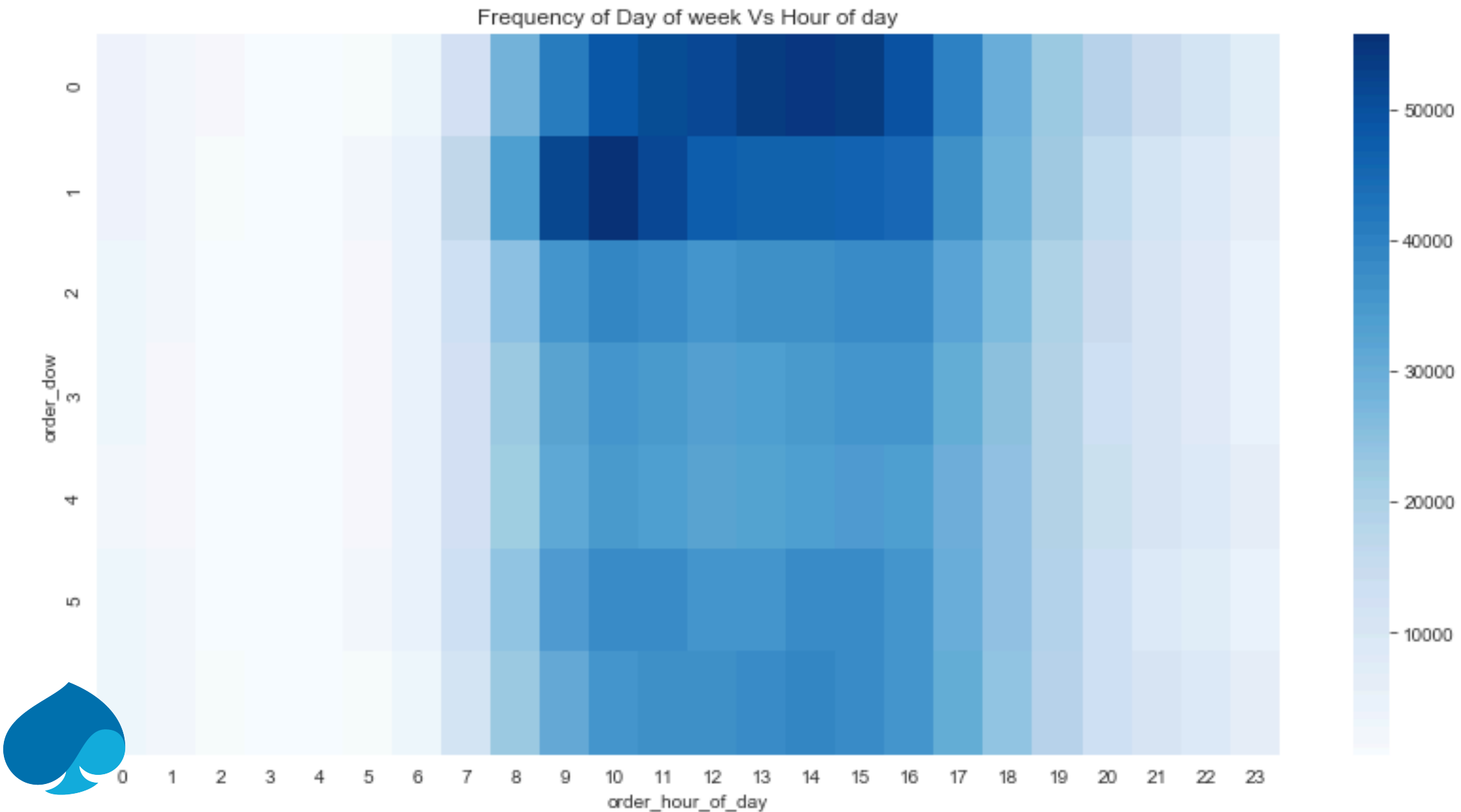


**Sunday and Monday** have relatively more orders during a week

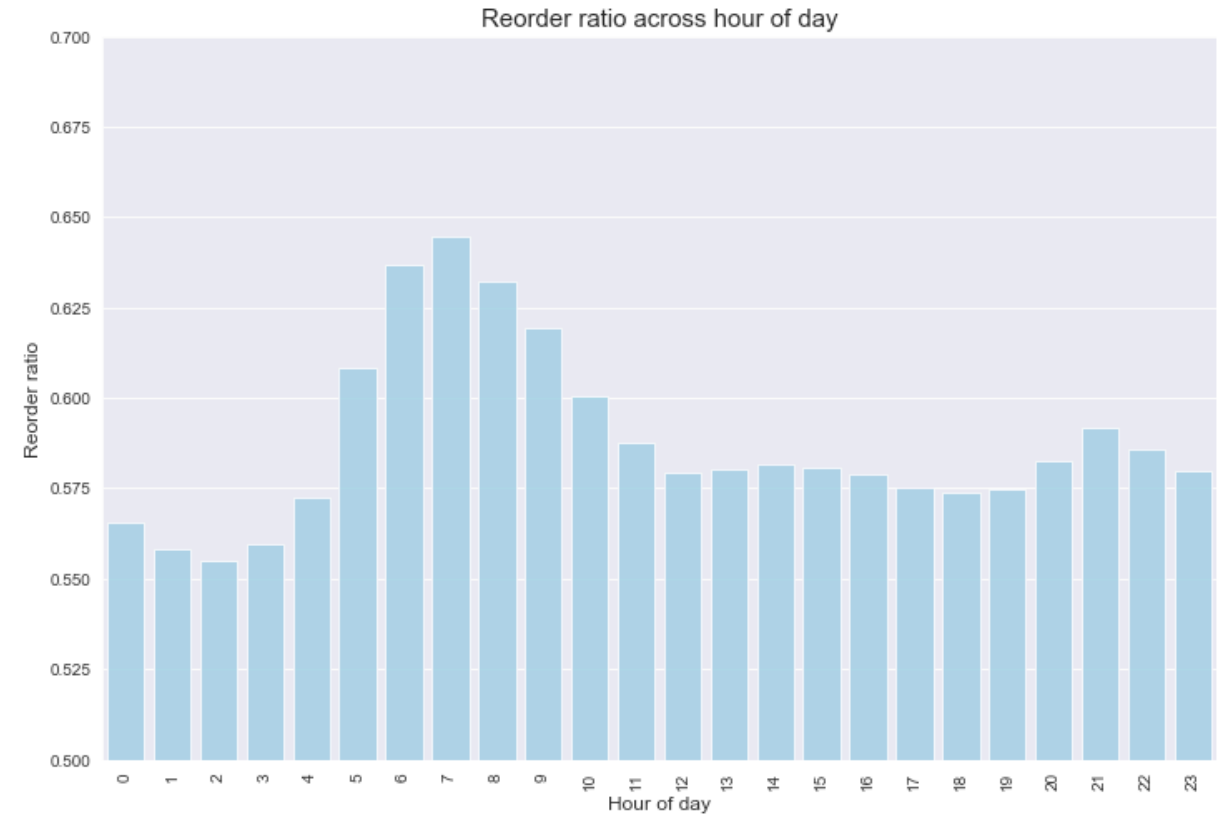
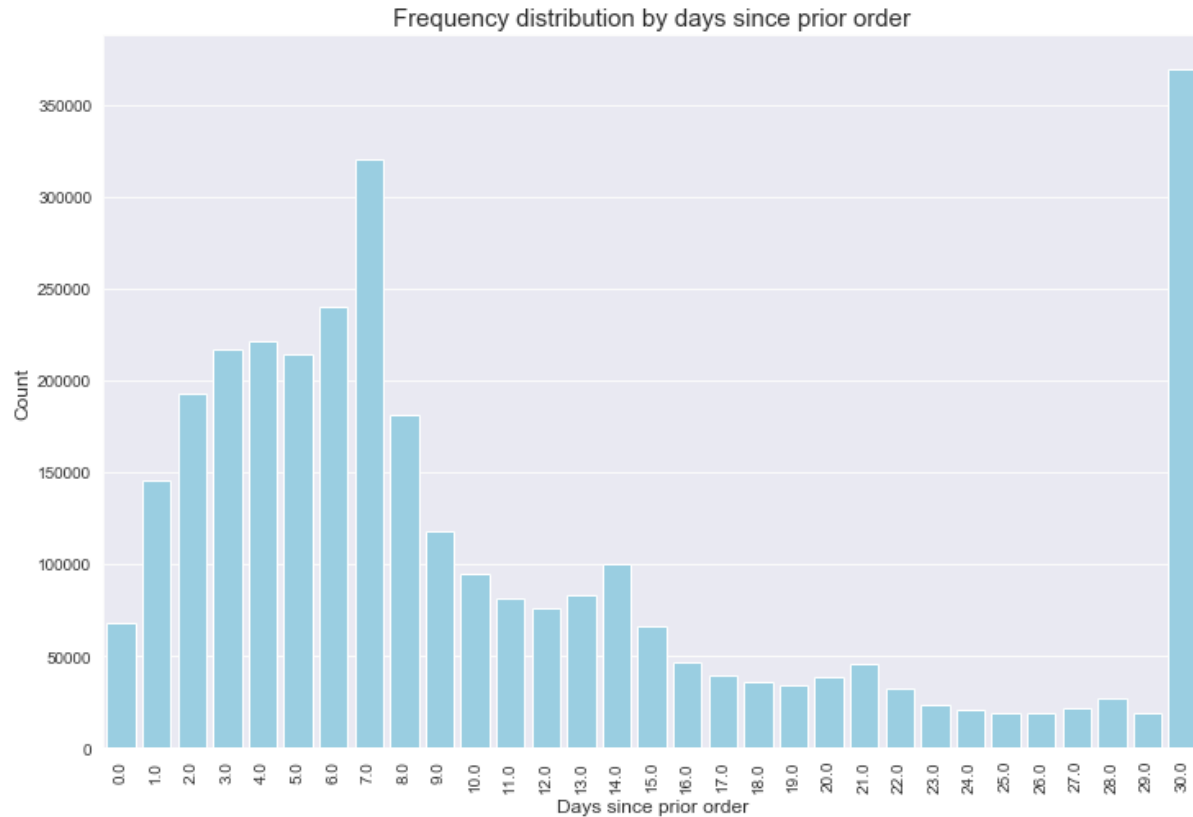
Most of orders occur during the time range 9-17 of a day



# Insights and Exploratory Analysis - Understanding the orders



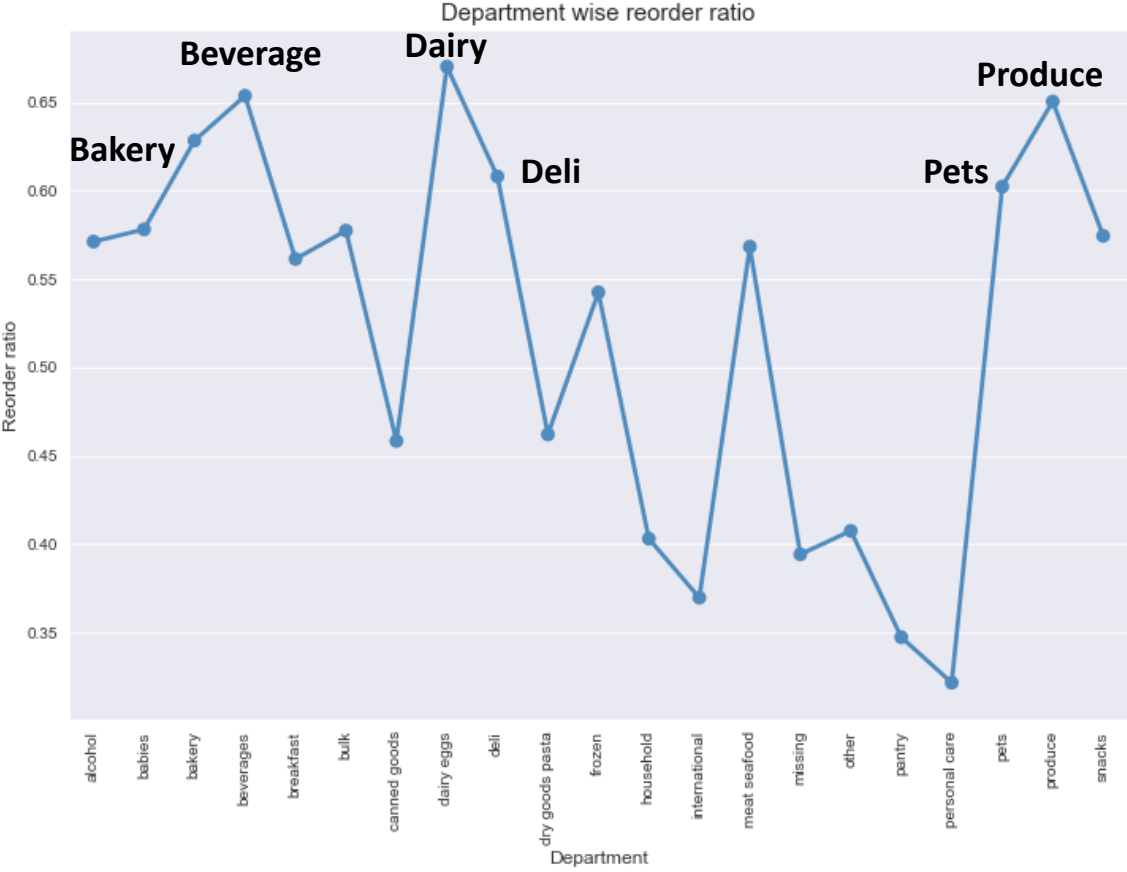
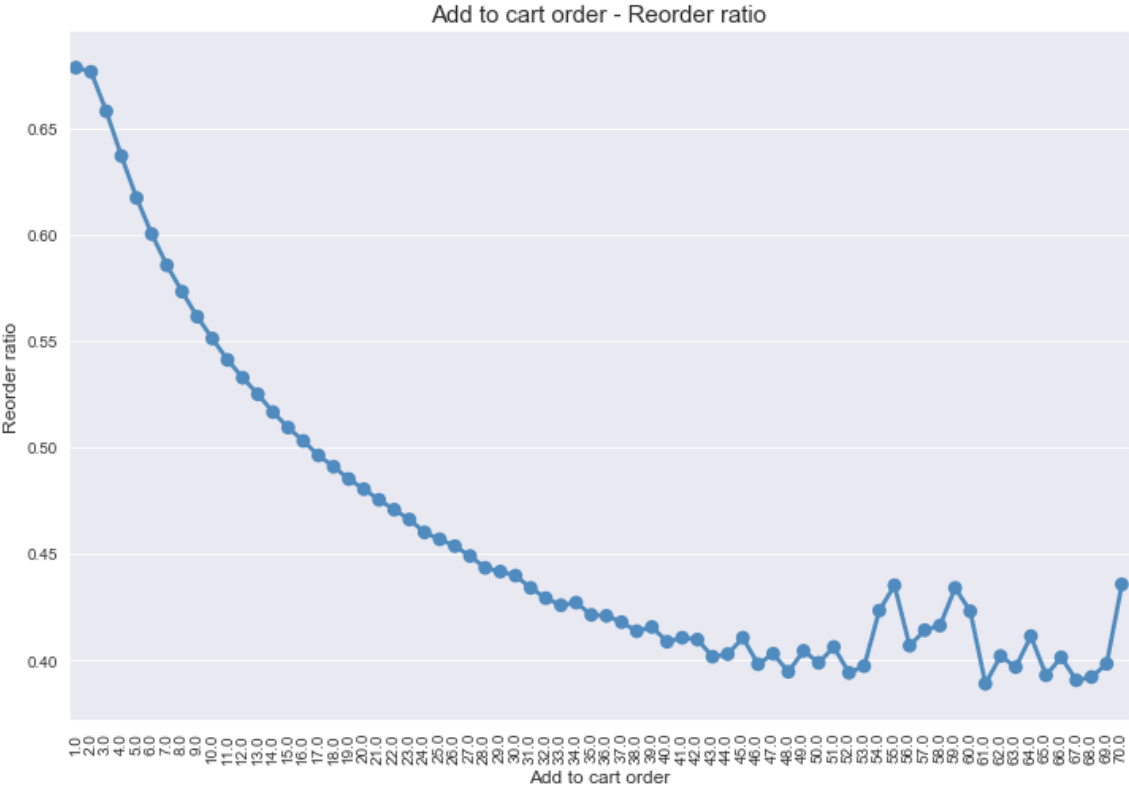
# Insights and Exploratory Analysis - Understanding the reorders



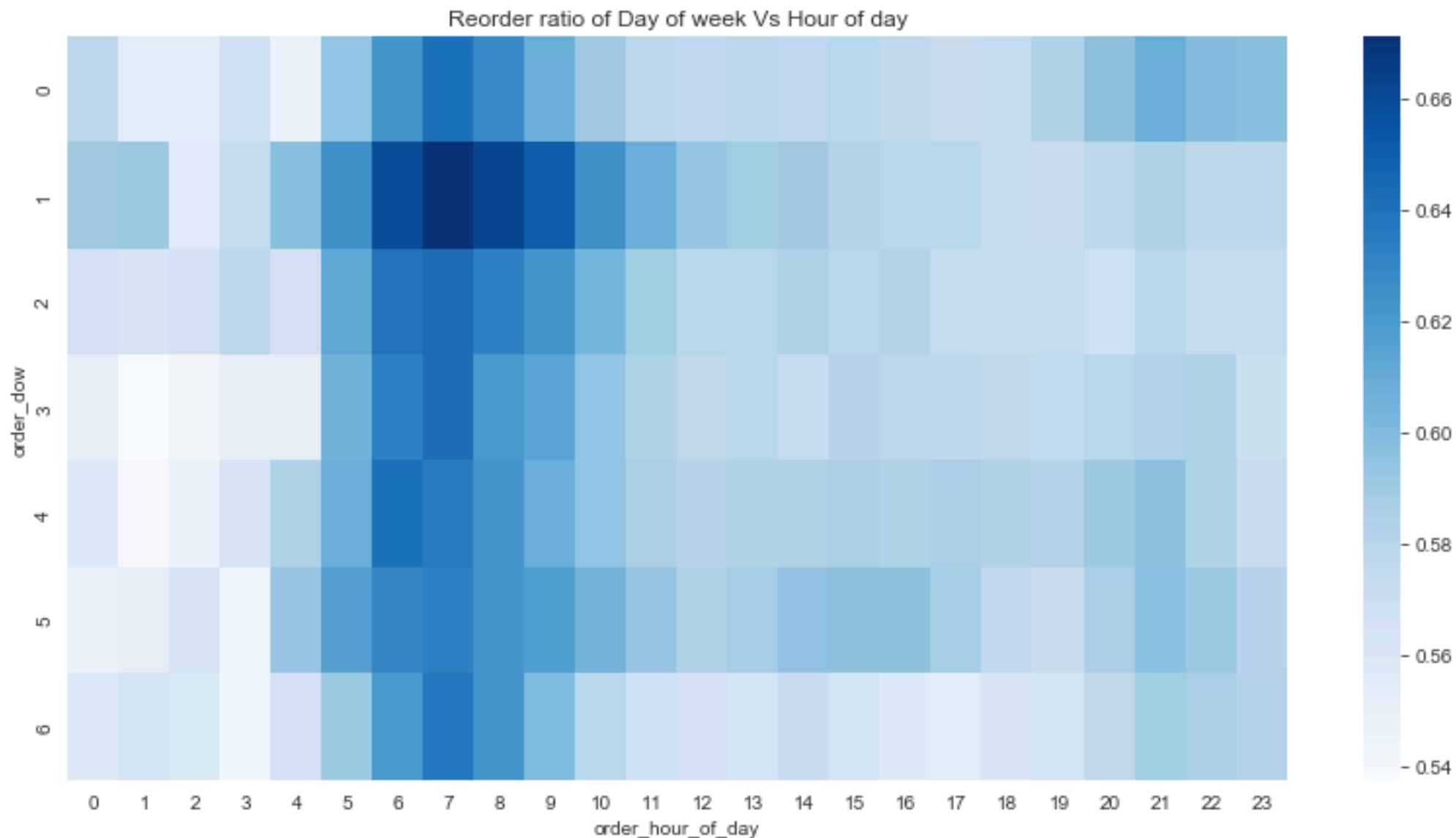
**The average times user buy a product is 3.31**  
**The maximum times product is 98**



# Insights and Exploratory Analysis - Understanding the reorders



# Insights and Exploratory Analysis - Understanding the reorders



# Feature Engineering & Clustering Analysis

Capgemini 

# Feature Engineering and Clustering Analysis

Use Merge and Group by agg function on order\_product data to create new features:



- product\_total
- **product\_reorder**
- product\_first\_order
- product\_second\_order
- product\_third\_order
- **product\_reorder\_pro**
- product\_triorder\_pro
- product\_reorder\_ratio
- product\_reorder\_times



- user\_total\_order
- user\_days\_since\_prior\_sum
- user\_days\_since\_prior\_avg
- user\_reorder\_ratio
- user\_total\_products
- user\_distinct\_products

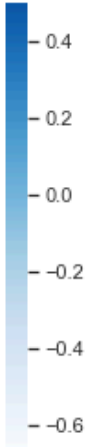
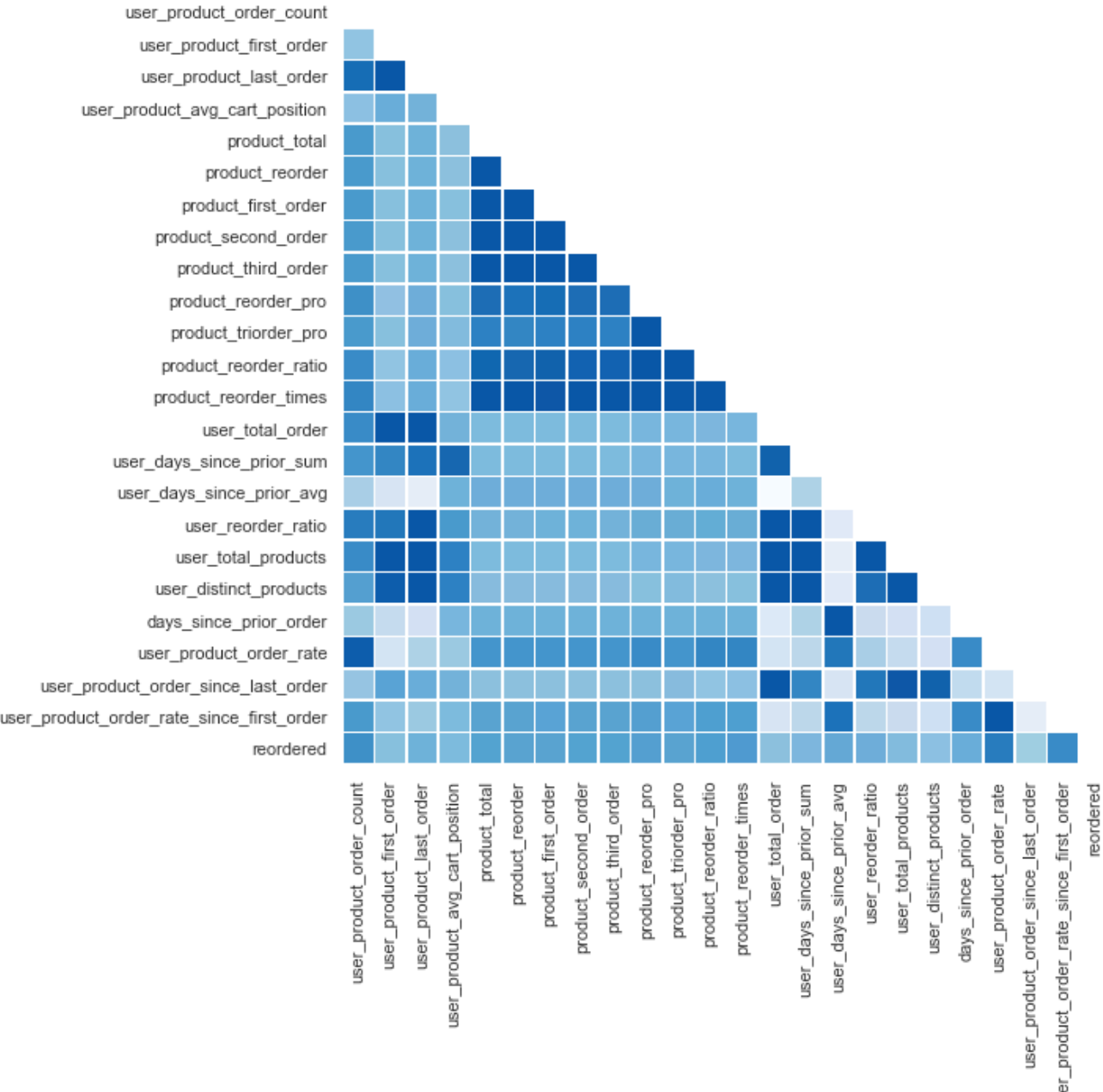


- **user\_product\_order\_count**
- user\_product\_first\_order
- user\_product\_last\_order
- user\_product\_avg\_cart\_position
- **user\_product\_order\_rate**
- user\_product\_order\_since\_last\_order
- **user\_product\_order\_rate\_since\_first\_order**

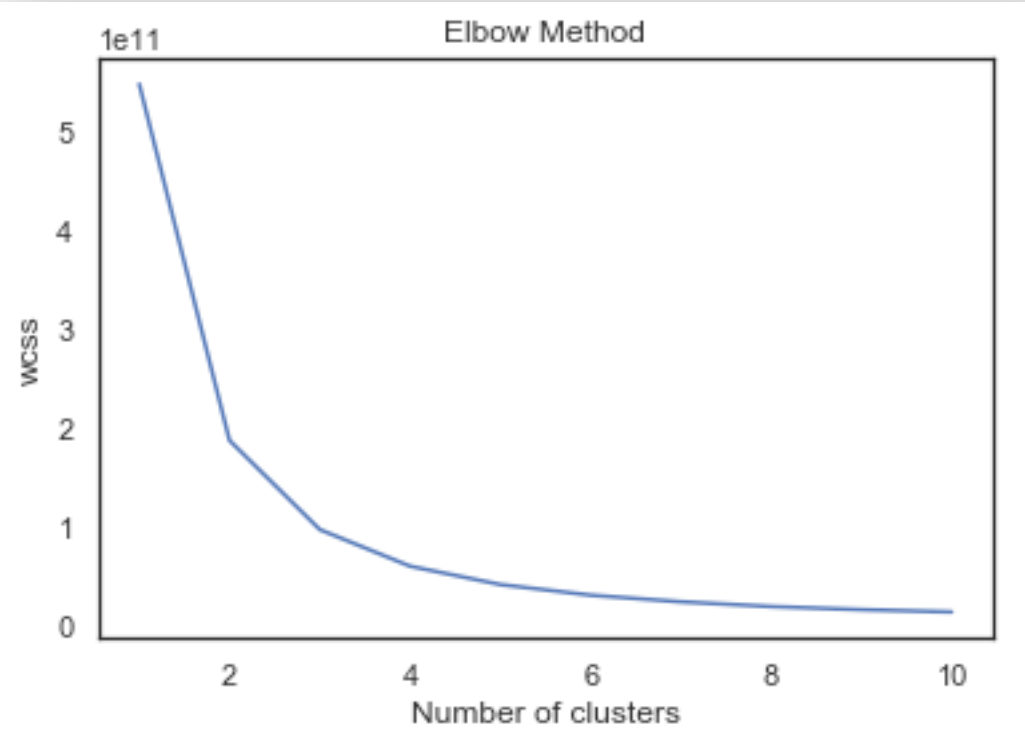




# Feature Engineering and Clustering Analysis



# Feature Engineering and Clustering Analysis



**Optimal Cluster Number: 4**



- **user\_total\_order**
- **user\_days\_since\_prior\_sum**
- **user\_days\_since\_prior\_avg**
- **user\_reorder\_ratio**
- **user\_total\_products**
- **user\_distinct\_products**

Cluster	Measure 1	Measure 2	Measure 3	Measure 4	Measure 5
1	8.31	16.42	0.42	52.30	31.24
2	31.05	12.98	0.64	393.18	135.84
3	19.82	15.31	0.53	176.51	79.30
4	37.28	11.44	0.73	759.99	199.17

# **Recommendation Engine Modeling & Diagnostics**



# Recommendation Engine Modeling and Diagnostics

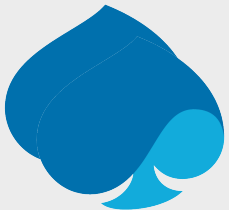
**01** Logistic Regression

**02** Lasso Regression

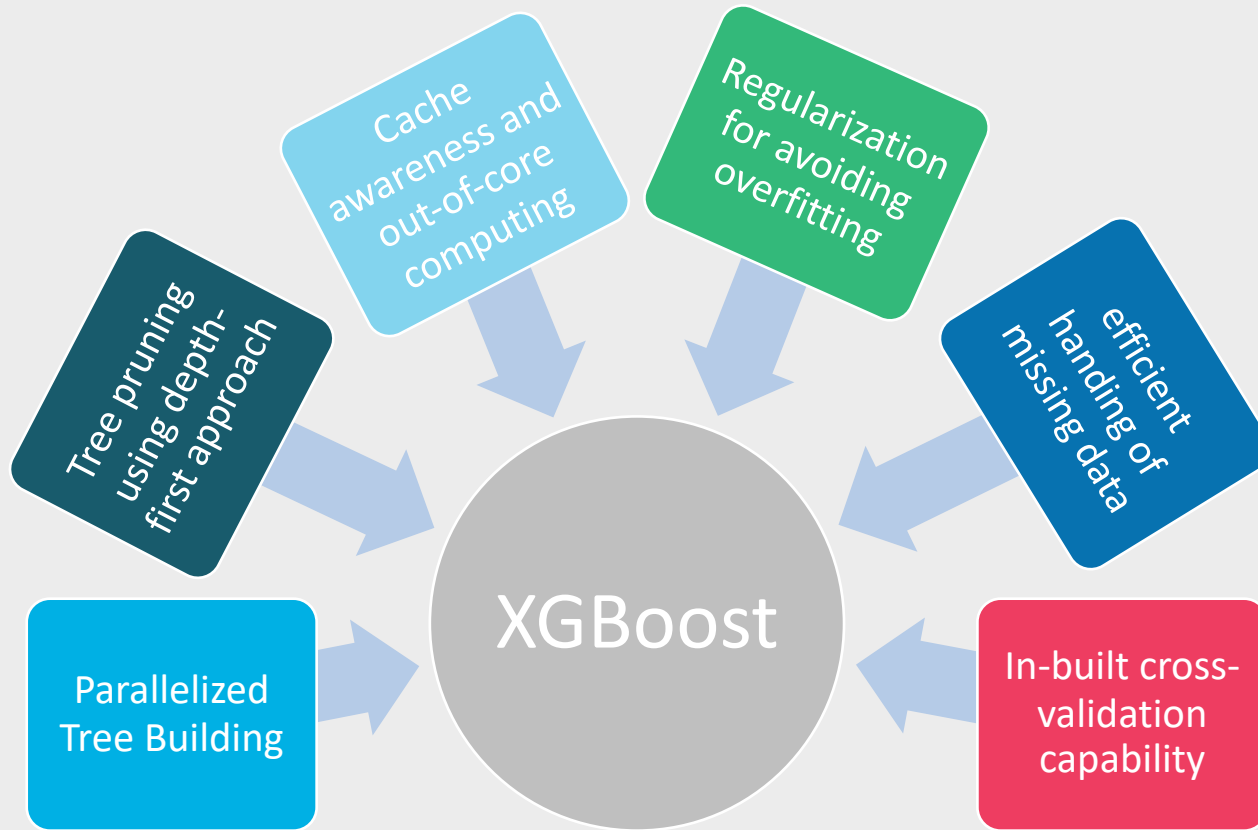
**03** XGBoost Model

**04** Optimize threshold

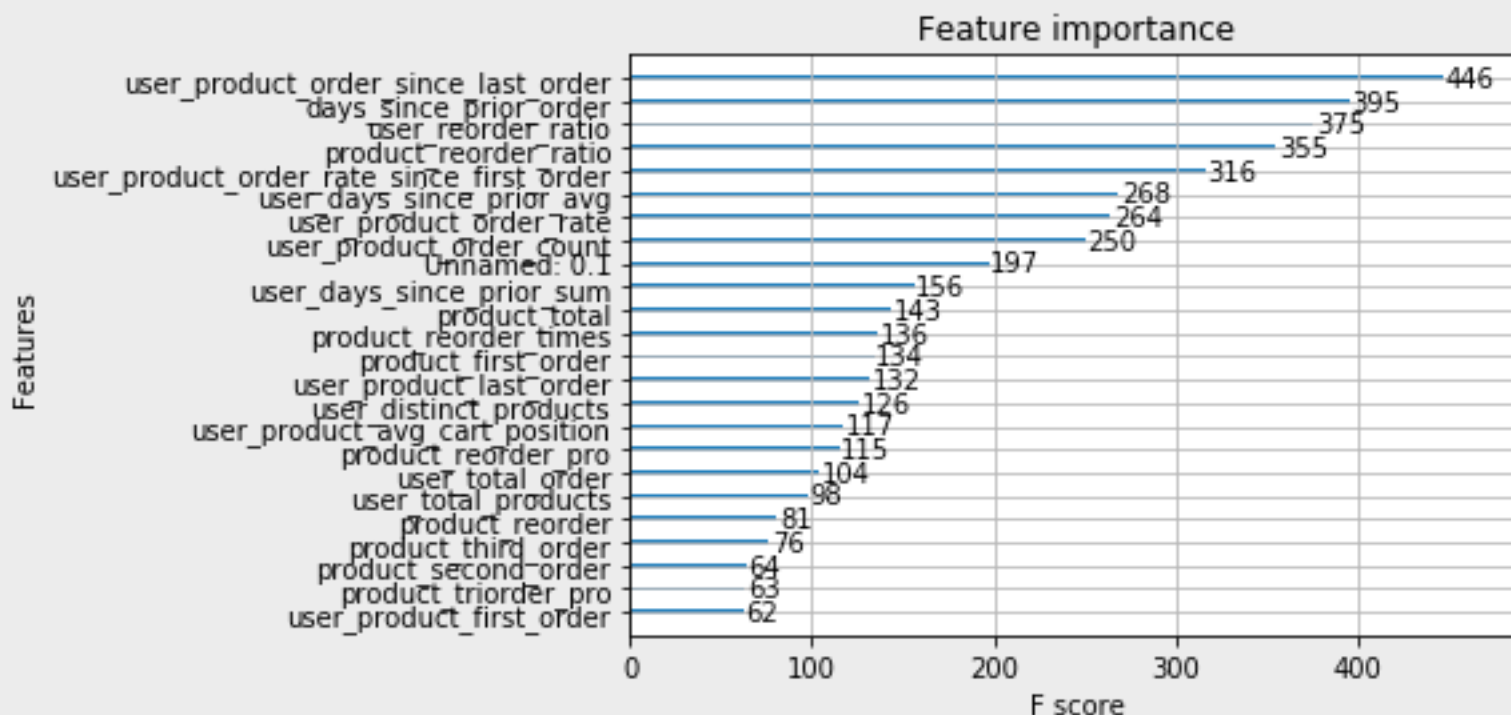
**05** Diagnostics



# Recommendation Engine Modeling and Diagnostics



# Recommendation Engine Modeling and Diagnostics



**Optimal Threshold for ROC Curve:** Best Threshold=0.073372

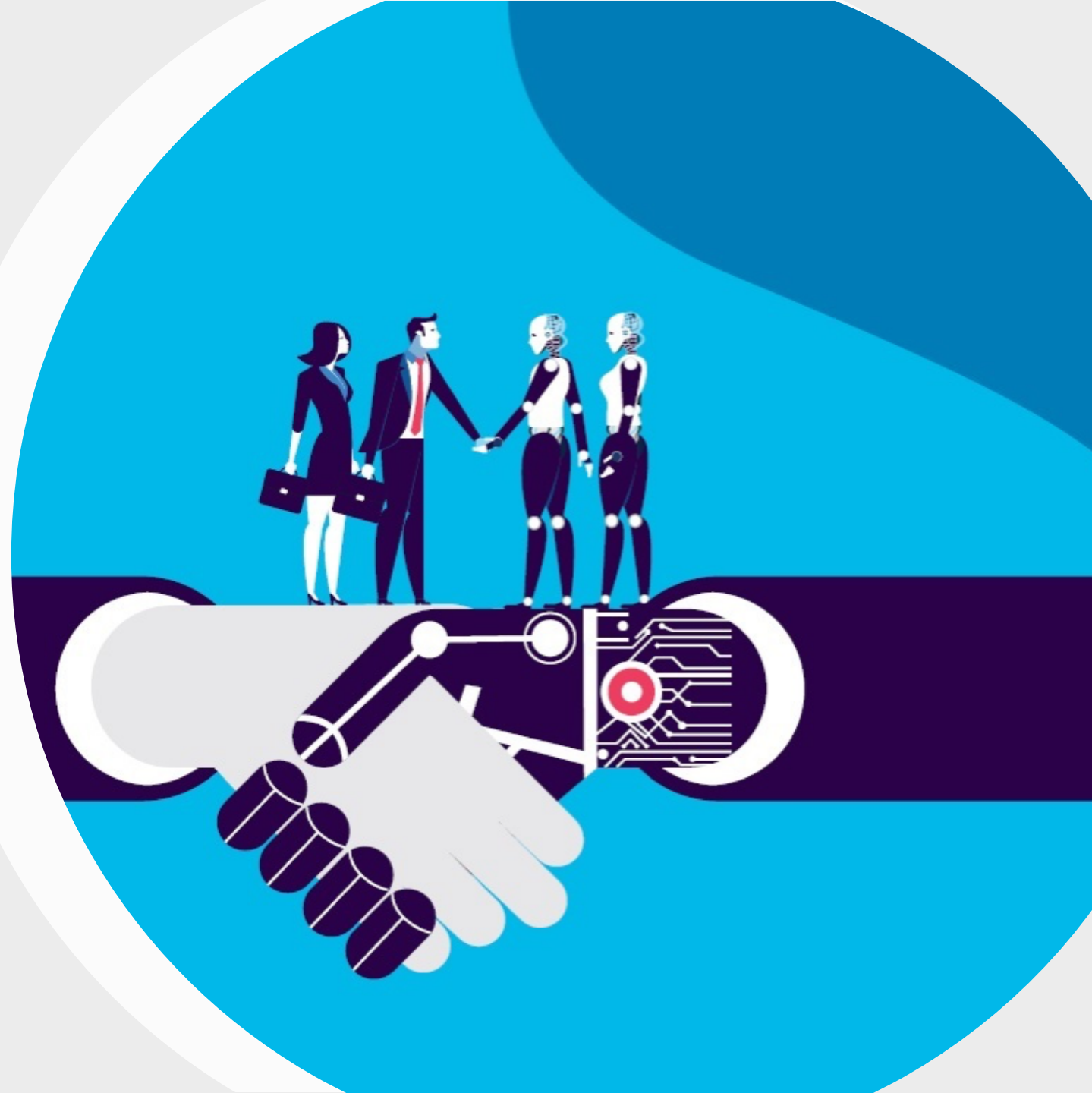
**Optimal Threshold for Precision-Recall Curve:** Best Threshold=0.184227

**Optimal Threshold Tuning on F-1 Score:** Best Threshold= 0.190





# Technology Used & Future Work



## Technology Used & Future Work

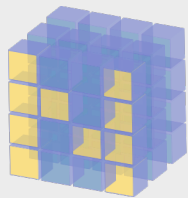


Seaborn

Pandas



matplotlib



NumPy



## Technology Used & Future Work

**01** Description of the Problem and Data Sets

**02** Insights and Exploratory Analysis

**03** Feature Engineering and Clustering Analysis

**04** Recommendation Engine Modeling and Diagnostics

**05** Technology Used and Future Work



Capgemini 





Capgemini 