The background is white with various colorful geometric shapes and illustrations. There are green circles, orange rectangles, and pink triangles. Scattered throughout are small black dots, zigzag lines, and simple line drawings of a pineapple, a carrot, and an apple. The title is centered in a large, bold, black font.

# InstaCart & Machine Learning

A Not-So “Insta” Analysis

Madison Dimaculangan  
June 2020



01

## Introduction

What is Instacart?  
Concepts and goals.



02

## Model Preparation

Dataset transformation.  
Data exploration.



03

## Unsupervised Learning


Using clustering to find patterns in the data.



04

## Supervised Learning

Using algorithms to make predictions.



The background is a vibrant yellow with a large diagonal split. The left side is white, and the right side is yellow. Various geometric shapes like triangles, circles, and lines are scattered across the background. On the left side, there are illustrations of a mushroom, a banana, and a broccoli. The text '01' is prominently displayed in the upper right, and 'Introduction' is in the center right. Below 'Introduction' is the subtitle 'Background and goals.'.

# 01

## Introduction

Background and goals.



# What is Instacart?

- alternative to traditional grocery experience
- on-demand grocery delivery service
- presence in 5,500 cities in US & Canada
- employs personal shoppers to fulfill & deliver
- partnerships with 350 retailers (over 25,000 locations)



# Project Goals



The main goal of this project is to determine whether machine learning techniques can be applied to the Instacart dataset to suggest new features to improve user experience.

This can be accomplished by learning more about the users themselves:

- Can the user population be divided into groups of users with similar characteristics?
- Can we make any predictions on future purchases based on the ordering history?

These and other questions, once answered, can lead to actions that allow for increased customer retention and product usage.





# Data Set




The data used for this project is from the 2017 competition that Instacart hosted on Kaggle.com.

The Instacart Market Basket Analysis competition data set consists of several csv files:

1. aisles.csv
2. departments.csv
3. orders.csv
4. products.csv
5. order\_products\_\*\_.csv (where \* = [prior, train])

The Model Preparation section of this slide deck will detail how the various files will be transformed and merged to create our model.

**Note:** While all work in this project was completed on data from 2017, the process can be repeated on any similarly collected data from any year.





# 02

## Model Preparation

Data merging and visualization.

# Dataset Inspection

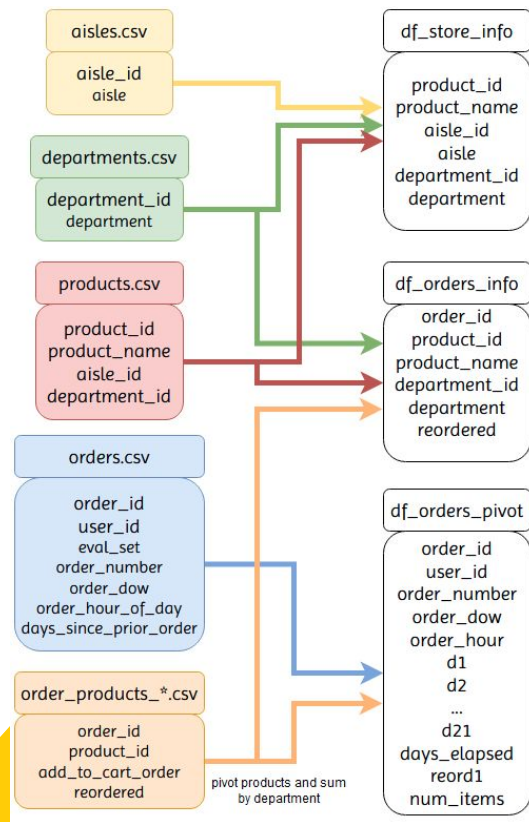
	Details
<b>aisles</b>	134 unique aisles, including 1 for “missing”
<b>departments</b>	21 unique departments
<b>products</b>	49,687 unique products
<b>orders</b>	3,421,083 unique orders
<b>users</b>	206,209 unique users



# Dataset Merging

Using PySpark and SQL commands:

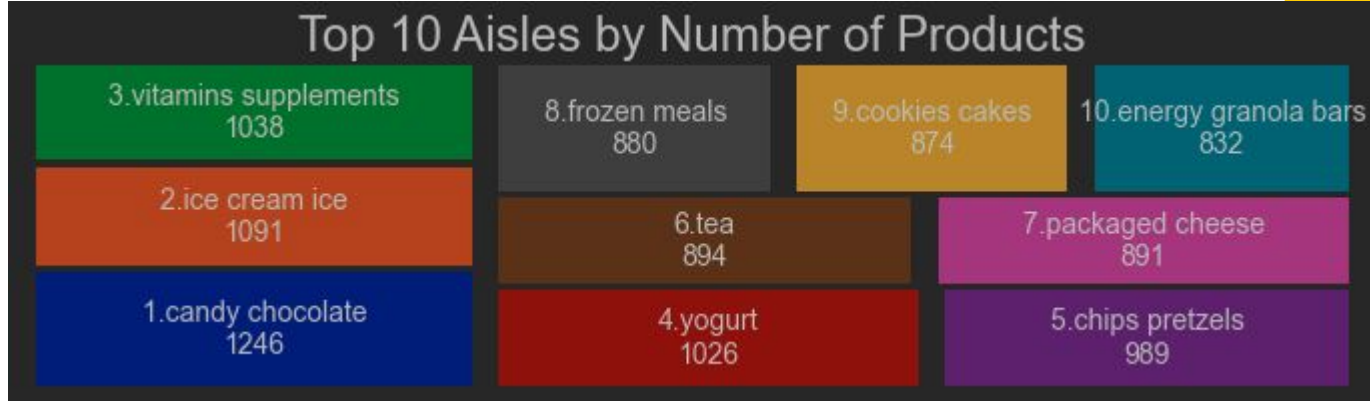
- from 6 csv files, created 3 dataframes to better associate order, user, and product information
- grouped products by their departments and recorded the sums as separate columns
- grouping and aggregating reduced observation count ten-fold, from ~32 million to ~3.4 million
- added a new column summing the number of items ordered per order\_id



# Aisles

In total, there are 133 unique aisles (plus 1 additional “aisle” for uncategorized products).

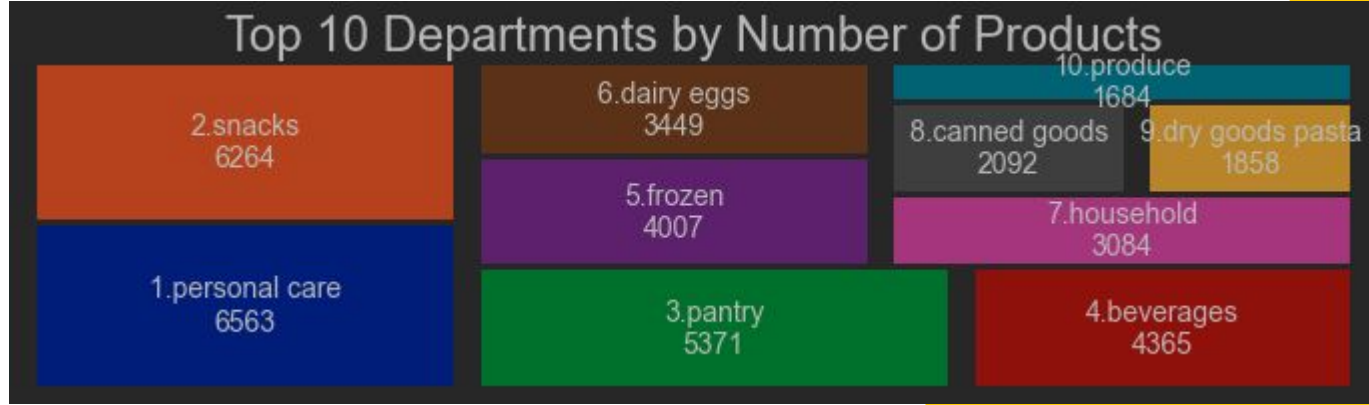
Below is a a graphic showing the 10 aisles with the greatest number of unique products.



# Departments

In total, there are 21 unique departments.

Below is a a graphic showing the 10 departments with the most number of products.



# Orders- by products

In total, there are 3,421,083 orders.

Below is a a graphic showing the 10 products ordered by users.

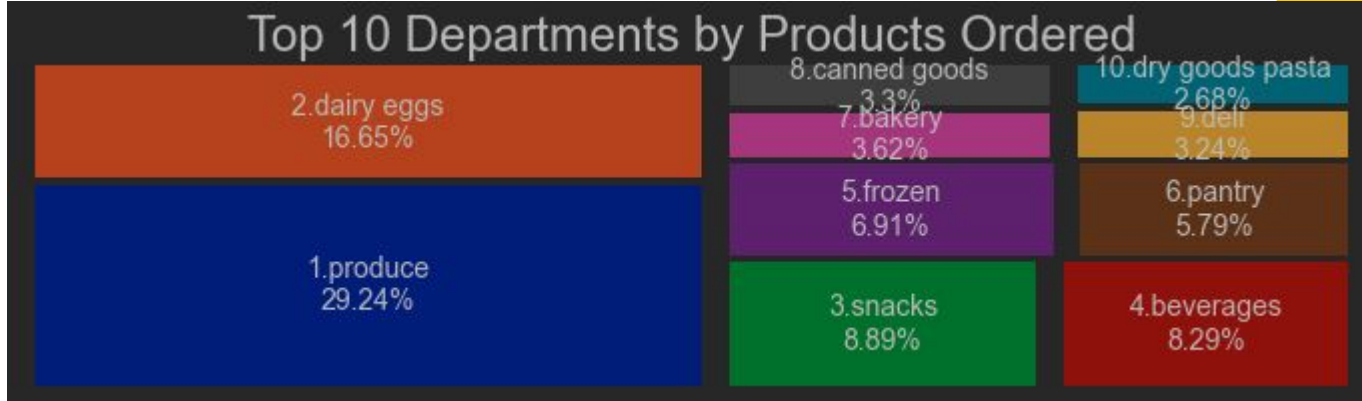
Instacart users love bananas!



# Orders – by departments

Below is a graphic showing the 10 departments with the most products ordered.

Perishable items are the most ordered with produce and dairy items and eggs accounting for over 45% of all products ordered.

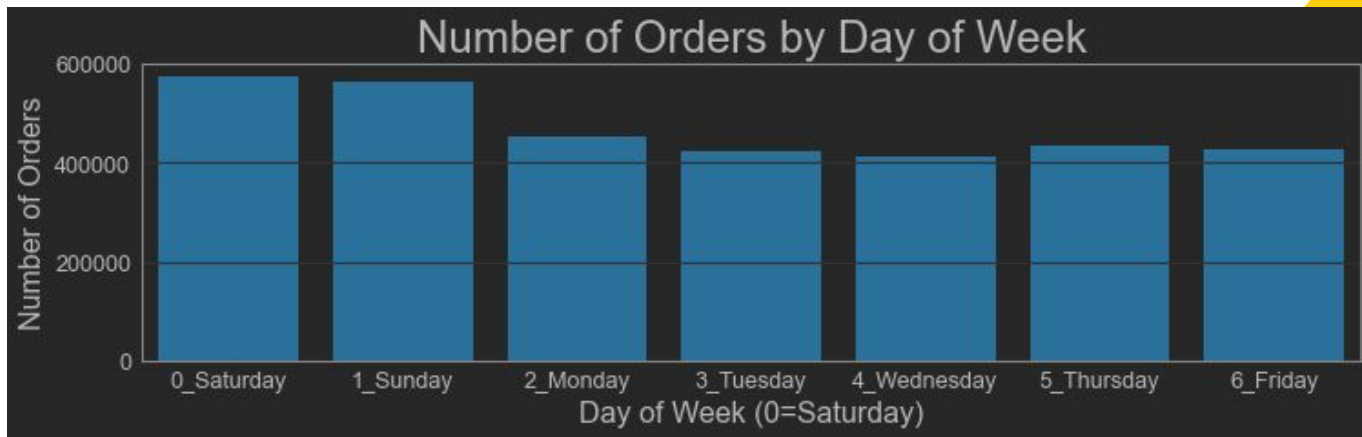




# Orders – By day of week



- Orders are most often placed during the weekends
- The number of orders decrease towards the middle of the workweek.

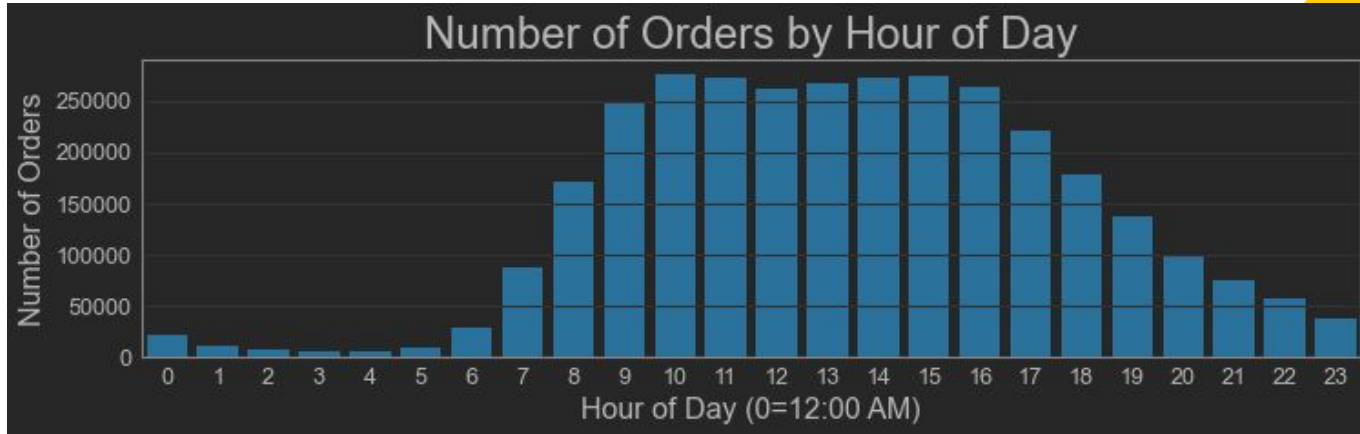




# Orders – By hour of day



- The peak hours are between 10 AM and 4 PM.
- The number of orders decrease gradually until midnight.
- There is a dead period of low number of orders between midnight and 6 AM.

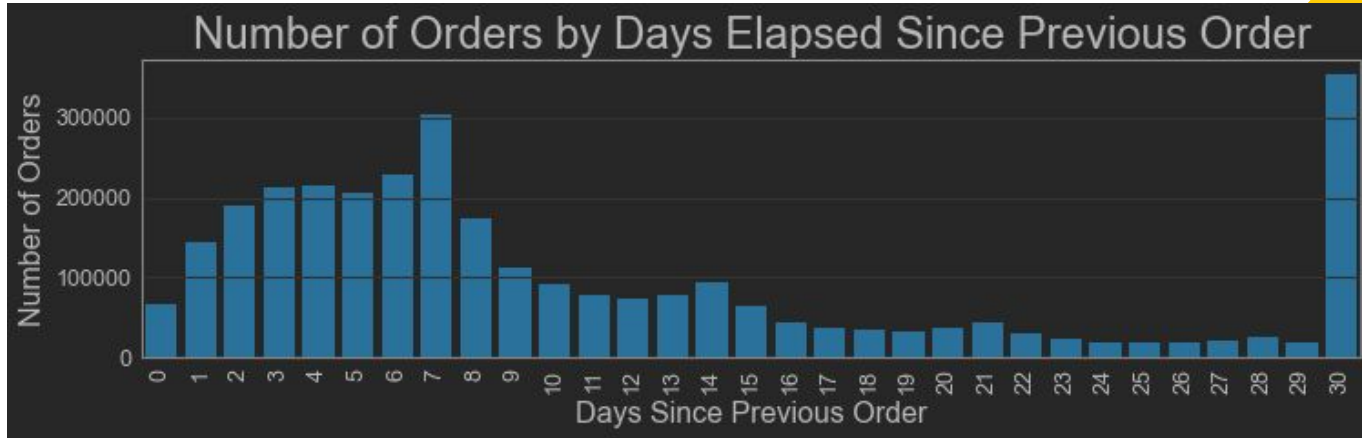




# Orders – By days since previous order



- The lag between orders ranges from 0 to 30 days.
- Excluding first orders, there are several distributions centered around:
  - every 3-4 days
  - every week
  - every other week
  - longer than 3 weeks





# Orders – By Order Number

- order\_number indicates which (1st, 2nd, 3rd, etc.) order it is for the customer
- 25% of all orders placed occur within the first 4 orders
- after 4 orders there is a dropoff, providing an opportunity for retention improvement

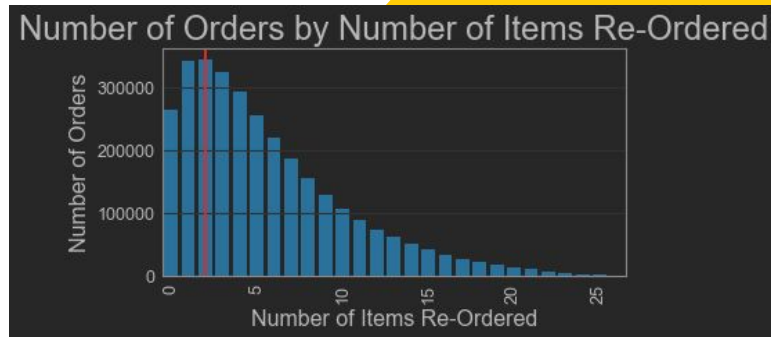




# ● Orders – By number of items ordered/re-ordered



- The number of items ordered is distributed around 5 items with tail of up to 26 items
- The number of items re-ordered is distributed around 2 items with a tail up to 26 items



# Correlations

- d4 (produce) and d16 (dairy/eggs) correlate the best with total number of items ordered
- these two departments had the most number of items ordered
- no variable correlates well with the order lag time, days\_elapsed

	index	corr
22	num_items	1.000000
21	reord1	0.731229
3	d4	0.653757
15	d16	0.593018
18	d19	0.401165
0	d1	0.390361
12	d13	0.381021
2	d3	0.353039
14	d15	0.346875
19	d20	0.335703

	index	corr
23	days_elapsed	1.000000
0	d1	0.033226
16	d17	0.032168
8	d9	0.026241
14	d15	0.022098
11	d12	0.016751
22	num_items	0.016646
10	d11	0.015481
19	d20	0.015340
13	d14	0.011909

The background is a vibrant yellow with a large pink triangle on the left side. Scattered throughout are various geometric shapes: small black dots, white triangles, and a zigzag line. Food items are also present: a small orange mushroom in the top left, a large yellow banana on the left, and a green broccoli floret in the bottom center. A thick orange horizontal bar is at the bottom left, and a vertical orange bar with a zigzag line is on the right.

# 03

## Unsupervised Learning

Segmenting users into clusters using order information to identify similarities between customers.

# Model Definition – feature selection

1. grouped data by user\_id
2. aggregated each feature as shown in the table below

	Features	Details
sum	d1, d2, d3,...,d21	indicates the number of products ordered from each department
mean	num_items	indicates the number of items ordered
	reord1	indicates the number of items re-ordered
	days_elapsed	indicates the number days since the previous order
last	order_number	Indicates the number of orders

# Model Definition – data transformation

Several features were transformed as shown in the table below:

	Features	Transformation
<b>Sparse data</b>	d1, d2, d3,...,d21	PCA dimensionality reduction First 5 components used (~53% of variance)
<b>Skewed data</b>	num_items	log transformation
	reord1	
	days_elapsed	
	order_number	



# Selecting the Algorithm



## Algorithm

The clustering algorithms below were tested to determine the optimal clusters. The k-means algorithm resulted in the greatest similarity scores.

- hierarchical clustering
- gaussian-mixture model
- DBSCAN
- **k-means** → best silhouette analysis results!

## Choosing k

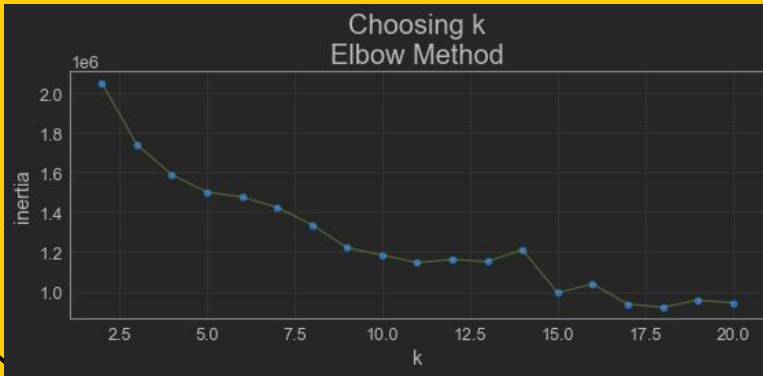
The number of clusters, **k**, was chosen based on the following methods:

- Elbow method
- Silhouette analysis



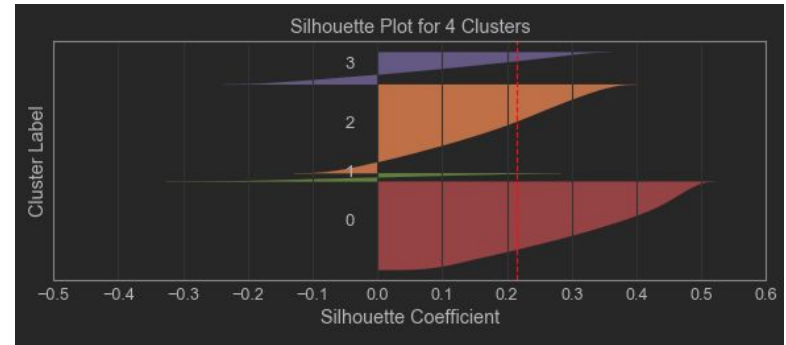
## Elbow Method

Elbow method narrowed search down to 4, 5, or 6 clusters.



## Silhouette Analysis

Silhouette analysis determined that 4 clusters results in the greatest intra-cluster similarity.







# Clustering Evaluation – number of items

Observations:

- **Cluster 1** users on average ordered and re-ordered the fewest number of items.
- **Cluster 2** users on average ordered and re-ordered the greatest number of items.



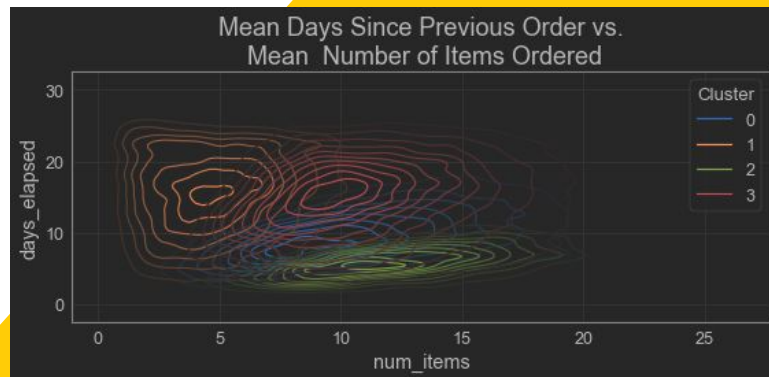


# Clustering Evaluation – days elapsed



## Observations

- **Cluster 1** users ordered the fewest items and had a wide range of order lag
- **Cluster 2** users placed orders with the least lag and had a wide range of number of items ordered.
- **Cluster 0** and **Cluster 3** have overlapping ranges of number of items ordered; however, Cluster 0 users have less order lag than Cluster 3 users.

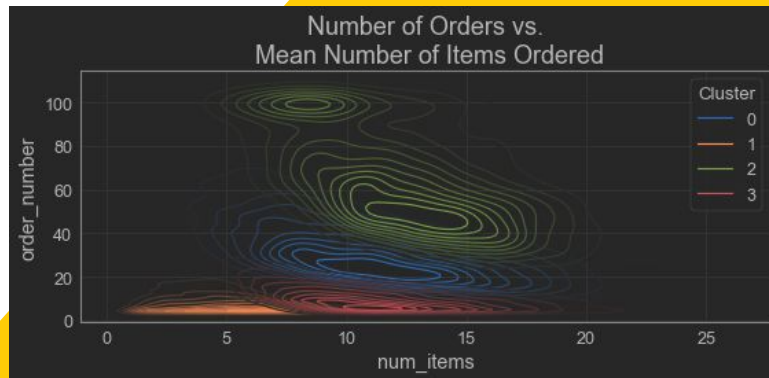




# Clustering Evaluation – number of orders

## Observations

- **Cluster 1** and **Cluster 3** users placed the fewest number of orders, though Cluster 3 users ordered more items on average
- **Cluster 2** users placed the greatest number of orders, though there is two densities of users within this cluster
  - Users that placed ~100 orders
  - Users that placed ~30~80 orders



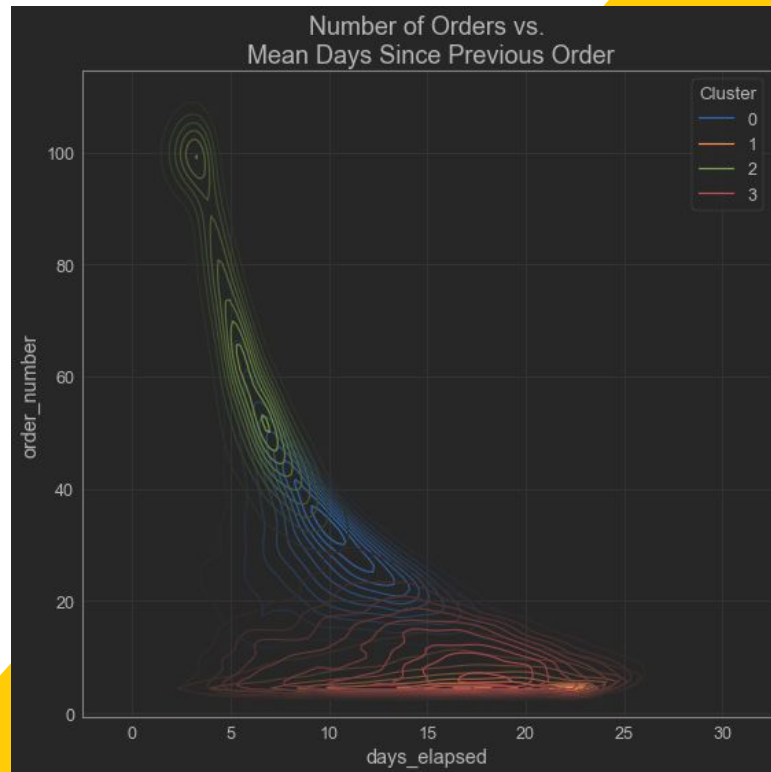


# Clustering Evaluation – number of orders



## Observations

- For **Cluster 0** and **Cluster 2** users, there is a negative correlation between the number of orders placed and the lag time between orders
- No such correlation is present for **Cluster 1** and **Cluster 3**



# Cluster Evaluation – summary

	Summary
Cluster 0	Semi-frequent shoppers who order a moderate number to many items
Cluster 1	Occasional shoppers who order very few items
Cluster 2	Frequent shoppers who order a moderate number to many items
Cluster 3	Occasional shoppers who order a moderate number of items



# 04

## Supervised Learning

Predicting the order frequency of each user using the order information.

# Model Definition – feature selection

1. grouped data by user\_id
2. aggregated each feature as shown in the table below
3. days\_elapsed, a continuous variable, selected as the output variable, requiring regression

	Features	Details
sum	d1, d2, d3,...,d21	indicates the number of products ordered from each department
mean	num_items	indicates the number of items ordered
	reord1	indicates the number of items re-ordered
	days_elapsed	indicates the number days since the previous order
last	order_number	indicates the number of orders

# Model Definition – data transformation

Due to issues faced, several features were transformed as shown in the table below:

	Features	Transformation
sparsity	d1, d2, d3,...,d21	PCA did not significantly impact regression results so for final analysis, PCA was skipped.
skew	num_items	log transformation
	reord1	
	days_elapsed	
	order_number	





# Selecting the Algorithm



## Algorithm

The supervised learning algorithms below were tested for sample sizes of up to 50%. The k-nearest-neighbor regressor was chosen as it resulted in the lowest RMSE.

- random forest
- gradient boosting
- support vector machine
- **knn** → lowest root mean squared error!

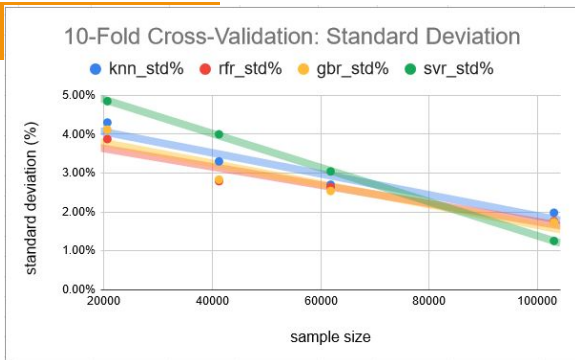
## Parameters

Searches were performed to determine the optimal parameter values.

- `n_neighbors = 50`
- `weights = 'distance'`
- `algorithm = 'ball tree'`



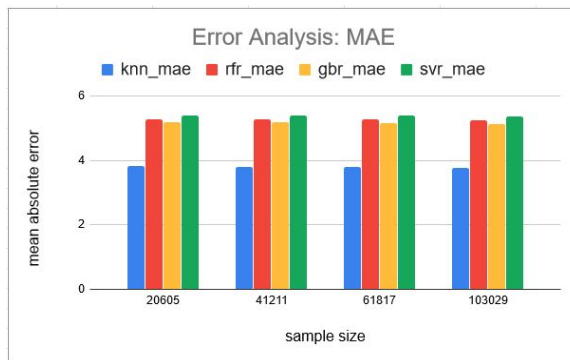
# Comparing Algorithms



## Standard Deviation

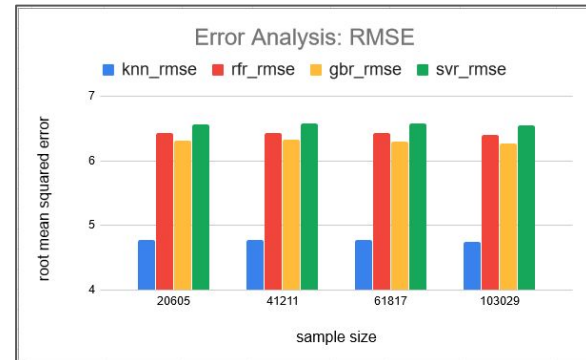
Standard deviation across folds decreases with sample size.

Similar values across algorithms (1.25% - 1.98%) at 50% sampling rate



## Mean Absolute Error

MAE is consistent across sample size with KNN algorithm having the lowest error.



## Root Mean Squared Error

RMSE is also consistent across sample size with KNN algorithm again having the lowest error.

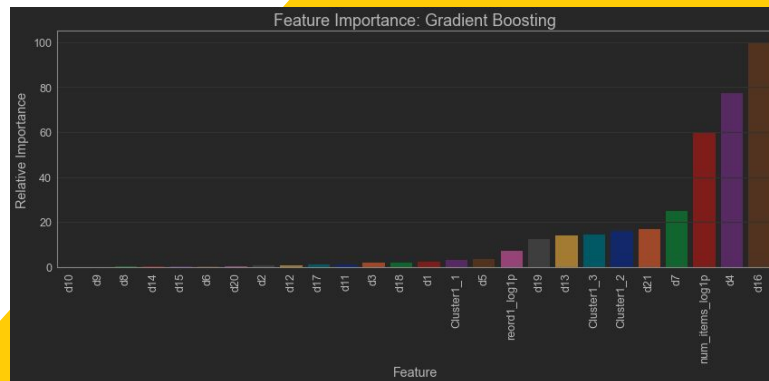
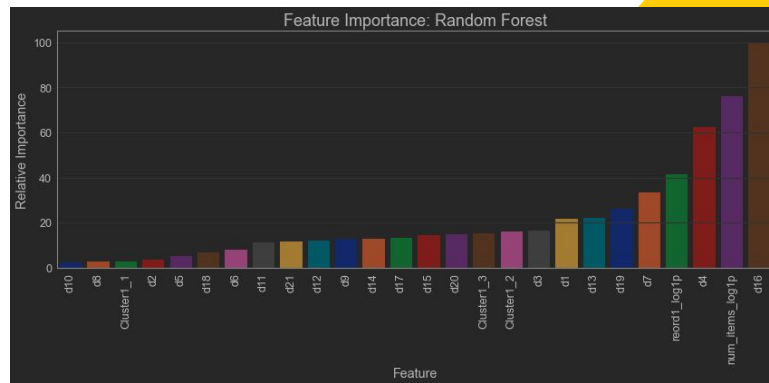


# Error Evaluation – feature importance



Observations:

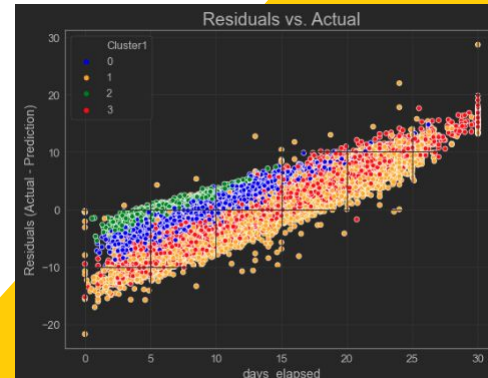
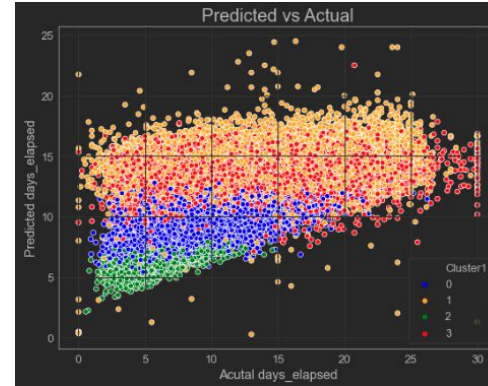
- Notably, two of the most important features are the departments with the most purchases (d4, d16)
- Both the random forest and gradient boosting algorithms consider d16, d4, and num\_items (log transformed) as the 3 most important features.



# Error Evaluation – residuals

Observations:

- **Cluster 2** and **Cluster 0** appear to have a slight linear relationship between predicted and actual values of order lag time (days\_elapsed)
- **Cluster 2** and **Cluster 0** also appear to have a narrower band of residual values relative to Cluster 1 and Cluster 3



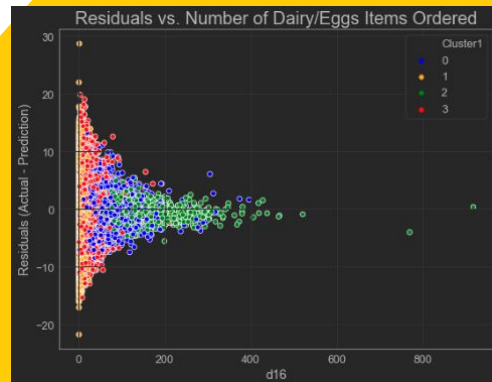
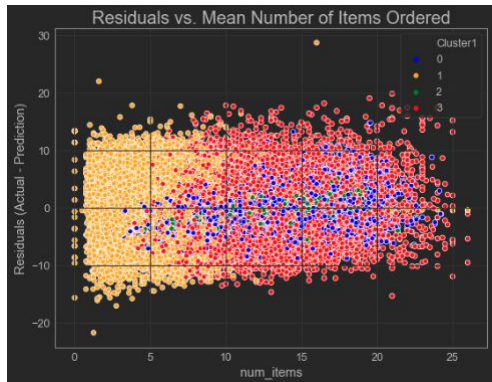


# Error Evaluation – residuals



Observations:

- The magnitude of the residuals is independent from the mean number of items ordered
- The magnitude of the residuals are smaller for greater total number of items ordered
- Data indicates users who order more items more frequently (Cluster 2) are easier to predict



The background is a vibrant collage of geometric shapes and food illustrations. A large yellow triangle occupies the right side, while a pink triangle is on the left. Scattered throughout are various elements: a yellow circle at the top, a green broccoli illustration at the bottom center, a yellow banana on the left, a small orange mushroom in the top left, and several small black dots and white geometric shapes (triangles, zig-zags) scattered across the white background.

# 05

## Conclusions

and future work

# Summary

	Summary	Ideas
<b>Cluster 0</b> (17.3%)	<ul style="list-style-type: none"><li>• Semi-frequent shoppers</li><li>• Shoppers order moderate number to many items</li><li>• Moderate difficulty in predicting lag time</li></ul>	<ul style="list-style-type: none"><li>• Provide rewards for every nth order to increase loyalty</li></ul>
<b>Cluster 1</b> (37.1%)	<ul style="list-style-type: none"><li>• Occasional shoppers</li><li>• Shoppers order very few items</li><li>• Most difficult to predict lag time</li></ul>	<ul style="list-style-type: none"><li>• Place second in priority</li><li>• Offer discounts or provide rewards for at least first 4 orders</li><li>• Target with weekly reminders and local specials</li></ul>
<b>Cluster 2</b> (4.7%)	<ul style="list-style-type: none"><li>• Frequent shoppers</li><li>• Shoppers order a moderate number to many items</li><li>• Least difficult to predict lag time</li></ul>	<ul style="list-style-type: none"><li>• Continue monitoring cluster for sudden changes in behavior</li></ul>
<b>Cluster 3</b> (40.8%)	<ul style="list-style-type: none"><li>• Occasional shoppers</li><li>• Shoppers order a moderate number of items</li><li>• Difficult to predict lag time</li></ul>	<ul style="list-style-type: none"><li>• Prioritize targeting this cluster</li><li>• Target with weekly reminders and local specials</li><li>• Repeat analysis on this cluster using products or aisles as features rather than departments for increased granularity in interests</li></ul>

# Future Work

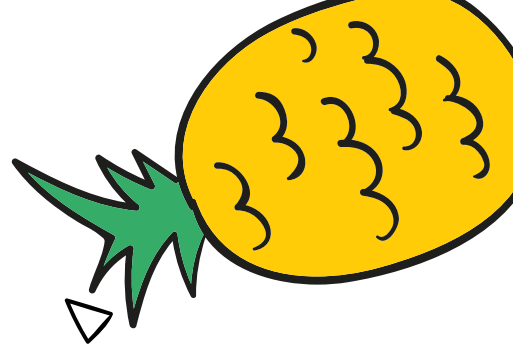
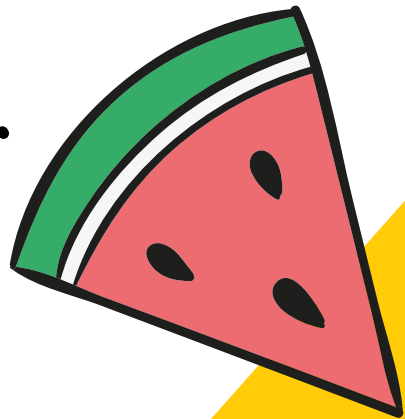
## Grocery List

- Address sparsity in features more effectively
- Repeat clustering and regression analysis on subset of data (Cluster 3)
- Repeat clustering and regression analysis using products or aisles as features instead of departments





**Questions?**



...

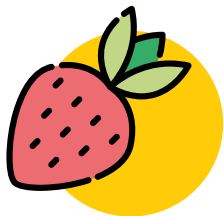
The background is a vibrant collage of geometric shapes and food illustrations. A large yellow triangle occupies the right side, while a pink triangle is on the left. Scattered throughout are various elements: a yellow circle at the top center, a green broccoli floret at the bottom center, a yellow banana on the left, a small orange mushroom in the top left, and several small black dots and white outlines of triangles and zig-zags.

06

**Backup Slides**

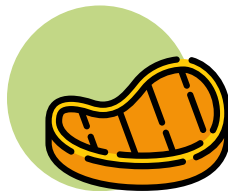
...

# What Is Machine Learning?



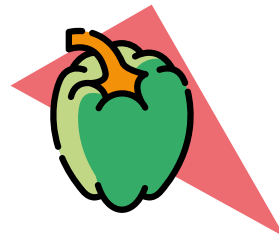
## Machine Learning

Computer algorithms that learn and improve from experience without being explicitly programmed.



## Unsupervised

Type of machine learning that searches for patterns in a data set with no pre-existing labels.

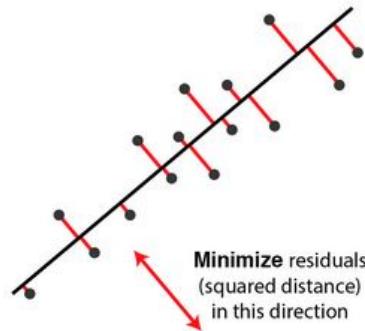
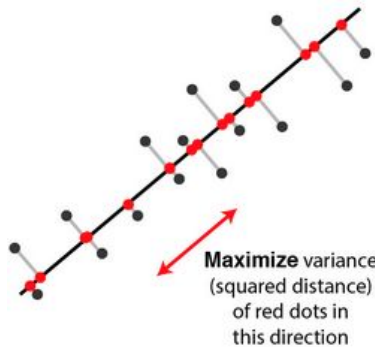


## Supervised

Type of machine learning that predicts an output given a set of inputs based on example input-output pairs

# PCA – Principal Components Analysis

- dimensionality reduction method that aims to minimize information loss / variance
- generates low-dimensional representations of high-dimensional data called **principal components**
- **PC1**, the first principal component is the directional line, or vector that:
  - **minimizes** the squared distances between data points and their projections onto line
  - **maximizes** the squared distances between the projected points and the origin point
- all successive component must be perpendicular to previous components
  - PC2 perpendicular to PC1
  - PC3 perpendicular to PC2 and PC1
  - etc.

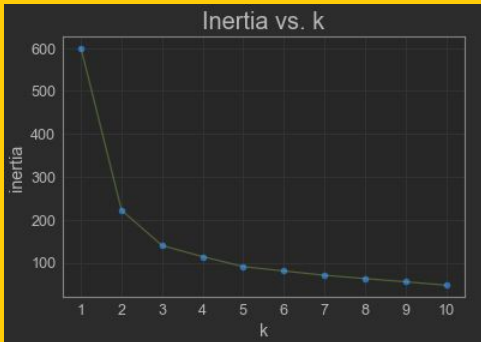


# Optimizing the Number of Clusters

## Elbow Method

The elbow method plots the **inertia**, the sum of squared distances of the samples from the cluster centers, vs. the number of clusters, **k**.

The optimal **k**, is the value at which the rate of decrease in inertia becomes more linear. Visually, this appears as the bend in the plot, much like the elbow of a bent arm, hence the name.



## Silhouette Analysis

The **silhouette coefficient** measures the similarity of data points within a cluster with one another. It is calculated as follows:

$$\frac{b_i - a_i}{\max(b_i, a_i)}$$

where for each data point **i**,

- **$a_i$**  = mean distance between **i** and all data points in its cluster
- **$b_i$**  = mean distance between **i** and all data points in neighboring clusters

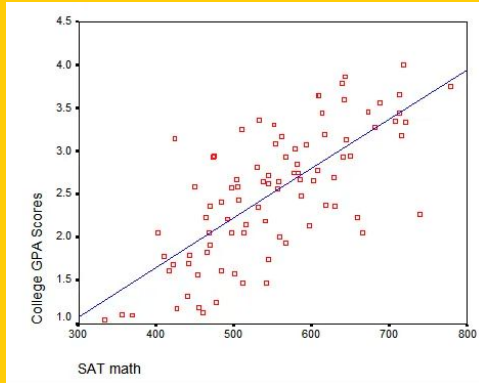
The **silhouette average** is the average of all silhouette coefficients of all of the data points.

# Types of Supervised Learning

## Regression

Type of supervised learning in which the task is to predict the values of a **continuous** outcome variable.

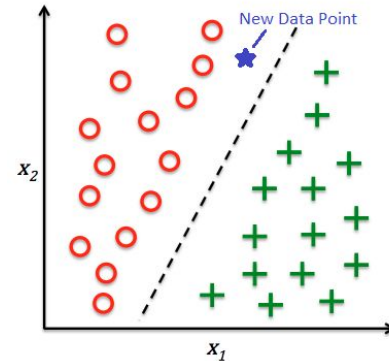
A continuous variable, such height or revenue, can take on an infinite number of values. In regression, we are trying to **quantify**.



## Classification

Type of supervised learning in which the task is to predict the values of a **categorical** outcome variable.

A categorical variable, such as hair color or car model, can take on only a limited number of values. In classification, we are trying to **select**.



# Thanks!

Do you have any questions?

[youremail@freepik.com](mailto:youremail@freepik.com)

+91 620 421 838

[yourcompany.com](http://yourcompany.com)

CREDITS: This presentation template was created by Slidesgo, including icons by Flaticon, and infographics & images by Freepik.

*Please keep this slide for attribution.*

