

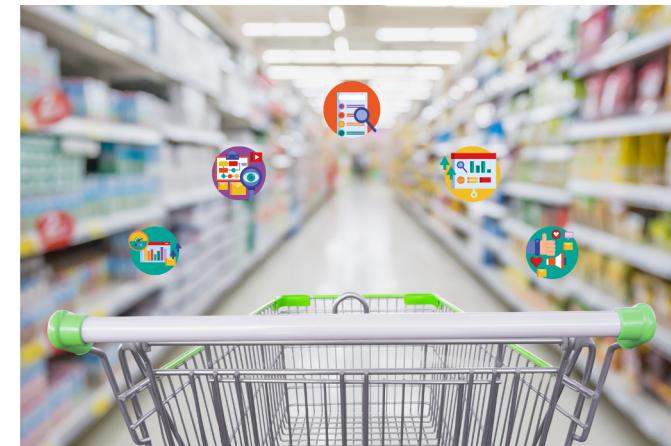


Market Basket Analysis

Group 5
12/09/2021

Story and Motivation

- Data mining technique used to better understand **customer purchasing patterns**
- Key techniques used to **uncover associations** between items
- **Analyzes customer buying habits** by finding associations between the different items that customers place in their “shopping baskets”

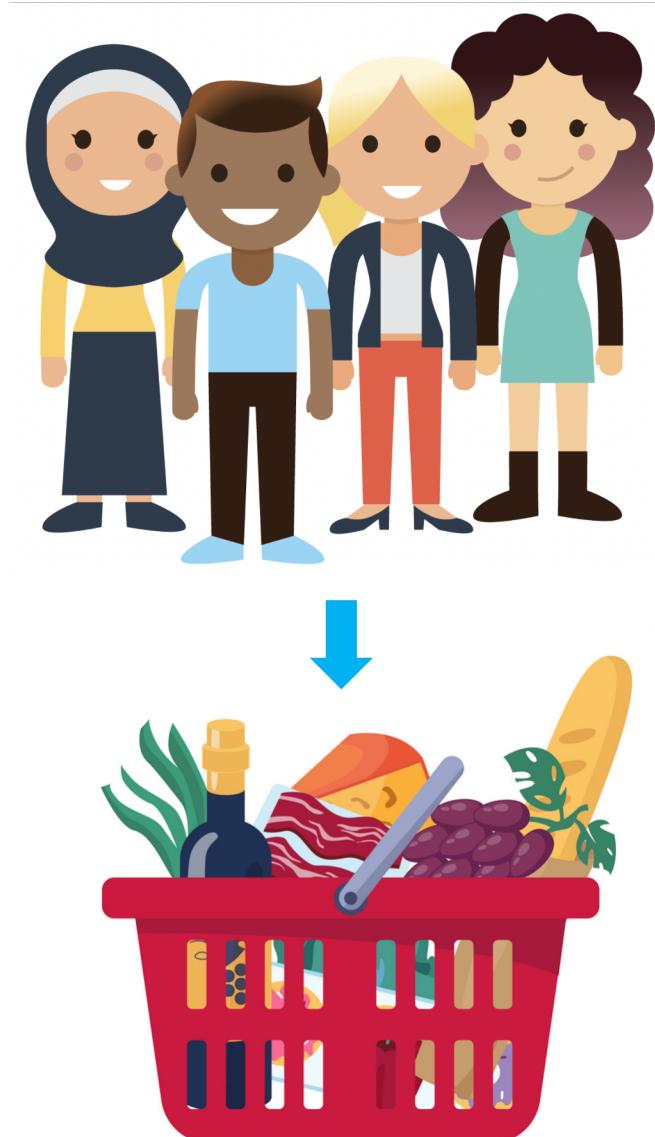


- Gives insight into which items are frequently bought together
- Helps improve marketing strategy
- Helps customer experience

Example: Instacart.com

The dataset

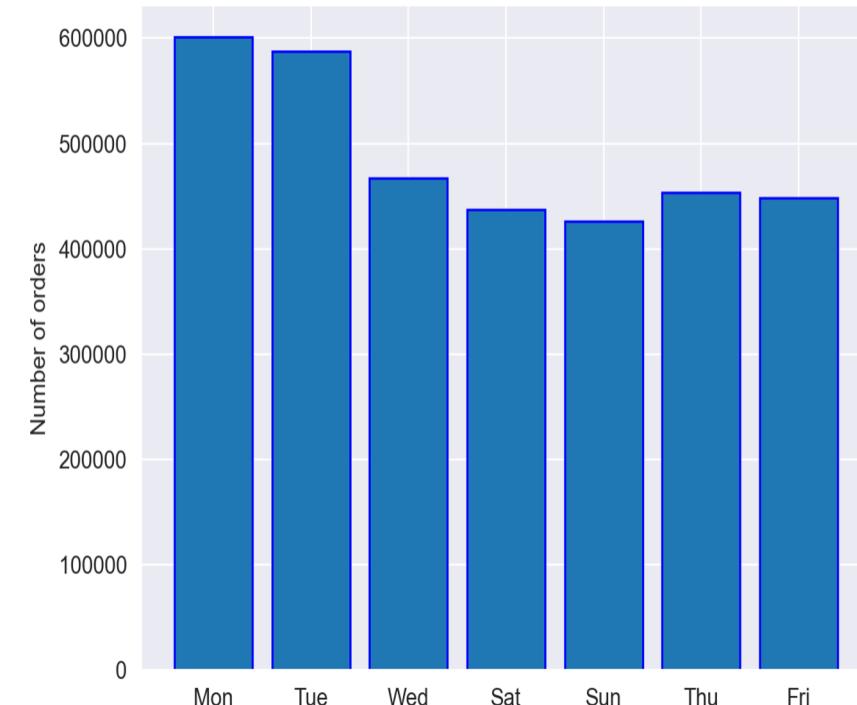
- Dataset used:
 - Instacart - market basket analysis
- Tables included:
 - Aisles
 - Products
 - Departments
 - Orders
 - Reorders product
 - Consumption habits (time)



Market basket analysis

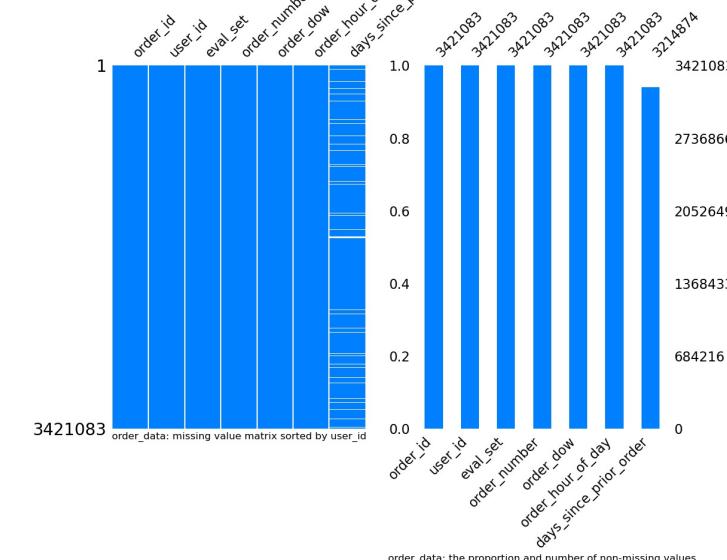
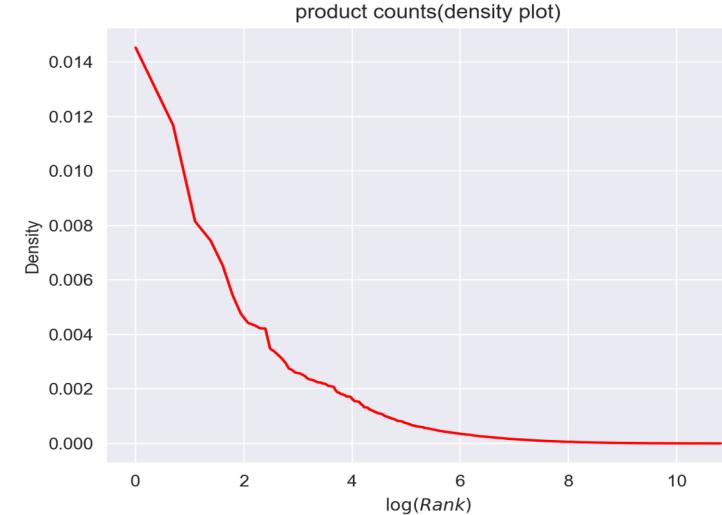
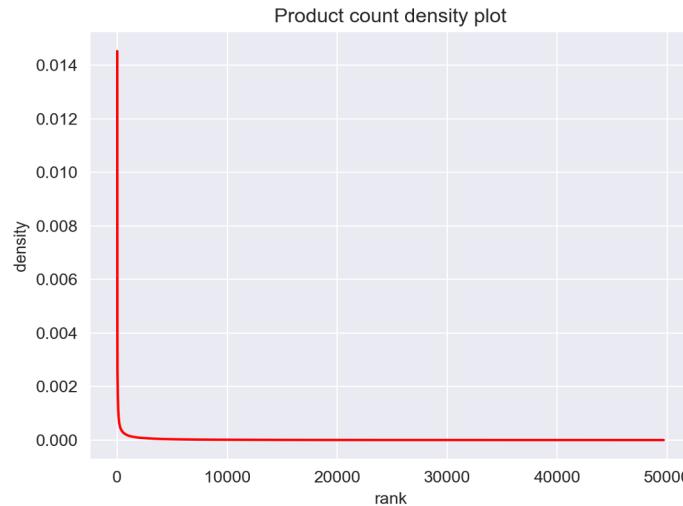
Exploratory Data Analysis

```
top 10 frequent products bought
product_name
0          Banana
1  Bag of Organic Bananas
2  Organic Strawberries
3  Organic Baby Spinach
4  Organic Hass Avocado
5      Organic Avocado
6          Large Lemon
7  Strawberries
8          Limes
9  Organic Whole Milk
10  Organic Raspberries
```



Market basket analysis

Exploratory Data Analysis



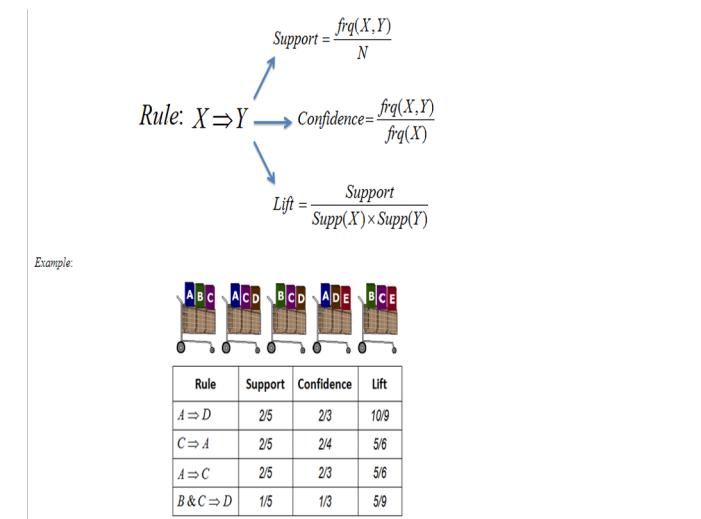
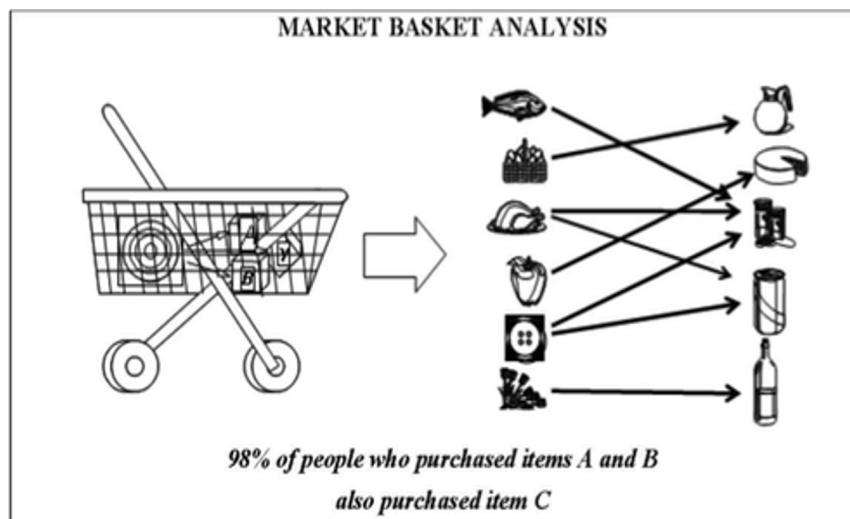
Check the missing data

- Association rule mining:
 - Technique to identify underlying relations between various items
- Most common approach: Apriori algorithm
- Main idea: “All non-empty subsets of a frequent itemset must also be frequent”
- Bottom-up approach
 - Start from every individual item
 - Generate candidates by self-joining
 - Itemsets that contain infrequent subsets are pruned
 - Repeat until no more successful itemsets are formed



Association Mining Rule

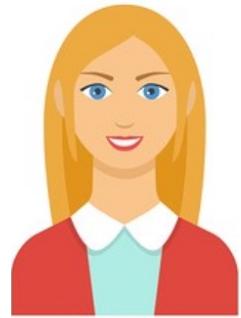
- One of the key techniques in Market Basket Analysis
- Market basket analysis tells you about items that are “frequently bought together”
- Eg., Amazon.com - “Customers also bought”,etc
- Idea of Market Basket Analysis: if item x is bought, item/ itemset y is bound to be or not be bought
- eg. , if one buys bread, chances of buying jam/butter is high



Market basket analysis dataset

LSH and KMeans clustering (Story Background)

Group1



Reorder products



?
may also like

Group2



Reorder products



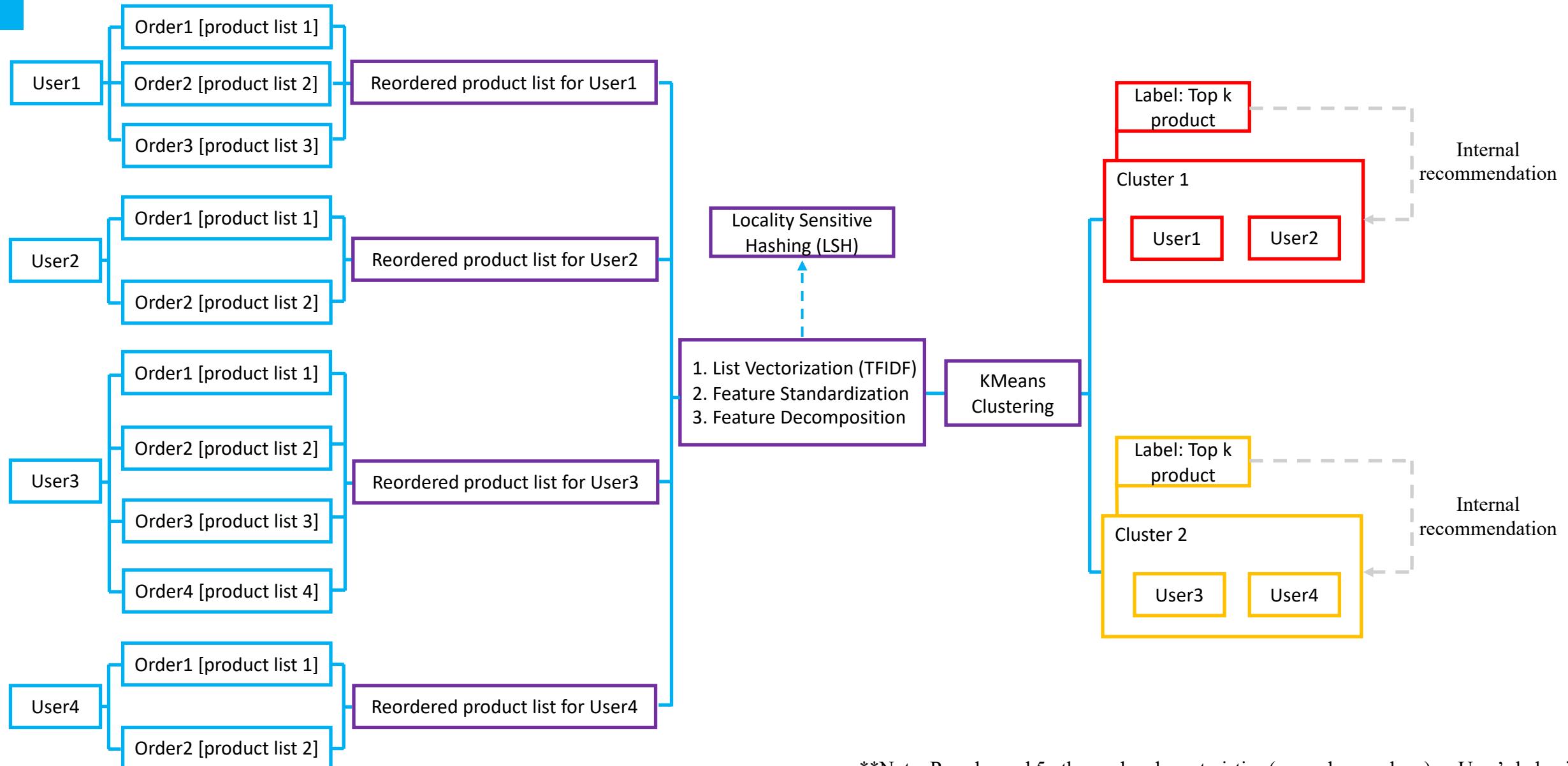
?
may also like

Each user has lots of orders.
The **reorder products** can be seen as user's favorites, which can reflect user behavior.

1. Based on the reorder product (user behavior) of each user, we can **classify user into different groups**;
2. Based on the **top-K favorite products of each group**, we can do an **internal recommendation**.

Market basket analysis dataset

LSH and KMeans clustering (workflow)



Market basket analysis dataset

LSH and KMeans clustering (results: Feature Engineering and LSH)

Feature Generation

```
*****  
1000_tfidf_features  
user_id 10 100 12 21 50 70 85 90 93 aa acai added ... yellow yerba yo yobaby yoghurt yogurt yokids yukon zbar zbars zero zucchini  
user_id ...  
1 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00000 0.0 ... 0.0 0.0 0.0 0.0 0.00000 0.00000 0.0 0.0 0.0 0.0 0.0 0.183683 0.0  
2 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.114357 0.0 ... 0.0 0.0 0.0 0.0 0.140415 0.123297 0.0 0.0 0.0 0.0 0.0 0.000000 0.0  
3 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00000 0.0 ... 0.0 0.0 0.0 0.0 0.00000 0.00000 0.0 0.0 0.0 0.0 0.0 0.000000 0.0  
4 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00000 0.0 ... 0.0 0.0 0.0 0.0 0.00000 0.00000 0.0 0.0 0.0 0.0 0.0 0.000000 0.0  
5 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.00000 0.0 ... 0.0 0.0 0.0 0.0 0.00000 0.206131 0.0 0.0 0.0 0.0 0.0 0.000000 0.0  
[5 rows x 1000 columns]  
*****  
LSH_similar_pairs_result_hashbucket  
hash_values similar_user_sets  
0 7c1bec451b36871428c8b1af0e9ebd04 [1, 20, 24, 57, 81, 83, 107, 117, 122, 138, 18...  
1 39ce64aef3c638d2fe0523aa4a08caf [2, 3, 7, 10, 47, 49, 59, 70, 101, 103, 144, 1...  
2 6fe27aea00bf1f35f01a604e97674ae0 [4, 19, 28, 61, 116, 291, 377, 389, 419, 431, ...  
3 69a4d6c6d8a2778f051127c13d9c671 [5, 48, 51, 54, 73, 76, 128, 153, 168, 198, 21...  
4 74b870fa26e3d49d62a395a32f2223a1 [6, 44, 74, 84, 100, 133, 186, 216, 217, 231, ...  
*****  
1000_tfidf_features_plus_5_user_behavior_features  
user_id order_dow order_hour_of_day days_since_prior_order order_num product_num 10 100 12 21 ... yobaby yoghurt yogurt yokids yukon zbar zbars zero zucchini  
1 -0.147273 -1.123285 0.591456 -0.299025 -0.333873 -0.109229 -0.485404 -0.15738 -0.144355 ... -0.119367 -0.162370 -0.682866 -0.182337 -0.153452 -0.083181 -0.08657 4.191568 -0.27855  
2 -0.892007 -1.248441 0.186580 -0.116386 -0.007183 -0.109229 -0.485404 -0.15738 -0.144355 ... -0.119367 3.605954 1.547926 -0.182337 -0.153452 -0.083181 -0.08657 -0.181701 -0.27855  
3 -1.696397 1.140663 -0.584408 -0.177265 -0.245918 -0.109229 -0.485404 -0.15738 -0.144355 ... -0.119367 -0.162370 -0.682866 -0.182337 -0.153452 -0.083181 -0.08657 -0.181701 -0.27855  
4 1.080598 0.586981 0.775148 -0.786065 -0.585174 -0.109229 -0.485404 -0.15738 -0.144355 ... -0.119367 -0.162370 -0.682866 -0.182337 -0.153452 -0.083181 -0.08657 -0.181701 -0.27855  
5 -1.479201 1.560088 0.025855 -0.664305 -0.503501 -0.109229 -0.485404 -0.15738 -0.144355 ... -0.119367 -0.162370 3.046620 -0.182337 -0.153452 -0.083181 -0.08657 -0.181701 -0.27855  
[5 rows x 1005 columns]  
*****  
decomposed_1005features_for_KMeans  
PC1 PC2  
user_id  
1 -4.237737 3.268695  
2 3.068878 2.148310  
3 0.112483 -1.852136  
4 -3.775791 -0.272410
```

Be further used for KMeans

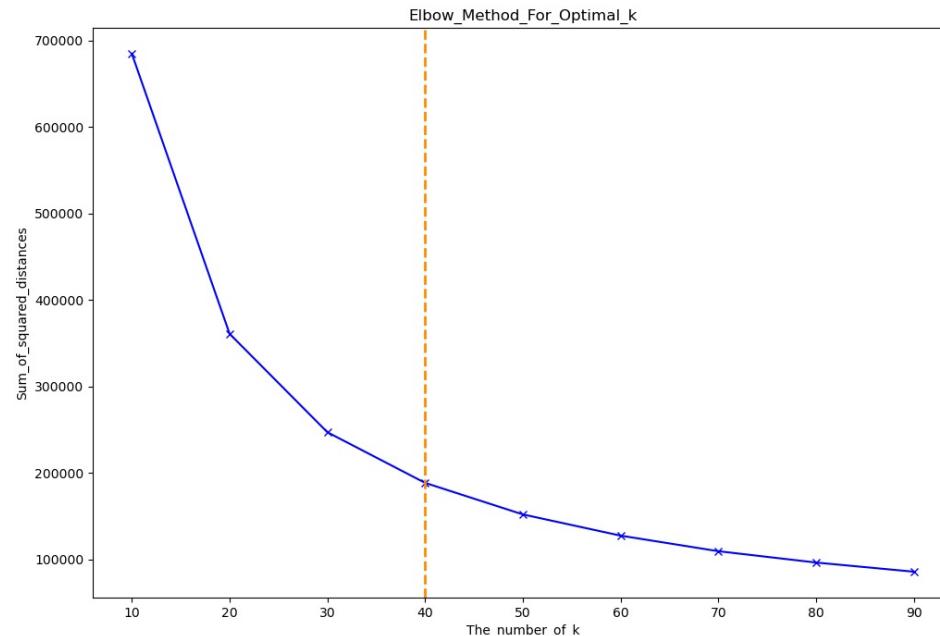
```
3 user_product_ori[user_product_ori.index==1].iloc[0, 0]  
'Soda, Original Beef Jerky, Pistachios, Organic String Cheese, Bag of Organic Bananas, Soda, Original Beef Jerky, Pistachios, Soda, Pistachios, Original Beef Jerky, Organic String Cheese, Cinnamon Toast Crunch, Soda, Original Beef Jerky, Pistachios, Organic String Cheese, XL Pick-A-Size Paper Towel Rolls, Organic Half & Half, Zero Calorie Cola, Organic String Cheese, Soda, Pistachios, Original Beef Jerky, Soda, Original Beef Jerky, Aged White Cheddar Popcorn, Soda, Zero Calorie Cola, Organic String Cheese, Pistachios, Cinnamon Toast Crunch, Original Beef Jerky, Original Beef Jerky, Soda, Pistachios, Organic String Cheese, Soda, Original Beef Jerky, Pistachios, Organic String Cheese'
```

```
1 user_product_ori[user_product_ori.index==57].iloc[0, 0]  
'Grated Parmesan Cheese, Whipped Cream Cheese, Whipped Cream Cheese, Roma Tomato, Grated Parmesan Cheese, Baked Aged White Cheddar Rice and Corn Puffs, Organic Roasted Dandelion Root Herbal Tea, Banana, Whipped Cream Cheese, Lemon Zinger Herbal Tea, Baked Aged White Cheddar Rice and Corn Puffs, Spaghetti, Organic Roasted Dandelion Root Herbal Tea, Whipped Cream Cheese, Banana, Distilled Water, Red Leaf Lettuce, 90% Lean Ground Beef, Artichoke Hearts, Whipped Cream Cheese, Baked Aged White Cheddar Rice and Corn Puffs, Banana, Whipped Cream Cheese, Roma Tomato, Organic Roasted Dandelion Root Herbal Tea, Lemon Zinger Herbal Tea, Whipped Cream Cheese, Lemon Zinger Herbal Tea, Red Leaf Lettuce, Banana, Whipped Cream Cheese, Grated Parmesan Cheese, Whipped Cream Cheese, Banana, Organic Roasted Dandelion Root Herbal Tea, Lemon Zinger Herbal Tea, Cucumber Kirby, Russet Potato, Whipped Cream Cheese, Grated Parmesan Cheese, Artichoke Hearts'
```

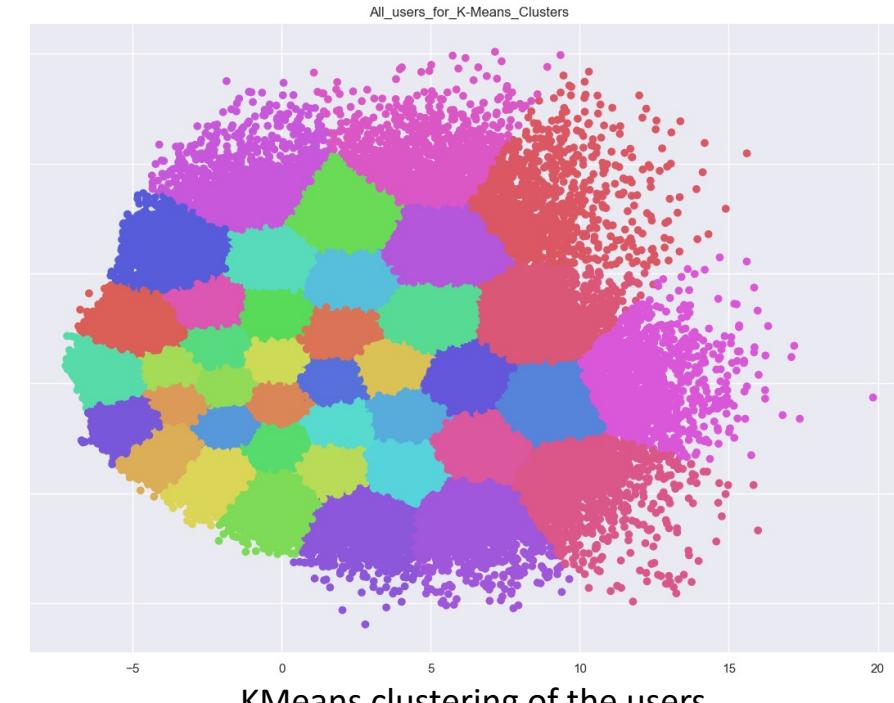
Roughly analysis by LSH:
The word "organic" and "cheese" may be two word put User1 and User57 together.

Market basket analysis dataset

LSH and KMeans clustering (results: Optimal K and KMeans Clustering)



Elbow Method: Screen for optimal-K



KMeans clustering of the users

Since we use **unsupervised Kmeans**, we have no idea about the true clustering, which means **how many clusters we need**.

We use **elbow method to look for the optimal K**:

1. Calculate the **Within-Cluster-Sum of Squared Errors** (WSS: predict cluster center) for different k;

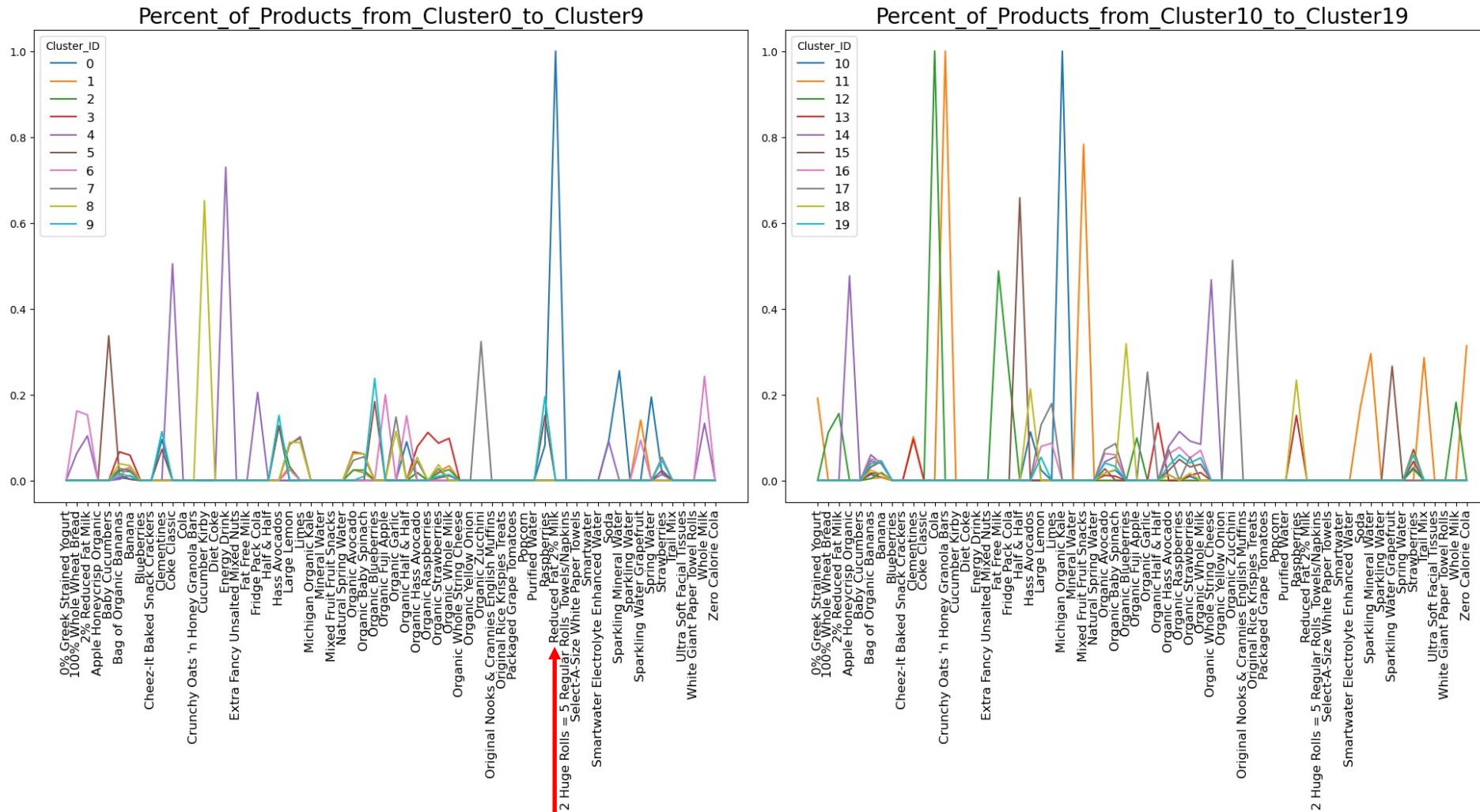
2. Choose the k for which **WSS starts to diminish (elbow)**. It could be a reasonable k for clustering.

Note: **Silhouette Analysis (the difference between its mean distance from points in the nearest cluster and its mean distance from points in its own cluster).

3. Finally, we **classify the users into 40 clusters** based on their behavior.

Market basket analysis dataset

LSH and KMeans clustering (results: TopK product for Internal Recommendation)



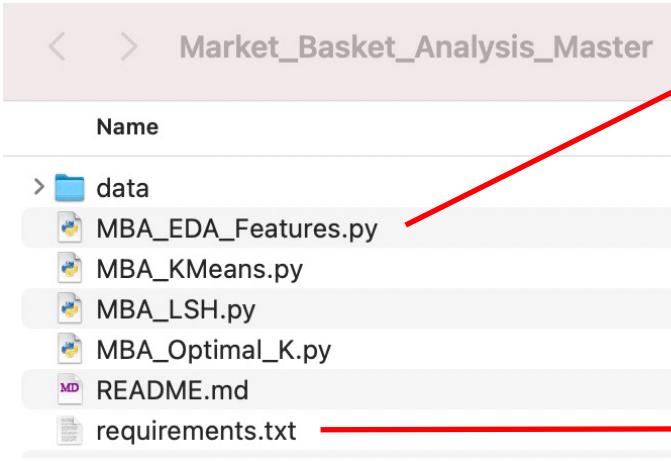
Top10 favorite products in each cluster comparing to other clusters

For example: In cluster 0: Reduced Fat Milk is the top 1 product, we could recommend this product to the other users of cluster 0 who have not bought it yet.

Market basket analysis dataset

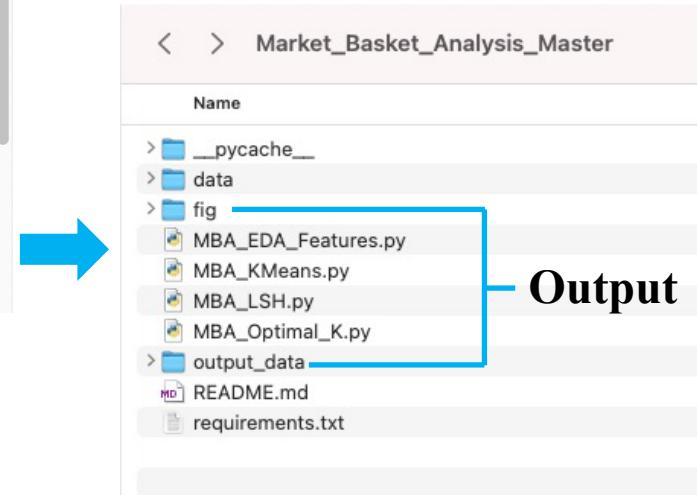
LSH and KMeans Scripts

Working directory for the Scripts



Instruction in each script

```
"""\nAUTHOR: Tong Wang, CENTER FOR COMPUTATIONAL AND INTEGRATIVE BIOLOGY\nDATE: 12/2021\nDESCRIPTION: This script is the first part used for Instacart-Market-Basket-analysis. Here we focus on\ntfidf features generation, and LSH analysis.\nUSAGE(In Macbook/Amarel): (please install anaconda 4.10.3 at first, the python version is 3.9.7)\n0. change path to the work directory\n1. Creating an environment with commands (note: replace myenv with the environment name):\n    conda create --name myenv\n2. Activate the environment with commands:\n    conda activate myenv\n3. Install packages in conda environment using requirements.txt:\n    conda install --file requirements.txt\n4. Execute the script:\n    python MBA_EDA_Features.py\n"""\n\n# requirements.txt\nmatplotlib==3.3.4\nnumpy==1.20.1\npandas==1.2.4\nscikit-learn==0.24.1\nseaborn==0.11.1\nmissingno==0.4.2\nIPython==7.22.0
```



Run the demo in the shell:

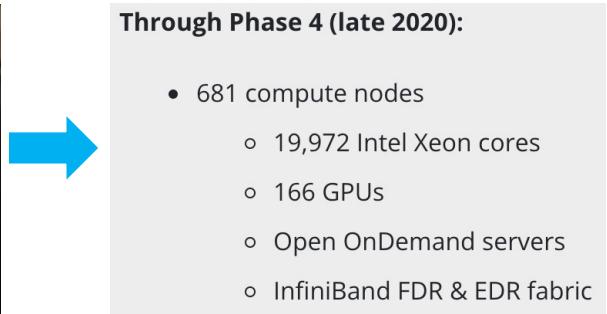


More details of the script instruction, please check the following Github link:

https://github.com/wtcruijff/Market-Basket-Analysis-Big-Data-Course/tree/main/Market_Basket_Analysis_Master

Market basket analysis dataset

Cloud Computing (Amarel Cluster)



Amarel (**owned by Rutgers**) is a shared community-owned advanced computing environment available to any investigator or student with projects requiring research computing resources.

<https://oarc.rutgers.edu/resources/amarel/>

Cloud Computing workflow on Amarel:

1. Upload the data and script on the work directory
2. Install anaconda and setting up environment by the shell

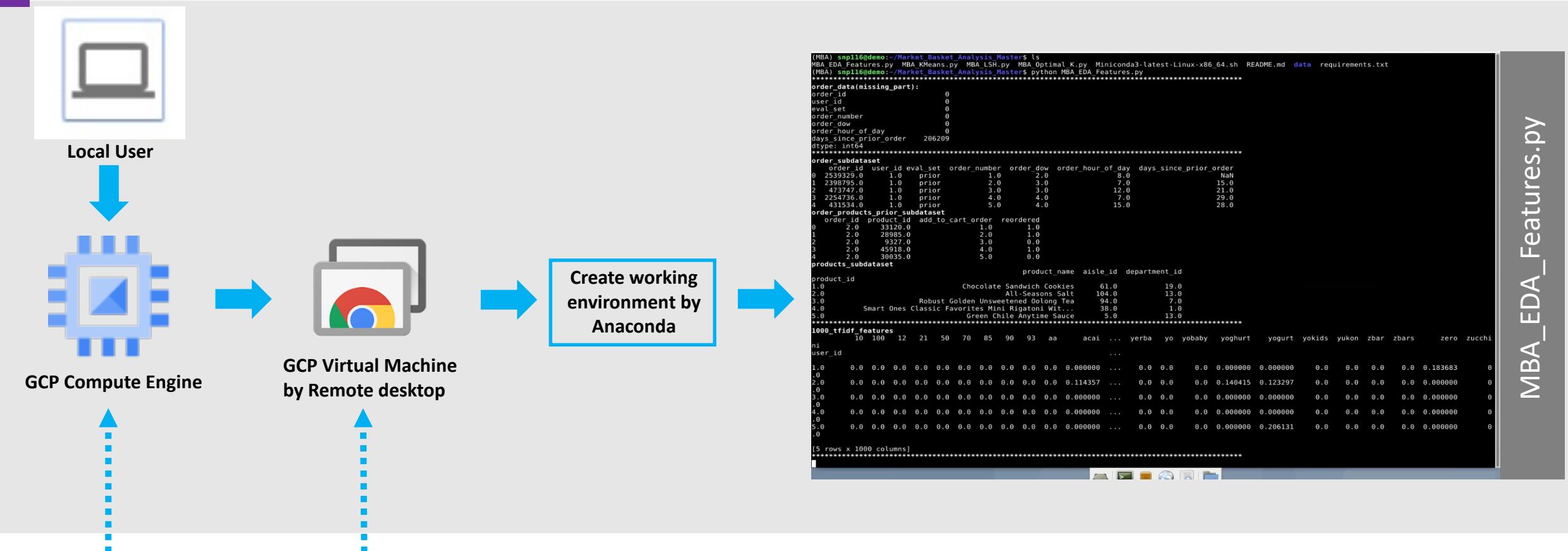
3. Create a VM with certain parameters (cores and memory)

```
(base) [tw543@ha08013 ~] $ conda activate marketbasket
(marketbasket) [tw543@ha08013 ~] $ cd ~/Market_Basket_Analysis_Master
(marketbasket) [tw543@ha08013 ~] $ cd Market_Basket_Analysis_Master
(marketbasket) [tw543@ha08013 ~] $ python MBA_EDA_Features.py
*****
order_datamissing_part:
user_id 0
eval_set 0
order_number 0
order_dow 0
order_hour_of_day 0
days_since_prior_order 296209
dtype: int64
*****
order_subdataset
order_id user_id eval_set order_number order_dow order_hour_of_day days_since_prior_order
0 2539329 1 prior -1.0 -2 8 NaN
1 2398795 1 prior 2.0 3 7 15.0
2 2254747 1 prior 3.0 3 12 21.0
3 2254736 1 prior 4.0 4 7 29.0
4 431534 1 prior 5.0 4 15 28.0
*****
order_products_prior_subdataset
order_id product_id add_to_cart_order reordered
0 2 33120 1 1
```

4. Run the script on Amarel

Market basket analysis dataset

Cloud Computing (Google Cloud Platform: GCP)



Live coding instruction, please check the following YouTube links:

- [1 https://www.youtube.com/watch?v=CM-wXdwAI5E](https://www.youtube.com/watch?v=CM-wXdwAI5E)
- [2 https://www.youtube.com/watch?v=ijxXGw66Urc](https://www.youtube.com/watch?v=ijxXGw66Urc)
- [3 https://www.youtube.com/watch?v=r1n6LMsV_RY&t=240](https://www.youtube.com/watch?v=r1n6LMsV_RY&t=240)



THANK YOU

Acknowledgement: Many Thanks to Dr. Shende for the kind guidance!

Group Member:



Iswarya Madhu Desikan



Shreya Patel



Tong Wang



Arkit Pawar