



*Data Project*

SR Data Analyst  
Bruno Bucalon Serra

# INDEX

1. Problem Definition
  - 1.1 Dataset Attributes
  - 1.2 Objective and Key Questions
  - 1.3 Technologies and Programming Language
2. First Steps
3. Data GAP
4. Data Analysis
  - 4.1 Valuable Users
  - 4.2 Channels & Valuable Users
  - 4.3 Conversion Rate
  - 4.4 Recommendations
5. Predictive Modeling

# 1. PROBLEM DEFINITION

# 1. PROBLEM DEFINITION

## 1.1 DATASET ATTRIBUTES

Acquisitions	Visits	Campaign (DIM)	Purchases
userID	userID	SEASON	userID
ACQUISITION_DATE	SEASON	CAMPAIGN_ID	CAMPAIGN_ID
ORDER_SEASON	CAMPAIGN_ID	CAMPAIGN_TITLE	CAMPAIGN_TITLE
PLAN_CODE	WINDOW_NUMBER	CAMPAIGN_TYPE	SKU
SAME_DAY_REFUND	CAMPAIGN_TITLE		CART_QUANTITY
CURRENT_SUBSCRIPTIONS_STATE	CAMPAIGN_TYPE		MSRP
AGE	WINDOW_START_TS_PST		COGS
USER_REGION	WINDOW_END_TS_PST		SALE_PRICE
LTV14	CHANNEL		CATEGORY
	NEW_OR_RETURNING_MEMBER		
	VISIT_CNT		

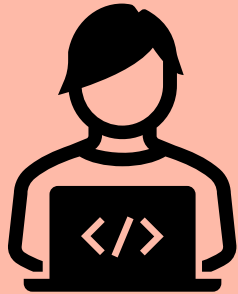
# 1. PROBLEM DEFINITION

## 1.2 OBJECTIVES AND KEY QUESTIONS

- Who are our most valuable members? How are they different from our less valuable members?
- Which channels are best at driving these members to sales?
- Where are purchase conversion rates strong?
- Based on your analysis, provide actionable recommendations to improve the conversion rate from visits to purchases. Explain the rationale behind each recommendation.

# 1. PROBLEM DEFINITION

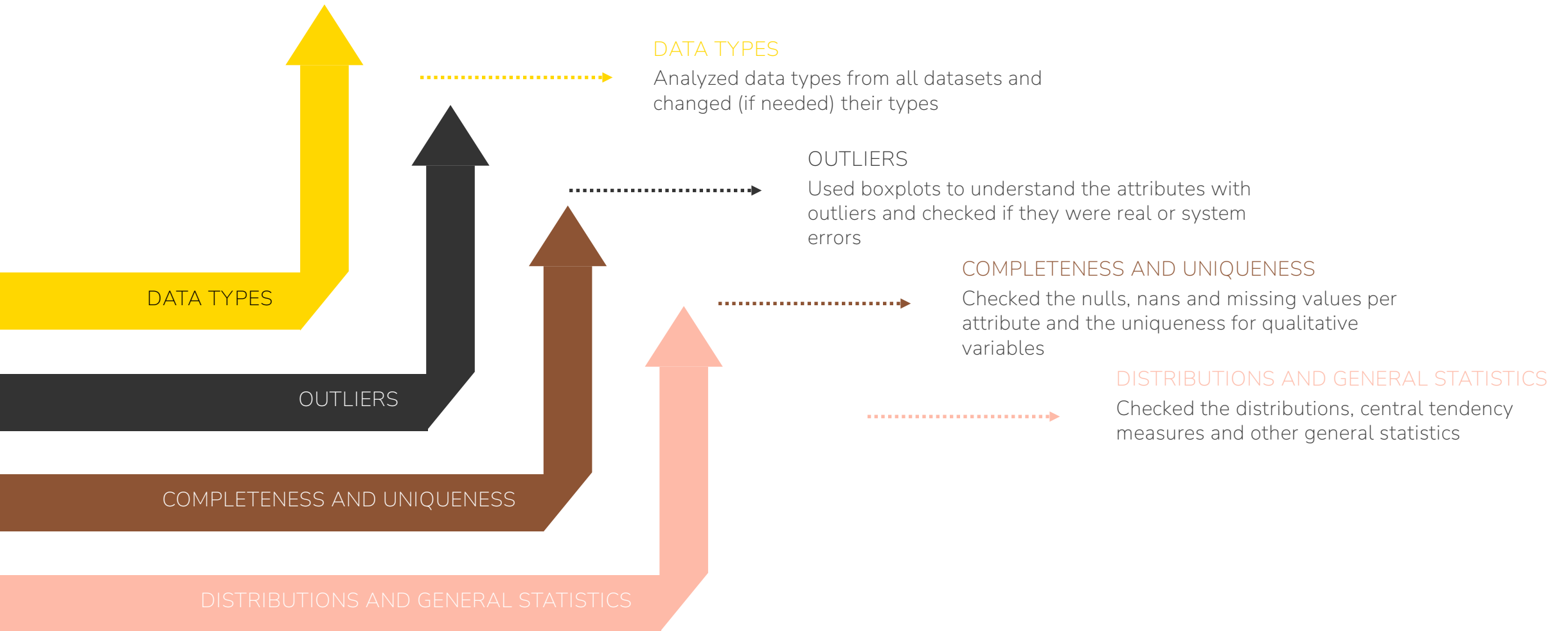
## 1.3 TECHNOLOGIES AND PROGRAMMING LANGUAGE



- Environment using Anaconda
- All analysis performed using Python
- Used libraries: NumPy, Pandas, Matplotlib, Sklearn
- Visual Studio Code used as IDE

## 2. FIRST STEPS

## 2. FIRST STEPS

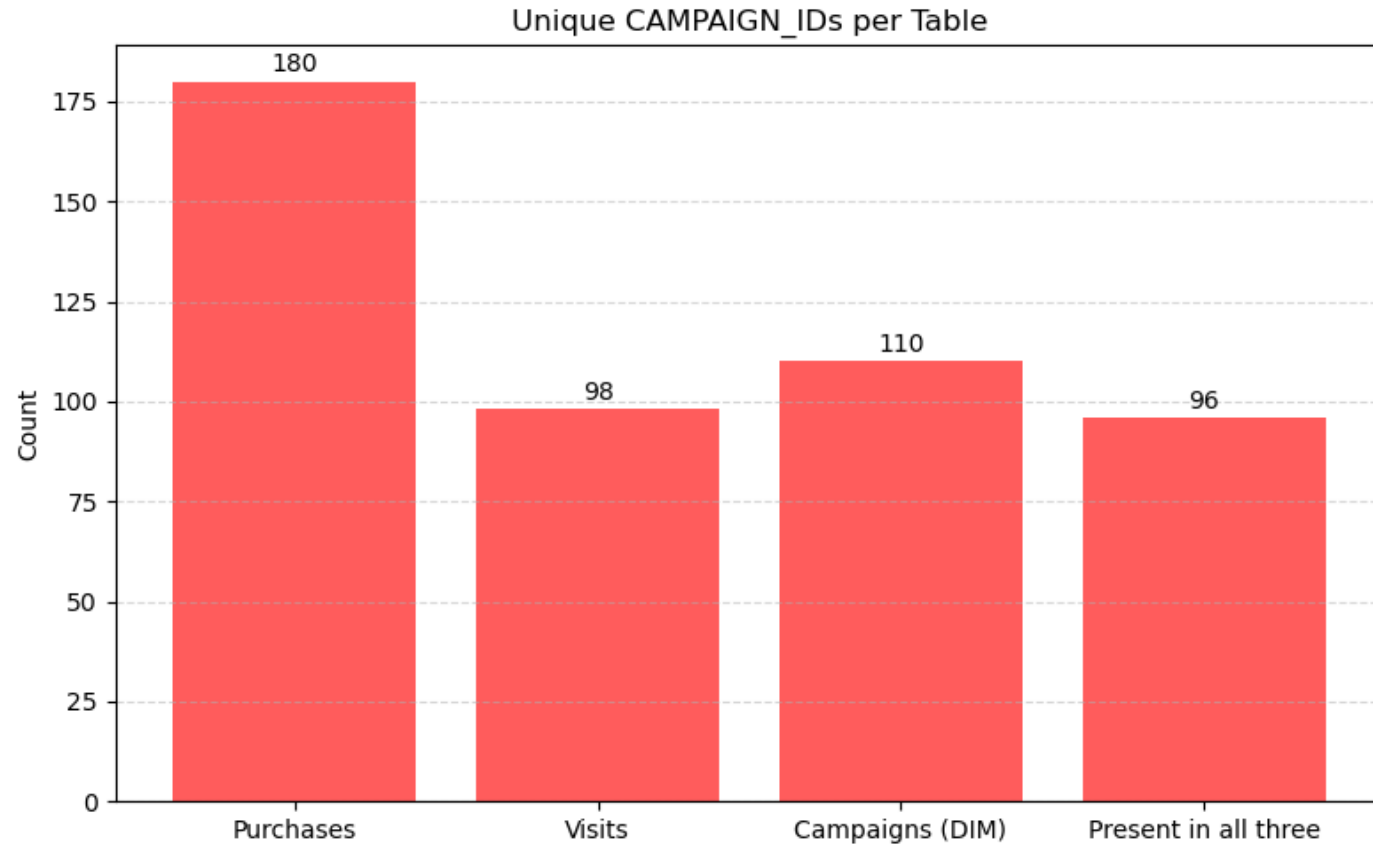




# 3. DATA GAPS

# 3. DATA GAP AND INCONSISTENCIES

## 3.1 CAMPAIGN ID

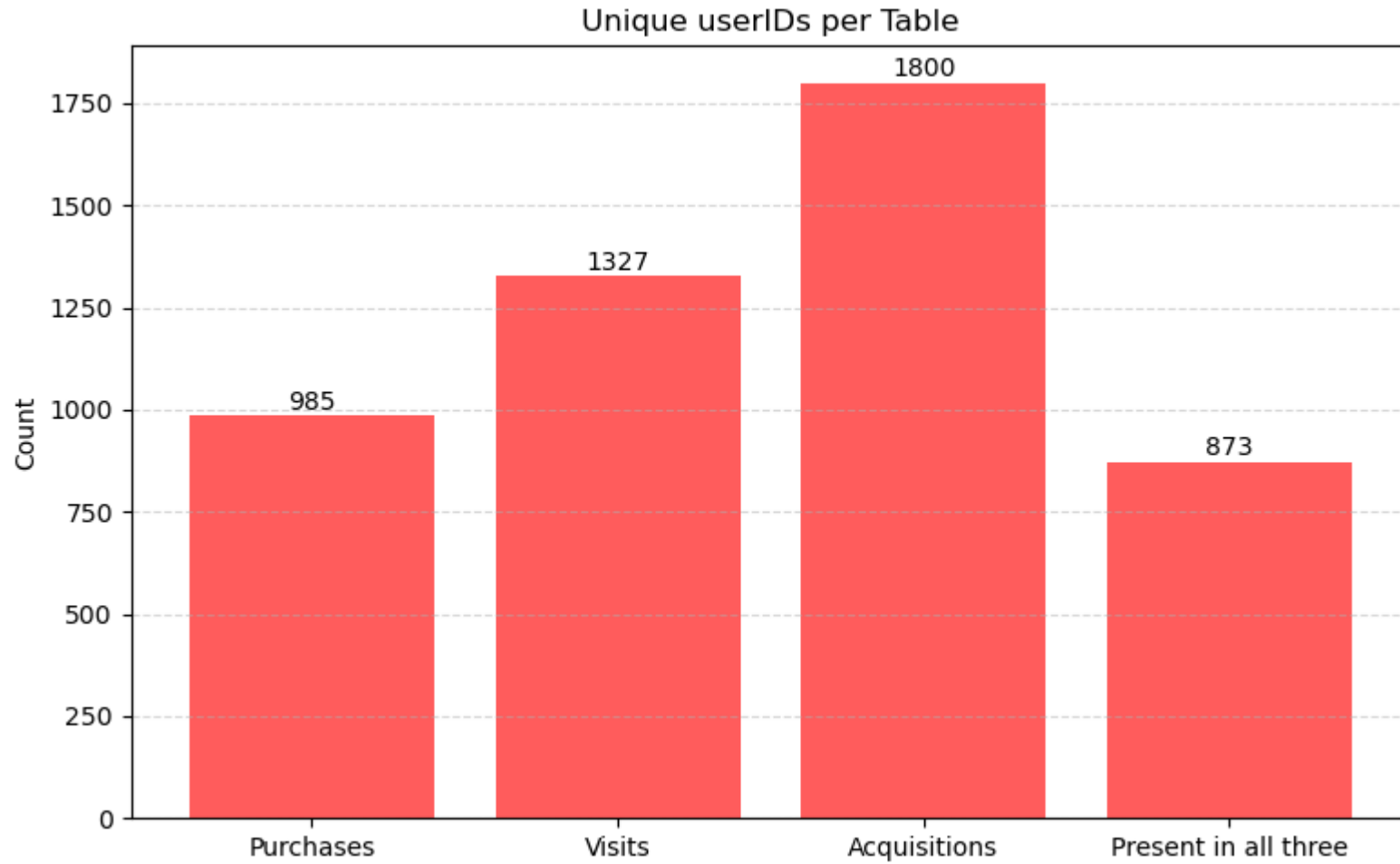


## INSIGHTS

Only 96 campaigns are available in the 3 datasets.

# 3. DATA GAP AND INCONSISTENCIES

## 3.2 USER ID



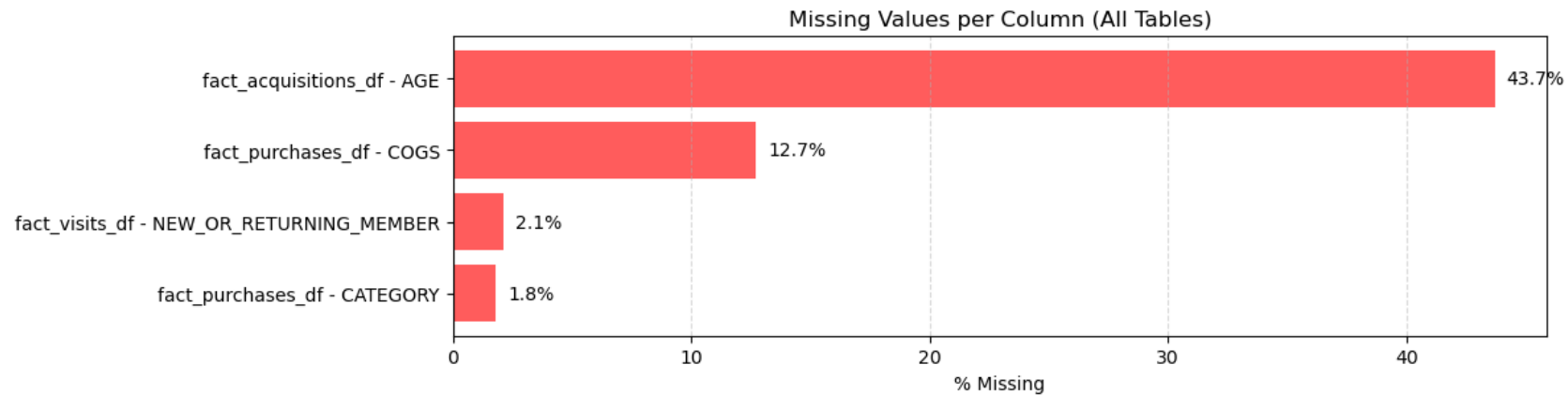
## INSIGHTS

Only 873 users are available in the 4 datasets.

*fabfitfun*

# 3. DATA GAP AND INCONSISTENCIES

## 3.3 AGE

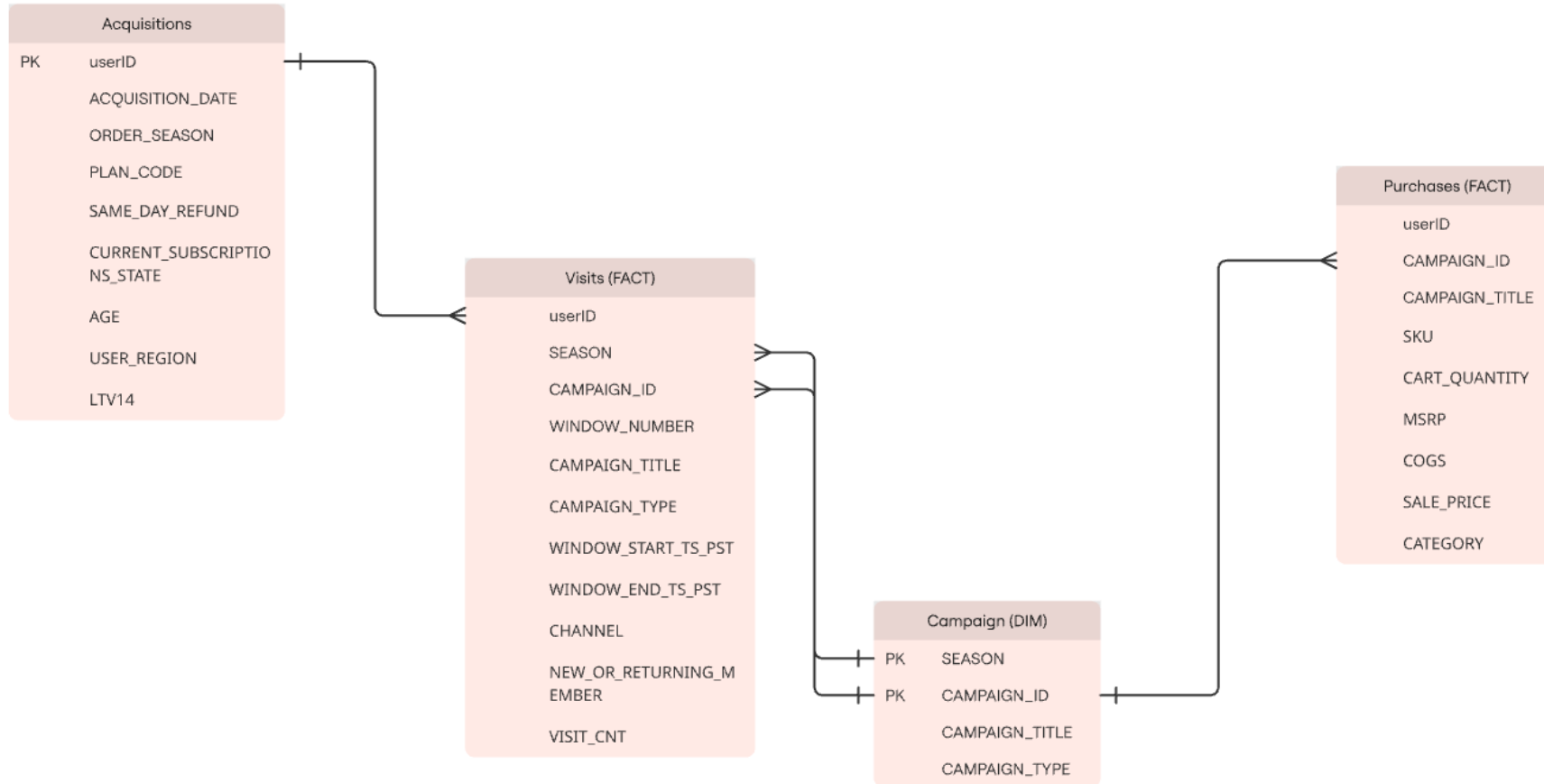


## INSIGHTS

Ages has the higher quantity of Nulls, preventing its usage for deeper insights about users

# 3. DATA GAP AND INCONSISTENCIES

## 3.4 MISSING PRIMARY KEY

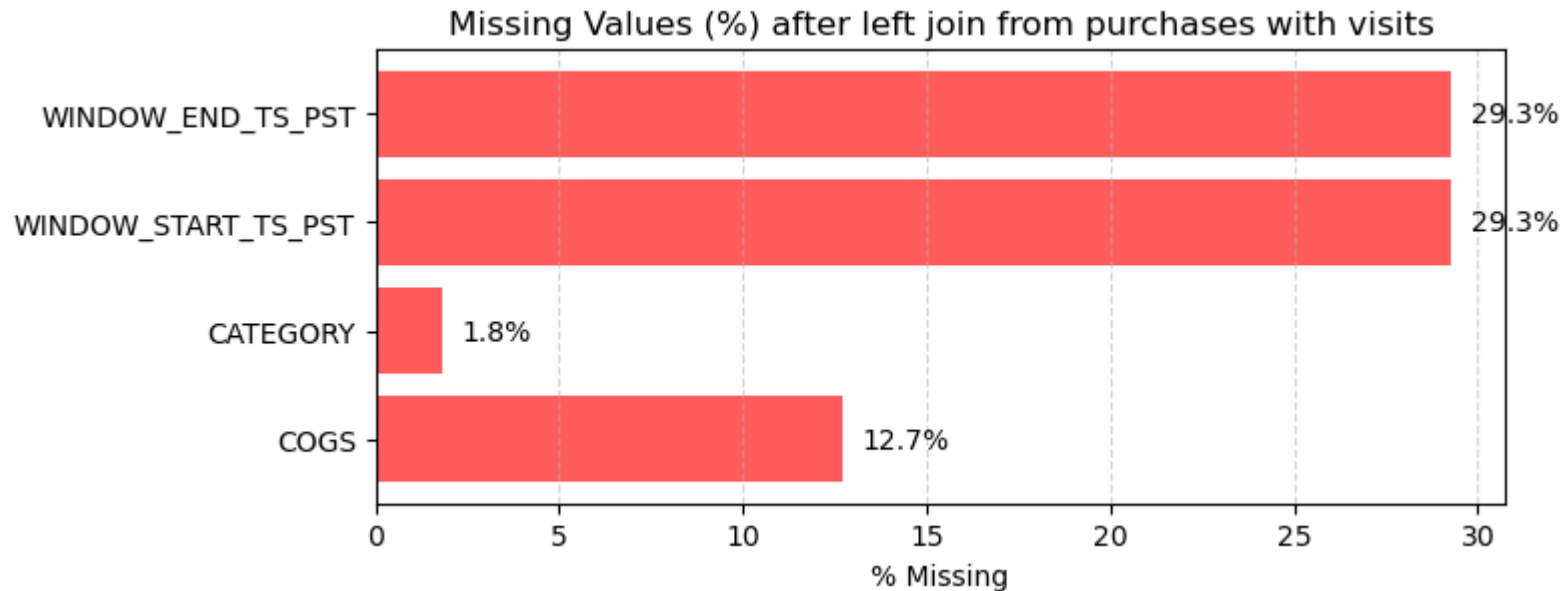


## INSIGHTS

- Relationship between Visits and Purchases table is N:N
- No PK in Visits and Purchases

# 3. DATA GAP AND INCONSISTENCIES

## 3.5 JOIN BETWEEN VISITS AND PURCHASES



## INSIGHTS

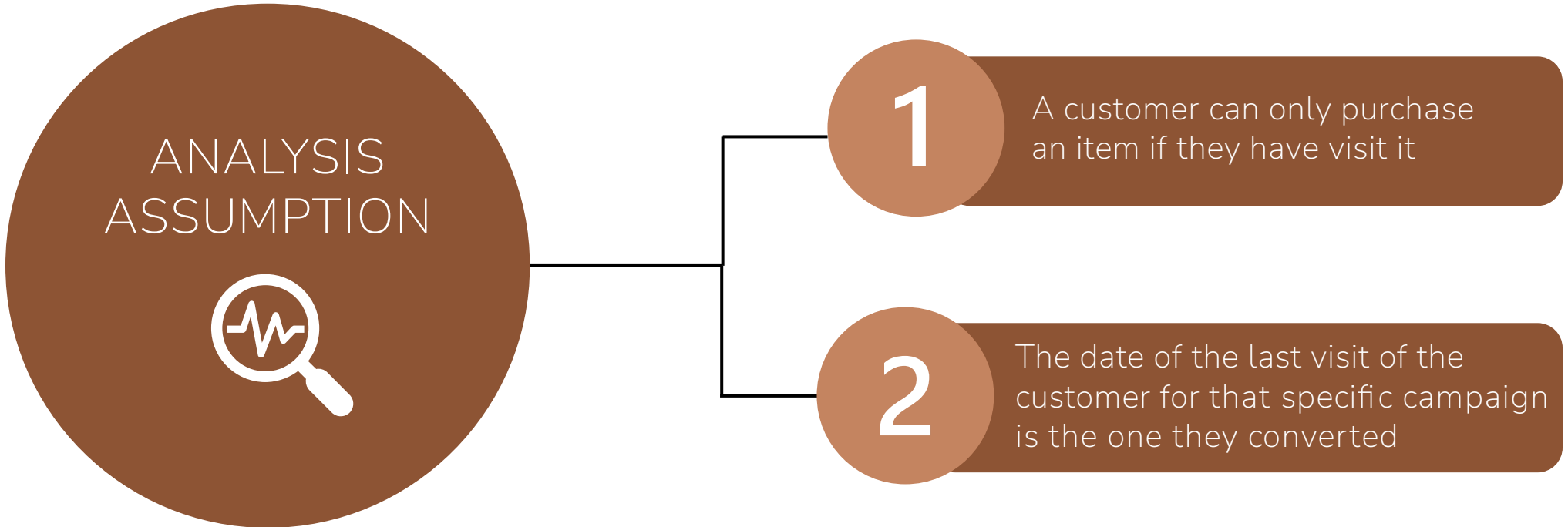
By joining visits and purchases using User ID and Campaign ID, we'll have a GAP of 29.3% of missing values for Window Start and Window End.

In addition, there's no window number in purchases.

# 4. DATA ANALYSIS

# 4. DATA ANALYSIS

## 4.1 VALUABLE USERS

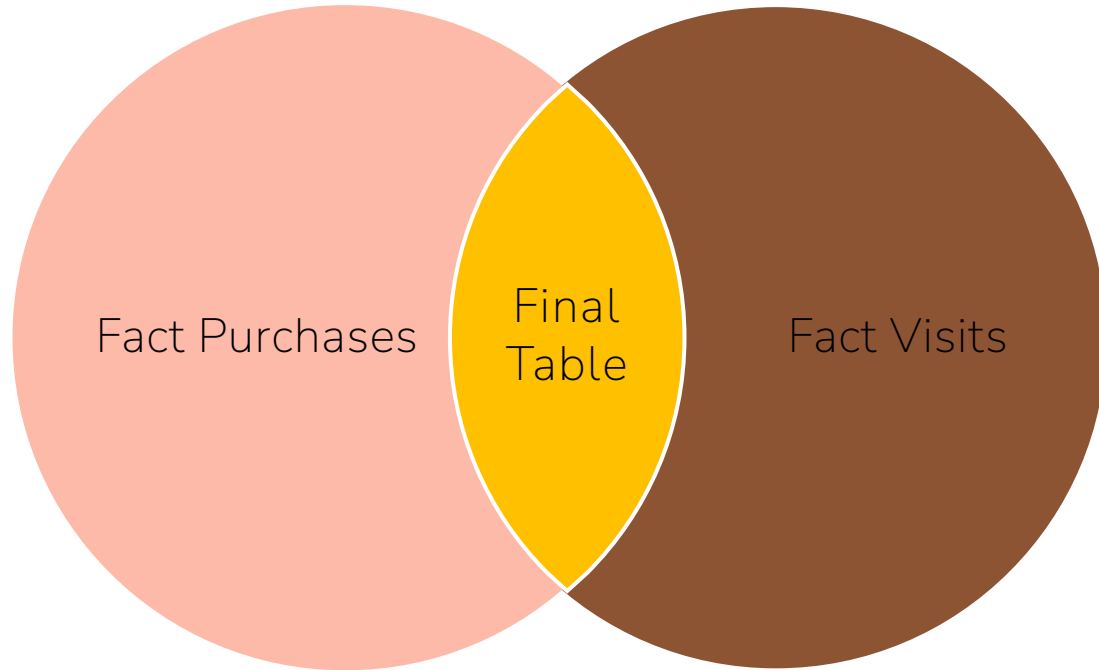


By performing this join, we kept 70.71% of the rows from Fact Purchases.



# 4. DATA ANALYSIS

## 4.1 VALUABLE USERS



- Sort values by *WINDOW\_END\_TS\_PST* and drop duplicates of *userID* and *CAMPAIGN\_ID*
- Inner join Fact Purchases with Fact Visits using of *userID* and *CAMPAIGN\_ID*
- Calculate *TOTAL\_REVENUE* by performing  $CART\_QUANTITY \times SALE\_PRICE$

By performing this join:

- Kept 70.71% of the rows from Fact Purchases
- Kept 77.77% of the users from Fact Purchases
- Kept 53.33% of the campaigns from Fact Purchases

# 4. DATA ANALYSIS

## 4.1 VALUABLE USERS

### Recency

How recently a customer made a purchase (in days).  
Used a *time window*<sup>1</sup> between:  
2023-08-06 and 2024-08-06  
per user.

### Frequency

How often or how many times  
the customer makes a purchase.  
Used a *time window*<sup>1</sup> between:  
2023-08-06 and 2024-08-06  
per user.

### Monetary

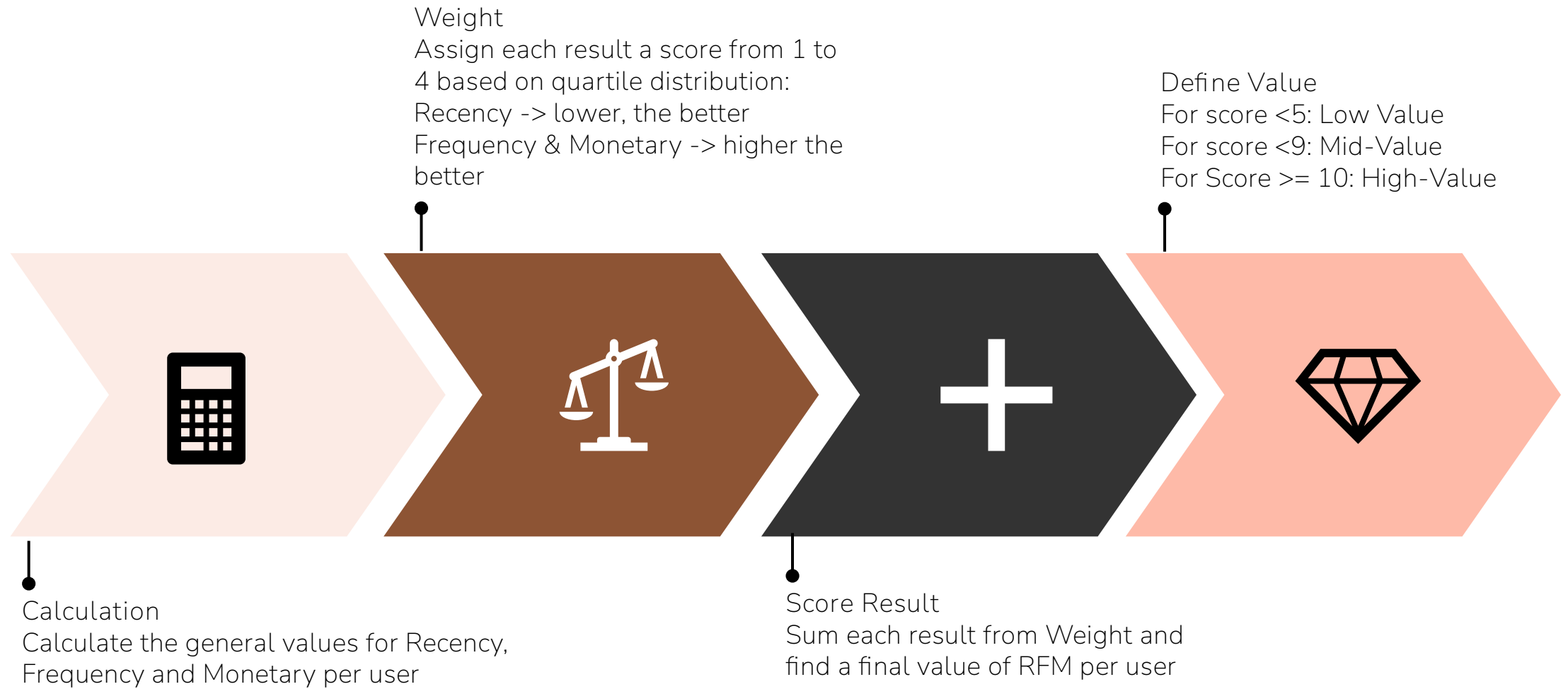
How much money the customer  
has spent.  
Used a *time window*<sup>1</sup> between:  
2023-08-06 and 2024-08-06  
per user.

*Time window*<sup>1</sup>:  
BETWEEN  
MAX(WINDOW\_END\_TS\_PST)  
AND  
MAX(WINDOW\_END\_TS\_PST) - 1 YEAR

*fabfitfun*

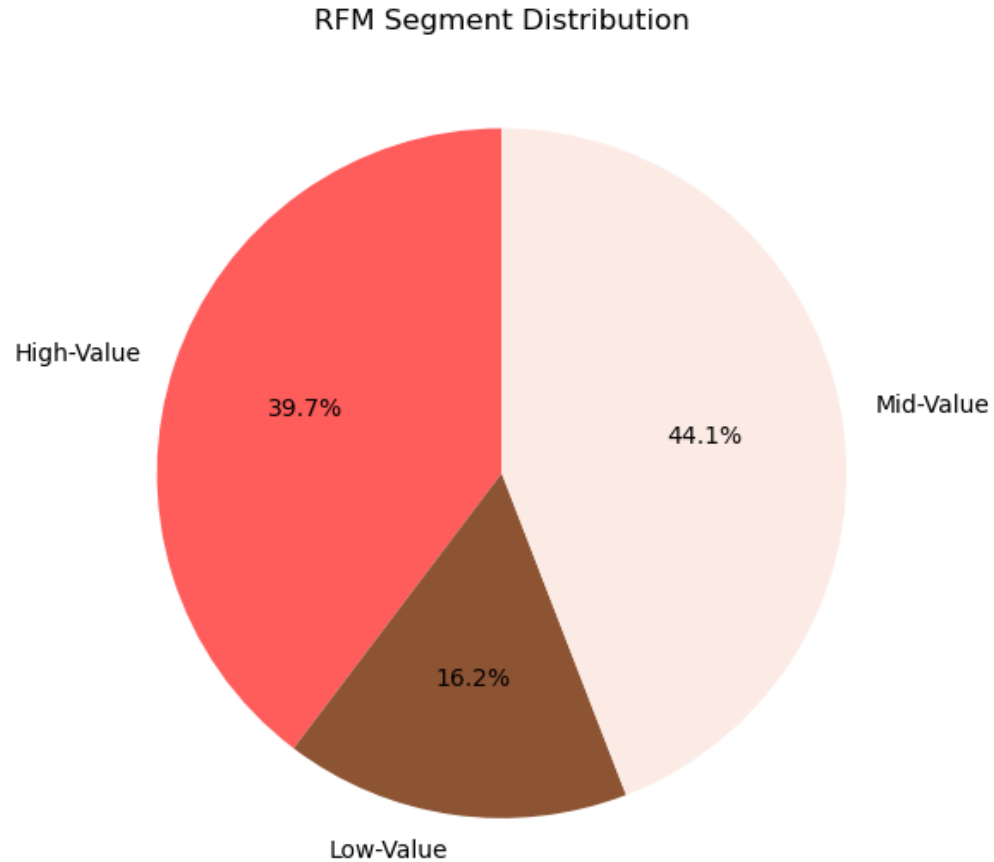
# 4. DATA ANALYSIS

## 4.1 VALUABLE USERS



# 4. DATA ANALYSIS

## 4.1 VALUEABLE USERS



Most valueble users represents 39.7%.

## INSIGHTS

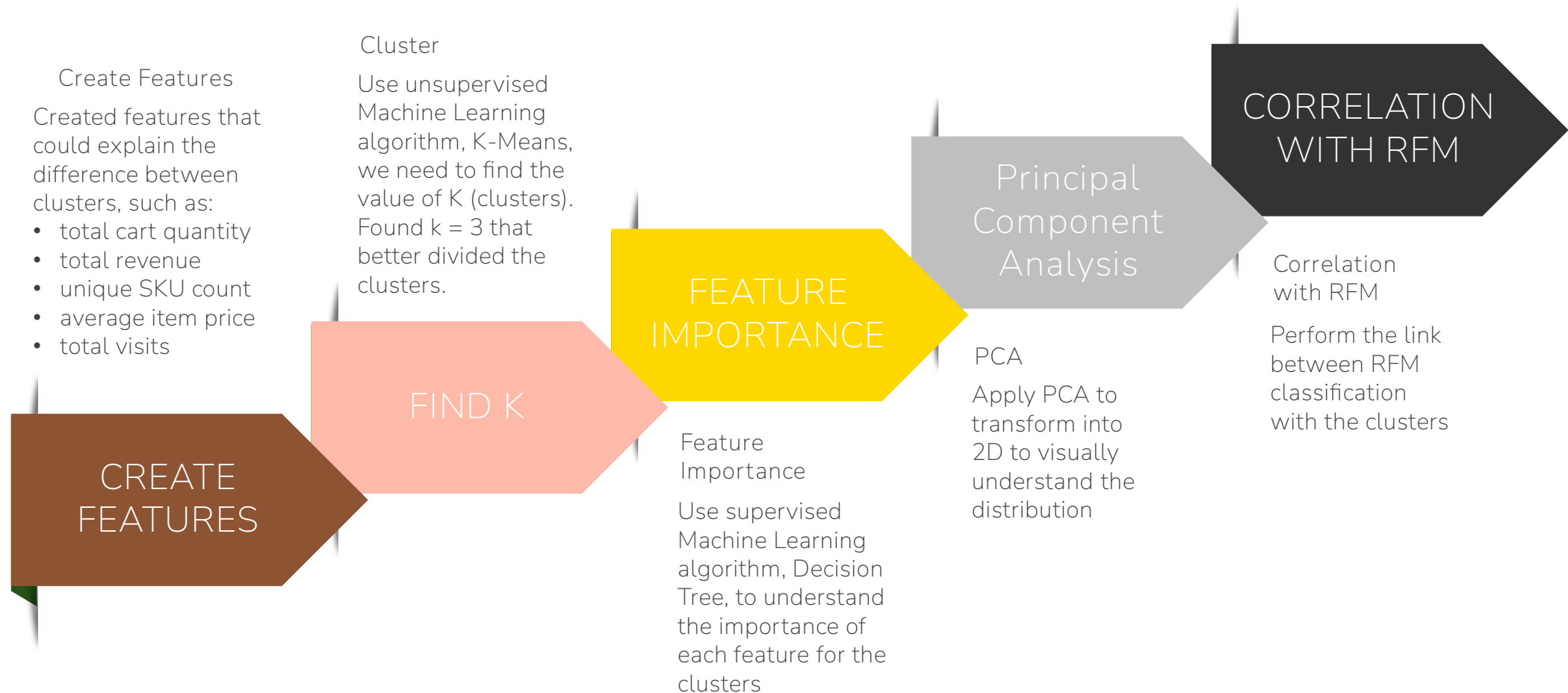
From 451 customers from 2023-08-06 to 2024-08-06:

- 39.7% have High Value
- 44.1% have Mid Value
- 16.2% have Low Value

*fabfitfun*

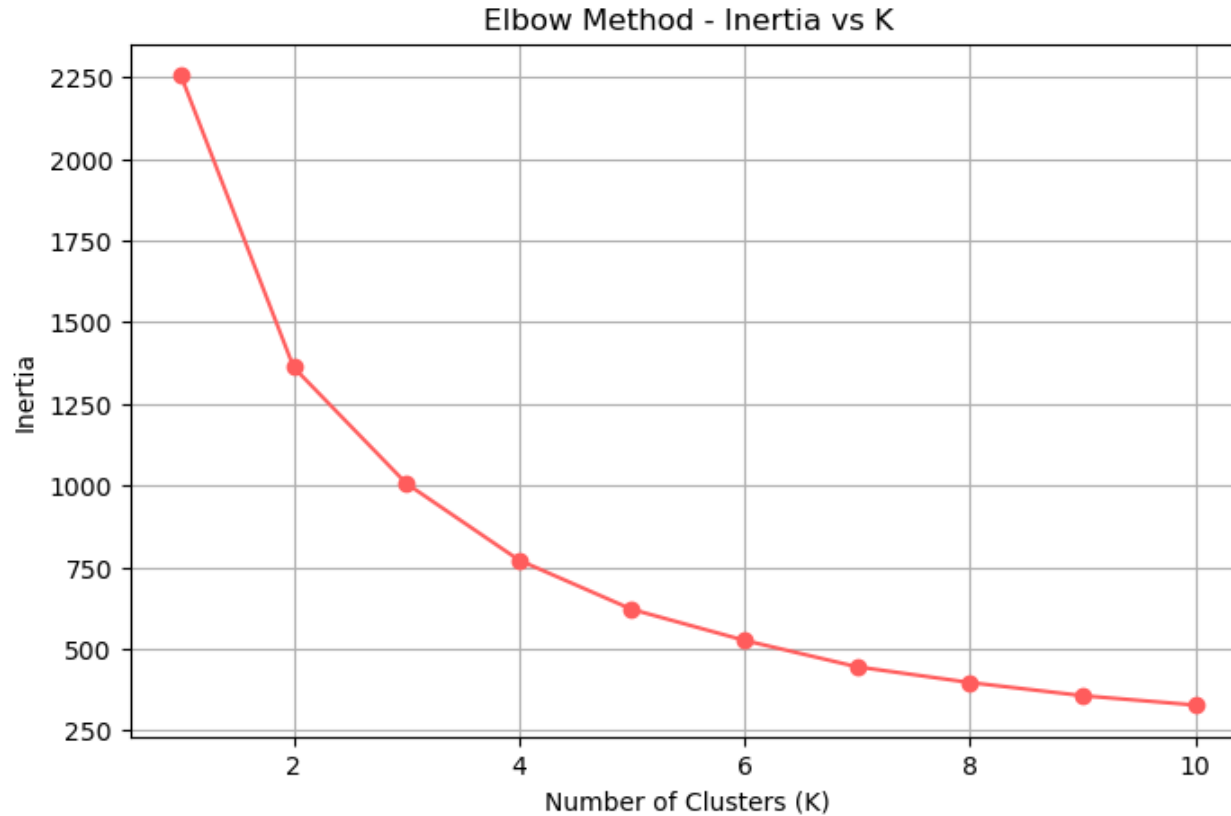
# 4. DATA ANALYSIS

## 4.2 DIFFERENCE BETWEEN VALUABLE USERS AND OTHER USERS



# 4. DATA ANALYSIS

## 4.2 DIFFERENCE BETWEEN VALUABLE USERS AND OTHER USERS

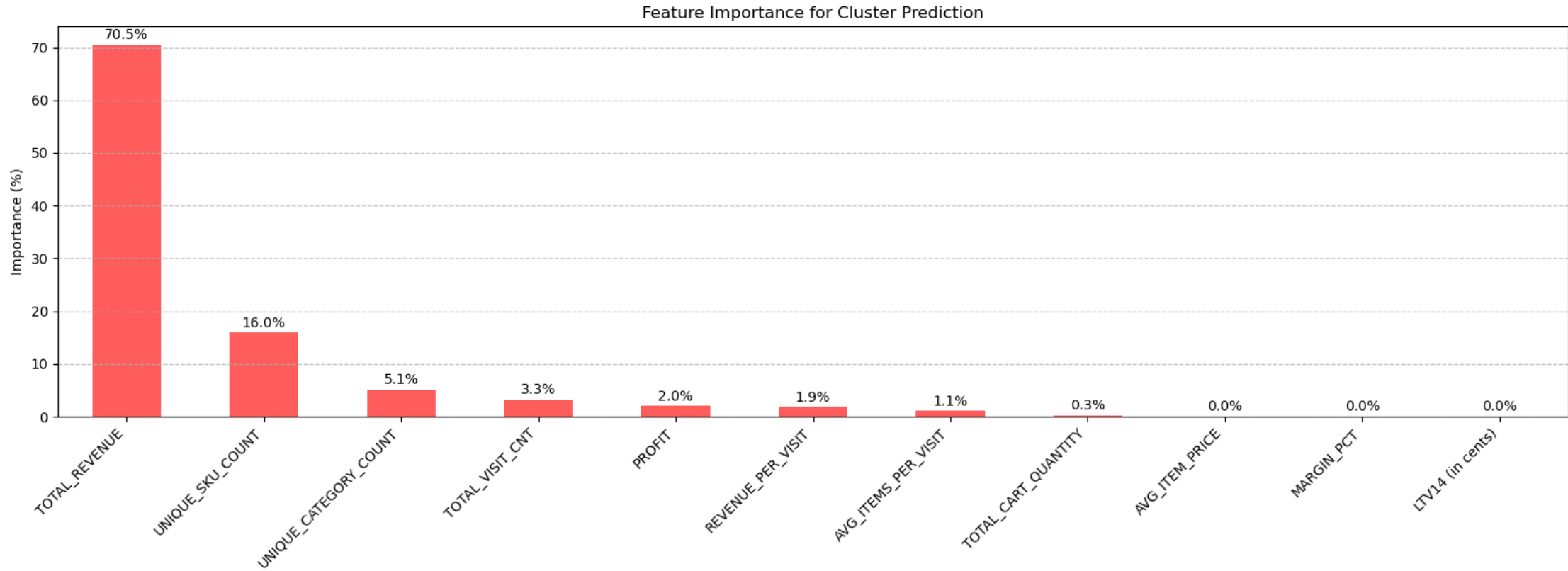


INSIGHTS

*fabfitfun*

# 4. DATA ANALYSIS

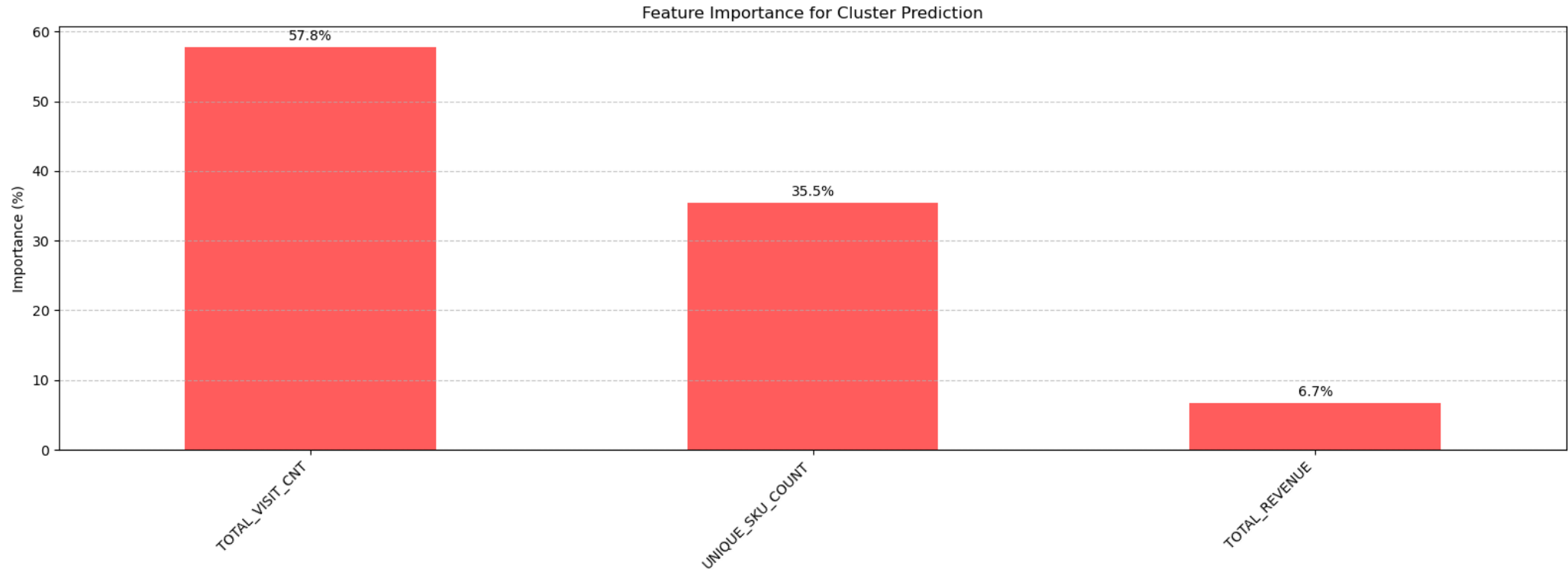
## 4.2 DIFFERENCE BETWEEN VALUABLE USERS AND OTHER USERS



Most of the features doesn't represent much of the differences between clusters. Let's create a thresh hold for over 3% to consider as an important feature

# 4. DATA ANALYSIS

## 4.2 DIFFERENCE BETWEEN VALUABLE USERS AND OTHER USERS

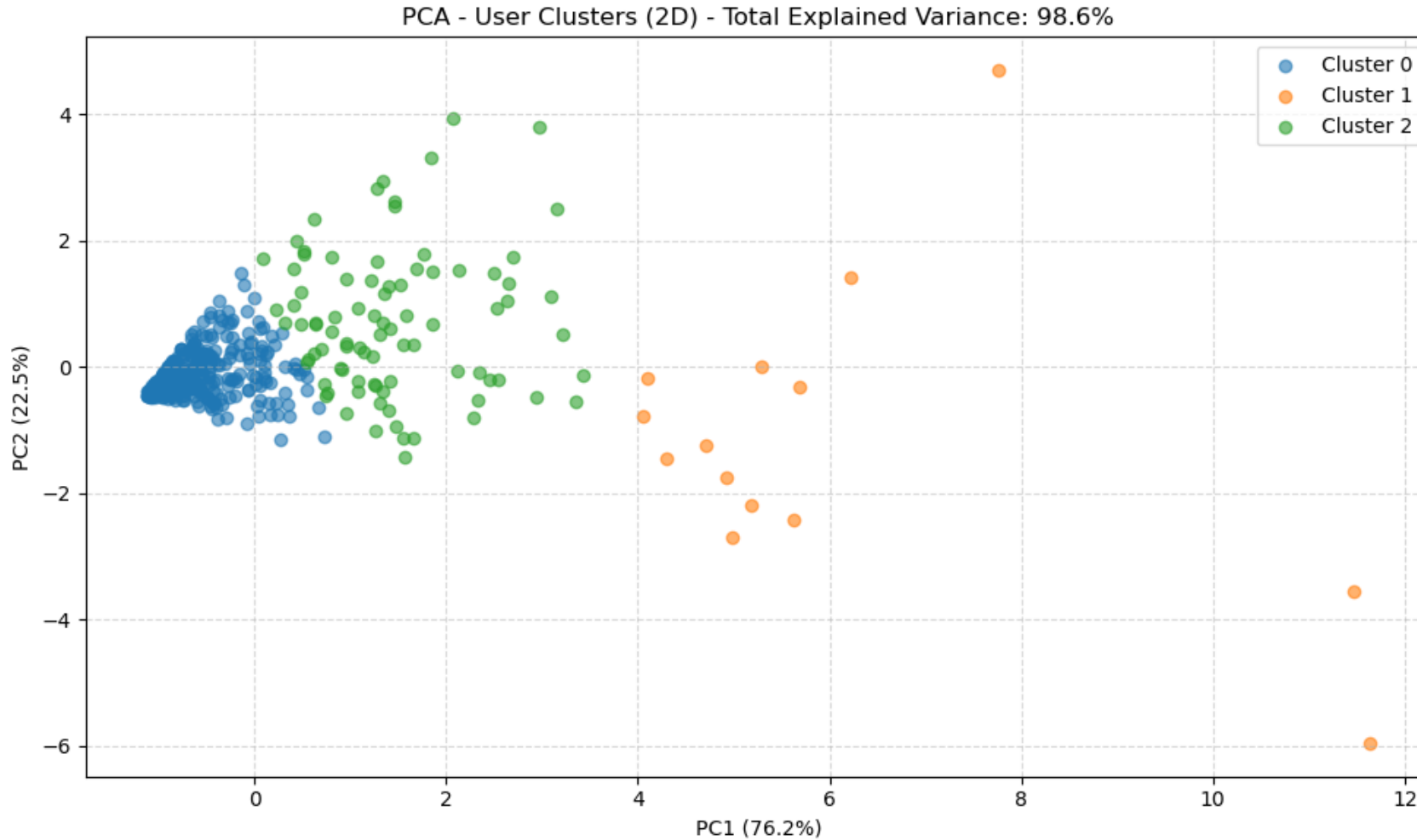


Recalculated the new feature importances for the algorithm.



# 4. DATA ANALYSIS

## 4.2 DIFFERENCE BETWEEN VALUABLE USERS AND OTHER USERS



## INSIGHTS

From 451 customers from 2023-08-06 to 2024-08-06:

- 77.61% are in cluster 0
- 03.10% are in cluster 1
- 19.29% are in cluster 2

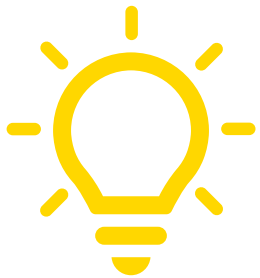
# 4. DATA ANALYSIS

## 4.2 DIFFERENCE BETWEEN VALUABLE USERS AND OTHER USERS

MEAN



CLUSTER	RFM SCORE	TOTAL REVENUE	UNIQUE SKU COUNT	TOTAL VISIT CNT	% Of High Value Users	% Of Mid Value Users	% Of Low Value Users
0	6.65	81.10	4.92	11.24	25.71%	20.86%	53.43%
1	11.29	1599.75	79.07	52.93	100%	0.00%	0.00%
2	10.67	398.96	22.56	44.21	86.21%	0.00%	13.79%



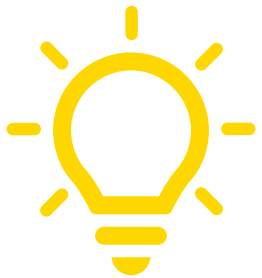
Most valuable users are in cluster 0 and 2. They differ from cart quantity, total revenue, unique skus and visits.

*fabfitfun*

# 4. DATA ANALYSIS

## 4.2 CHANNELS

CHANNEL	High-Value	Low-Value	Mid-Value
Facebook	158	3	18
Non-Attributed	1644	67	405
Rakuten	20	0	1
crm_email	550	9	75
crm_sms	143	1	46
mobile_android	166	3	29
mobile_ios	1339	13	198



The channels that best drives these members to sales are: Mobile IOS and E-mail

# 4. DATA ANALYSIS

## 4.3 STRONG CONVERSION RATE

user_region	users_campaign_visited	users_campaign_bought	conversion_rate
canada	536	150.0	27.99
military	3	0.0	0.0
non-continental us	22	7.0	31.82
other	44	5.0	11.36
us-midwest	1040	253.0	24.33



The locations with a stronger conversion rates are, respectively: Non-continental us (31.82%), Canada (27.99%), us-midwest (24.33%).

## INSIGHTS

The locations with a stronger conversion rates are, respectively:

- Non-continental us = 31.82%
- Canada = 27.99%
- Us-midwest = 24.33%

*fabfitfun*