



TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN  
TRUNG TÂM TIN HỌC

# ĐỒ ÁN TỐT NGHIỆP DATA SCIENCE PROJECT CUSTOMER SEGEMENTATION

Nguyễn Nhật Tố Trân vs Nguyễn Vũ Mai Phương



# TABLE OF CONTENTS

**01** **Business Understanding**

**02** **EDA**

**03** **RFM**

**04** **Manual Segmentation**

**05** **Unsupervised Modeling**

**06** **SOLUTIONS**

01

# **BUSINESS UNDERSTANDING**

# BUSINESS UNDERSTANDING



## Problem

- Một cửa hàng tạp hóa bán sản phẩm thiết yếu và khách hàng là người mua lẻ
- Cần tăng doanh thu, cải thiện mức độ hài lòng của khách hàng



## Solution

- Xây dựng hệ thống phân cụm khách hàng
- Đề xuất phương án tiếp cận phù hợp dựa trên đặc điểm của từng nhóm khách hàng





# PROCESS

1

## EDA

Tiền xử lý dữ liệu  
và khám phá insights

2

## RFM

Phân nhóm khách hàng dựa trên  
thông tin mua hàng trước đó

3

## SOLUTIONS

Đề xuất giải pháp tiếp  
cận khách hàng phù hợp

### MANUAL SEGMENTATION

### UNSUPERVISED MODELING

KMeans (PySpark)

GMM

DBSCAN

Hierarchical  
Clustering

02

**EDA**

```
<class 'pandas.core.frame.DataFrame'>
Index: 38521 entries, 0 to 38764
Data columns (total 8 columns):
#   Column                Non-Null Count  Dtype
---  -
0   productId              38521 non-null  int64
1   productName            38521 non-null  object
2   price                  38521 non-null  float64
3   Category               38521 non-null  object
4   Member_number          38521 non-null  int64
5   Date                   38521 non-null  datetime64[ns]
6   items                  38521 non-null  int64
7   purchase_amount        38521 non-null  float64
dtypes: datetime64[ns](1), float64(2), int64(3), object(2)
memory usage: 2.6+ MB
```

Transactions timeframe from 2014-01-01 00:00:00 to 2015-12-30 00:00:00  
 0 transactions don't have a member number  
 3,898 unique Member\_number

## Data Pre-processing

- 38.521 entries (có duplicates, không có missing values, chuyển data type của “Date”)
- 8 features (“purchase\_amount” được tạo mới)
- Dữ liệu được ghi nhận trong 2 năm
- 3898 khách hàng

```
df[num_cols].describe()
```

	price	items	purchase_amount
count	38521.000000	38521.000000	38521.000000
mean	4.305367	1.996729	8.590893
std	4.320088	0.817539	9.952770
min	0.300000	1.000000	0.300000
25%	1.250000	1.000000	2.400000
50%	2.500000	2.000000	5.200000
75%	6.100000	3.000000	10.500000
max	28.500000	3.000000	85.500000

```
df[cat_cols].describe()
```

	productName	Category
count	38521	38521
unique	167	11
top	whole milk	Fresh Food
freq	2455	11426

## Data Pre-processing

- Phần lớn sản phẩm có giá thấp, và một số ít sản phẩm có giá rất cao
- Số lượng mỗi sản phẩm mỗi khách hàng mua là từ 1-3
- 167 sản phẩm, và “whole milk” được khi nhận nhiều nhất
- 11 ngành hàng, phổ biến nhất là Fresh Food



```
# Total purchases
total_purchases = df['items'].sum()
print(total_purchases)
```

76916

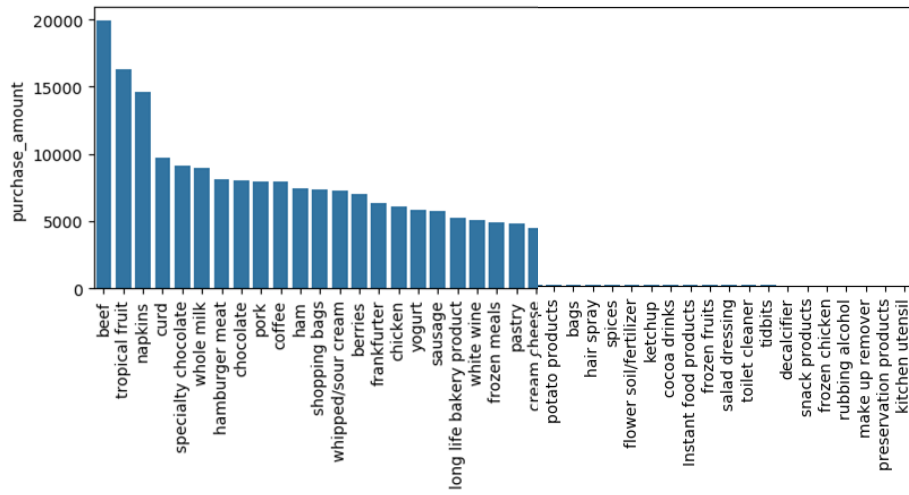
```
# Total revenue
total_revenue = (df['items'] * df['price']).sum()
print(total_revenue)
```

330929.800000000005

## Insights

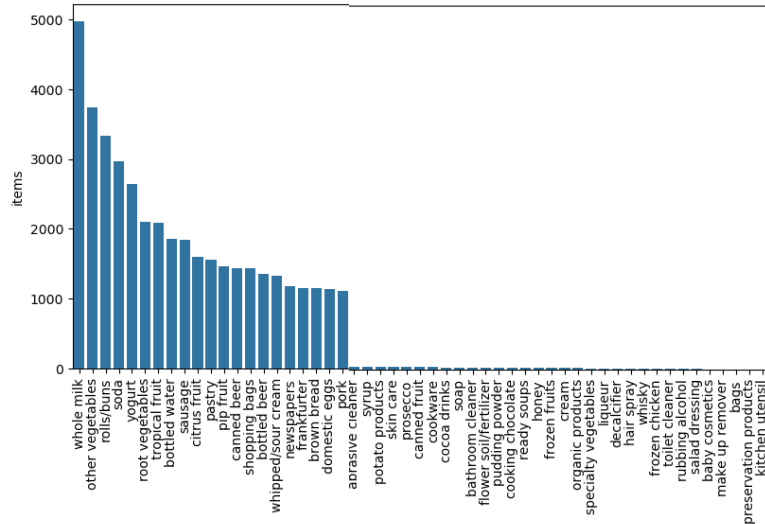
Trong vòng 2 năm

- Tổng sản phẩm được bán: 77.000
- Tổng doanh thu : 331.000\$



## Insights

- Sản phẩm mang lại doanh thu nhiều nhất : Thịt bò, trái cây nhiệt đới, khăn giấy, phô mai tươi, sô cô la đặc sản
- Sản phẩm bán chạy nhất : Sữa tươi nguyên chất, rau củ khác, bánh mì cuộn, soda, sữa chua
- Sản phẩm bán chậm nhất : Gà đông lạnh, cồn sát trùng, nước tẩy trang, sản phẩm bảo quản, dụng cụ nhà bếp



items purchase\_amount percent\_amount

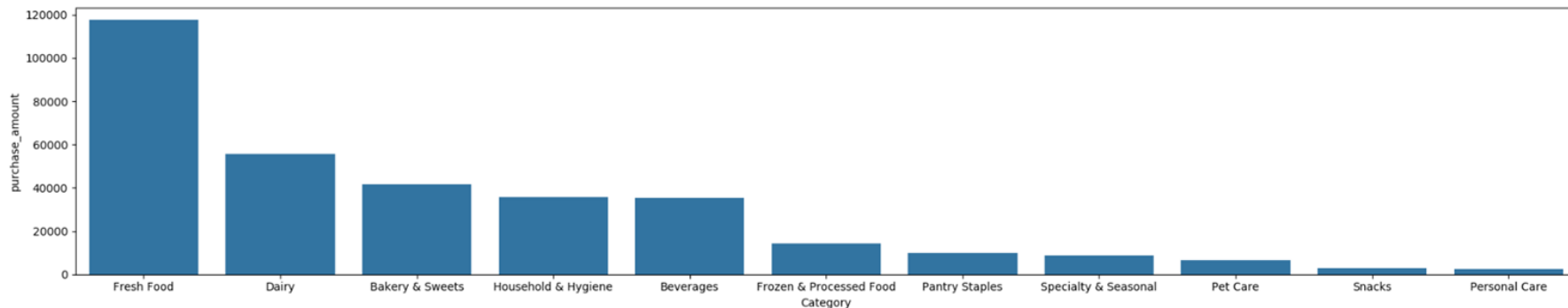
Category

Fresh Food	22786	117566.20	35.526024
Dairy	14968	55748.24	16.845941
Bakery & Sweets	11796	41780.78	12.625270
Household & Hygiene	4339	35555.16	10.744019
Beverages	11963	35391.10	10.694443
Frozen & Processed Food	3689	14371.92	4.342891
Pantry Staples	3134	9857.45	2.978713
Specialty & Seasonal	2360	8864.30	2.678604
Pet Care	641	6425.80	1.941741
Snacks	928	2727.40	0.824163
Personal Care	312	2641.45	0.798190

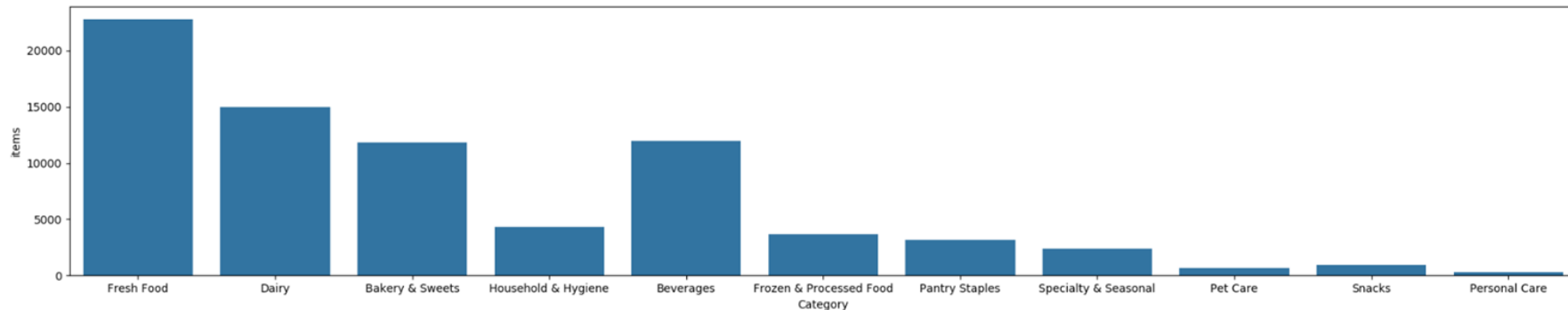
## Insights

- Fresh Food đóng góp 1/3 doanh thu
- Ngành hàng bán chạy nhất (60% doanh thu) : Thực phẩm tươi sống, Sản phẩm từ sữa, Bánh mì & đồ ngọt
- Sản phẩm chăm sóc thú cưng, đồ ăn vặt, chăm sóc cá nhân là các ngành hàng bán chậm nhất của cửa hàng

Biểu đồ doanh thu từng ngành hàng



Biểu đồ sản lượng mua từng ngành hàng



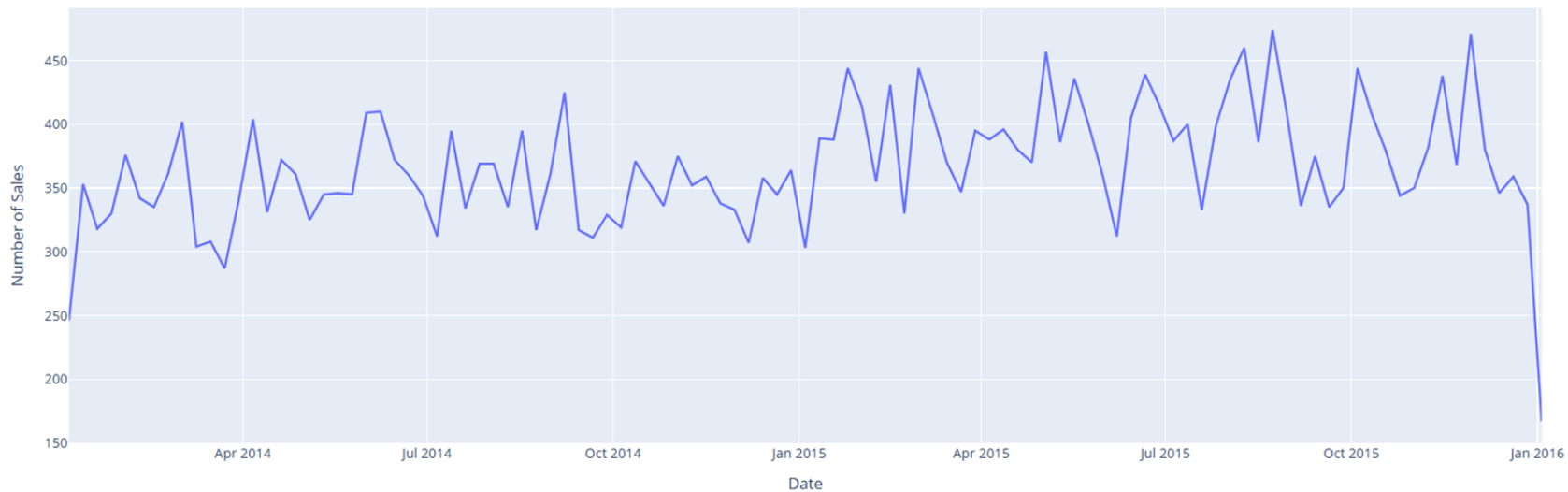
	items	purchase_amount
Member_number		
2193	63	361.45
1793	52	345.10
3289	63	334.15
2433	57	316.81
2743	41	312.46
...	...	...
1560	5	1.90
1221	2	1.70
4029	2	1.60
1250	2	1.30
4565	1	1.10

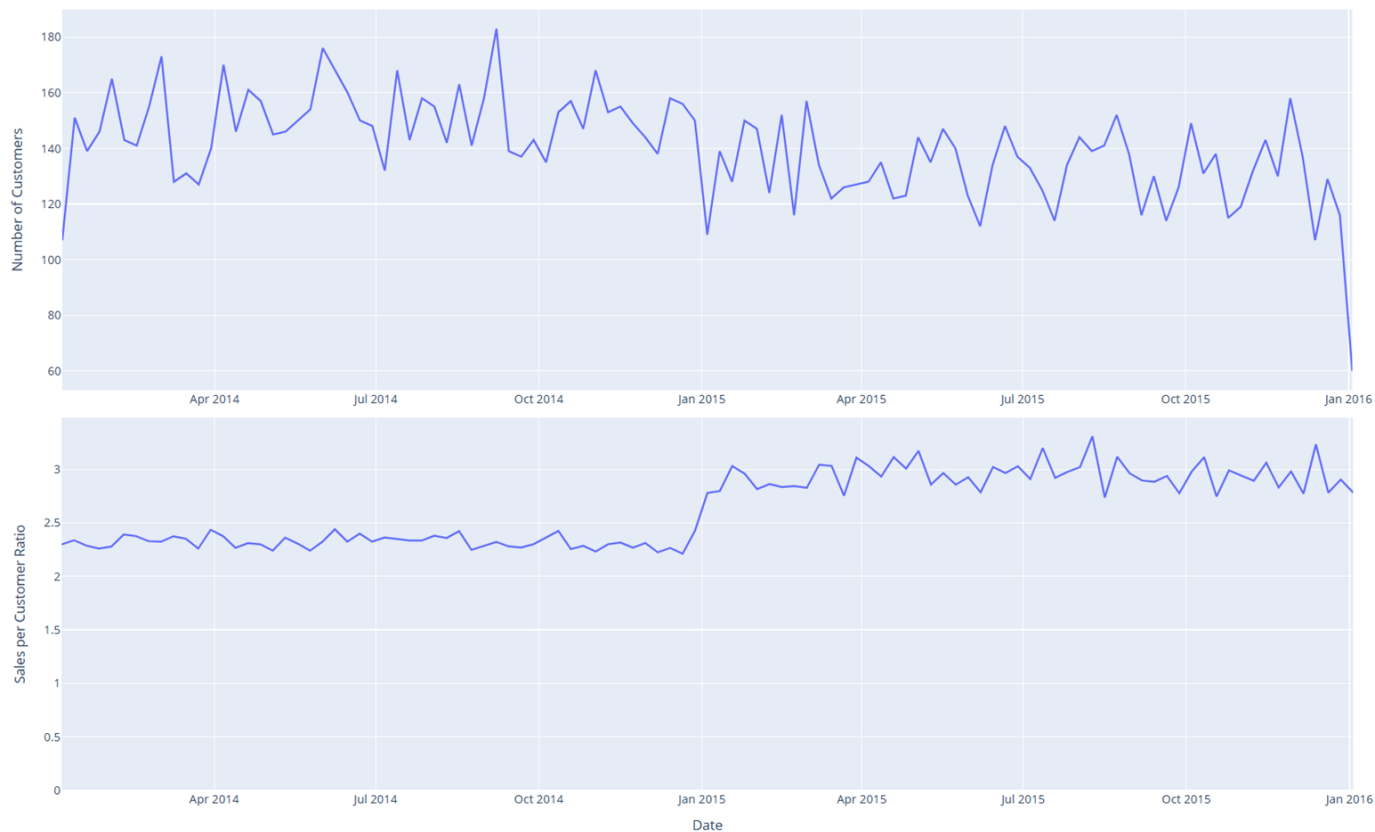
Danh sách sản lượng và doanh thu  
theo từng khách hàng

## Insights

- Khách hàng chi tiêu nhiều nhất với **63** sản phẩm trong 2 năm tương ứng tổng số tiền là **361\$**
- Xác định được khách hàng giá trị cao

Biểu đồ doanh thu hàng tuần





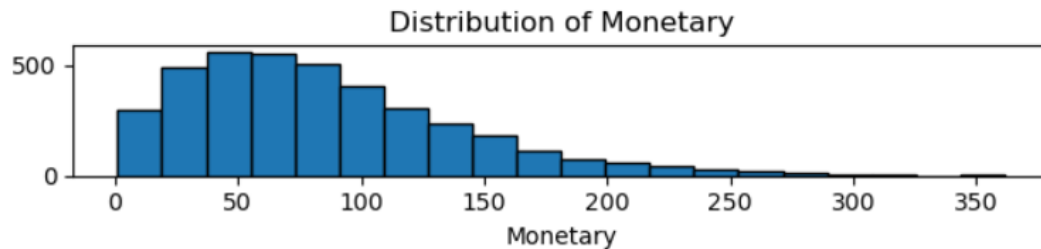
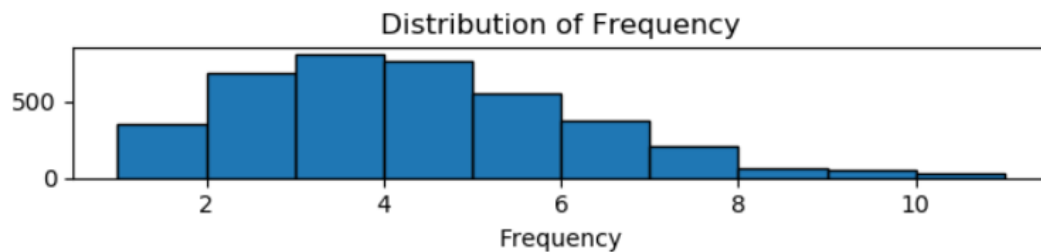
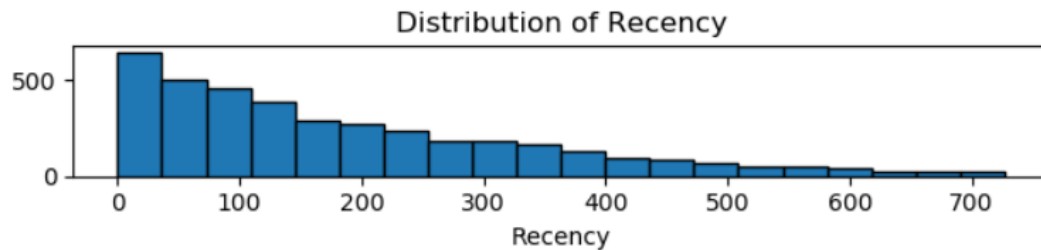
**Biểu đồ số lượng và doanh thu mỗi khách hàng hàng tuần**

03

**RFM**



# RFM

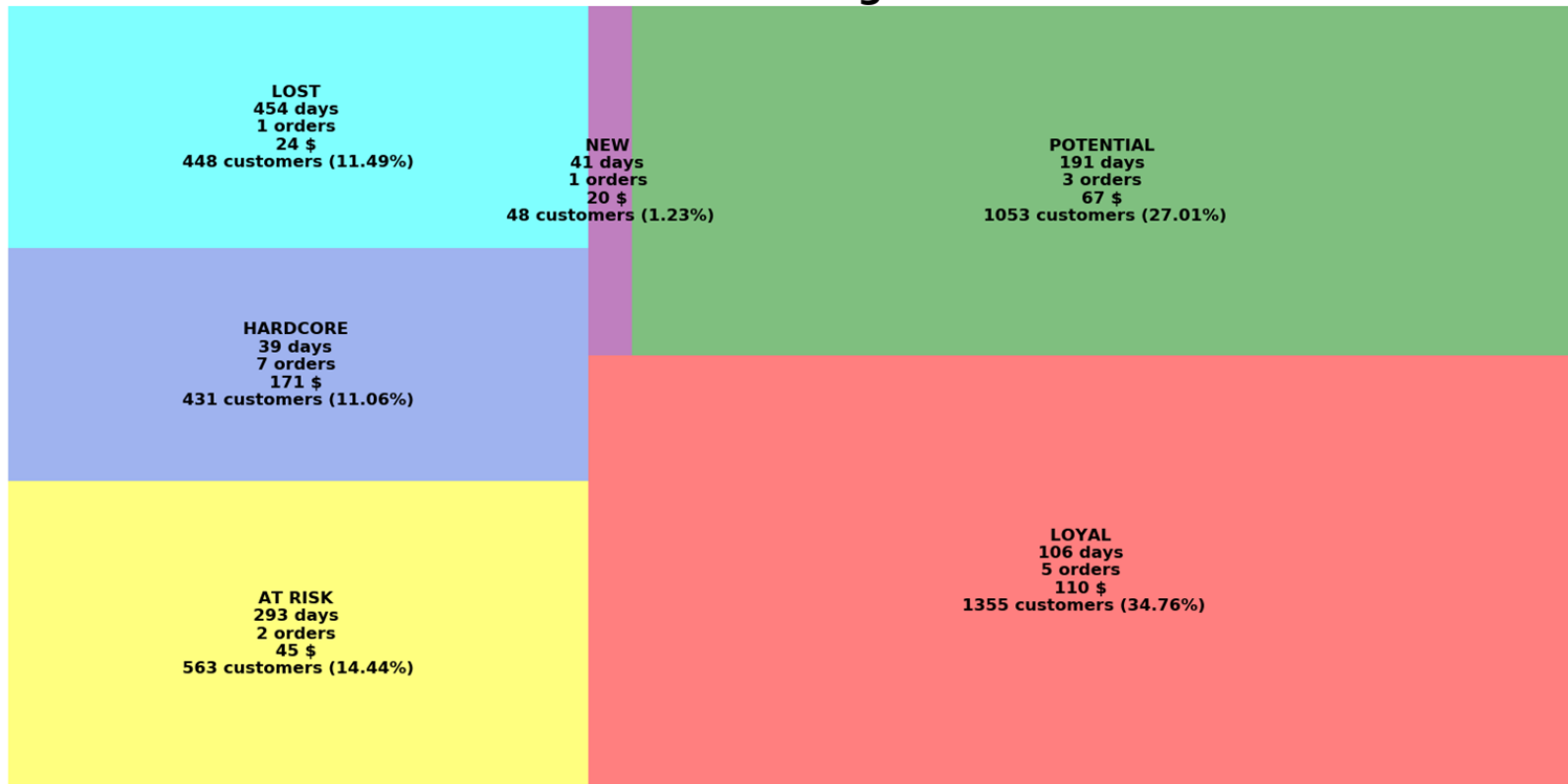


04

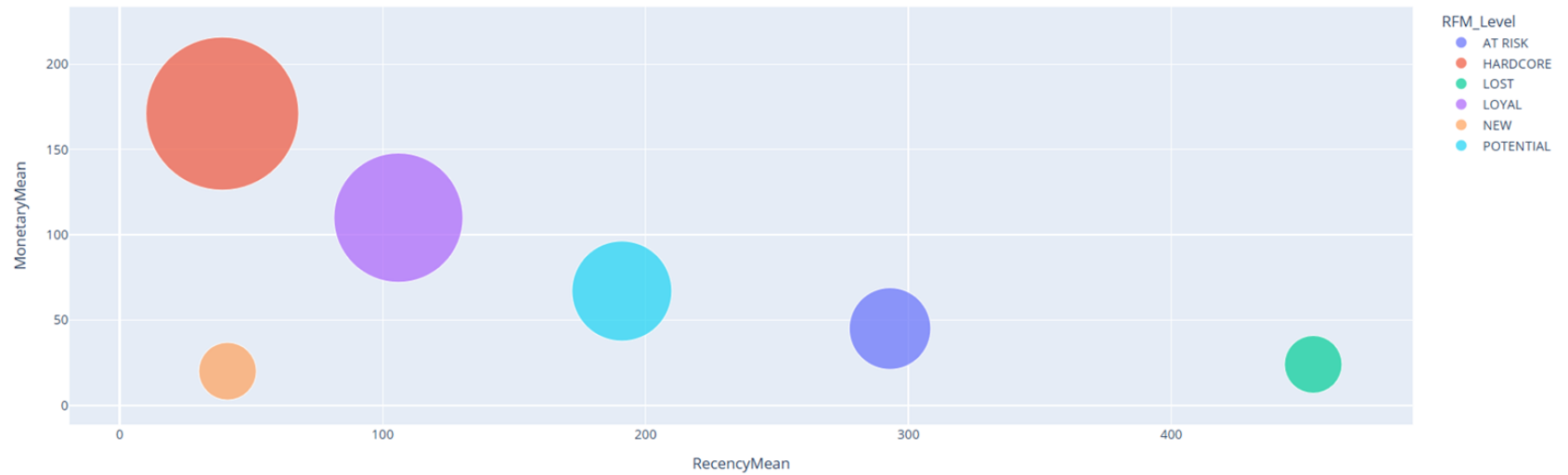
## **MANUAL SEGMENTATION**

# MANUAL SEGMENTATION

## Customers Segments



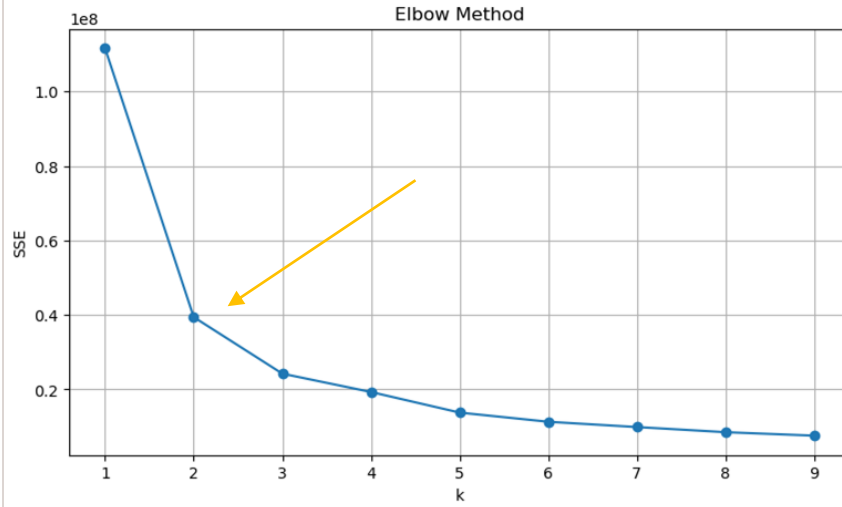
# MANUAL SEGMENTATION



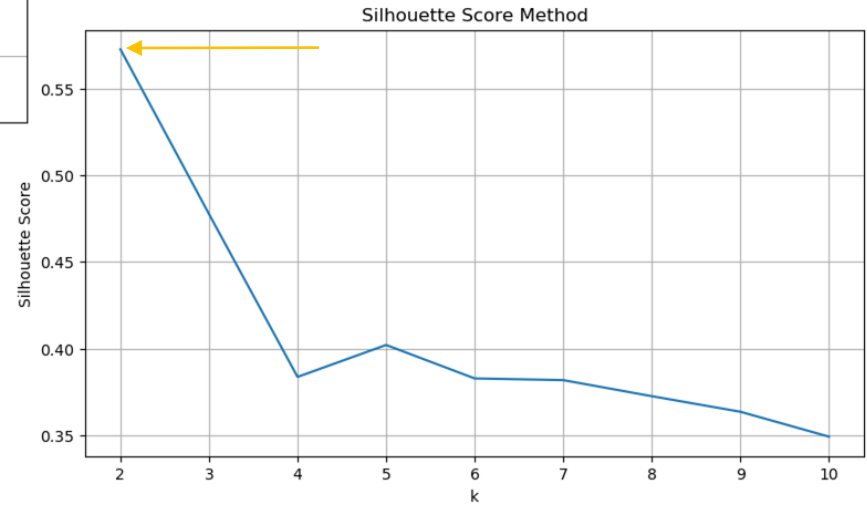
05

# **UNSUPERVISED MODELING**

# KMEANS



$k=2$



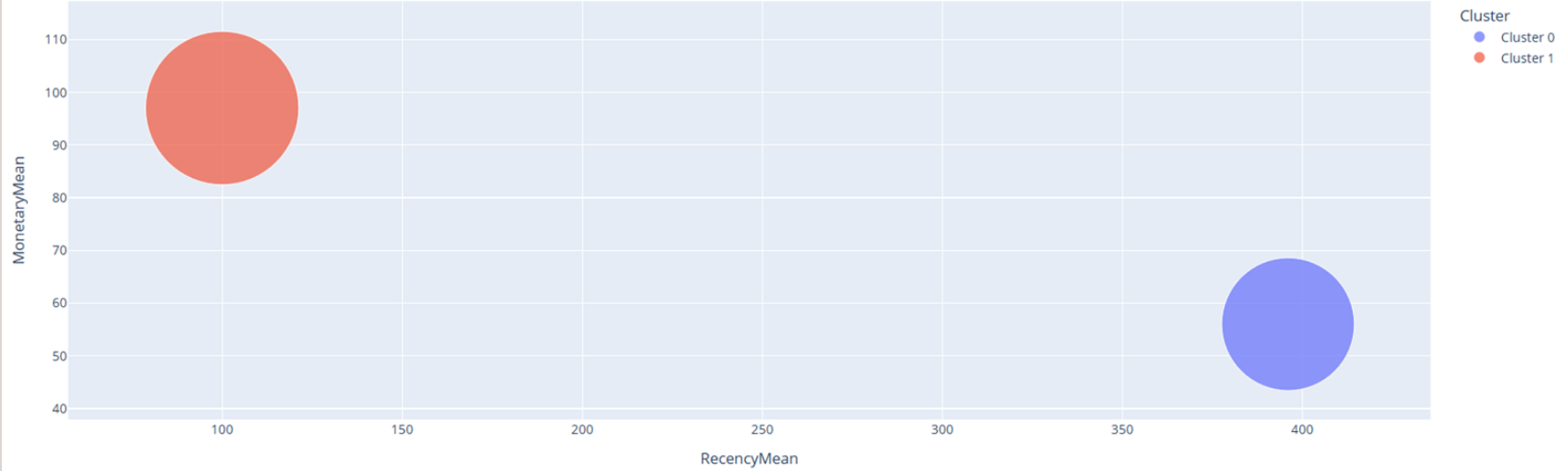
# KMEANS

## Customers Segments

**Cluster 1**  
**100 days**  
**4 orders**  
**97 \$**  
**2745 customers (70.42%)**

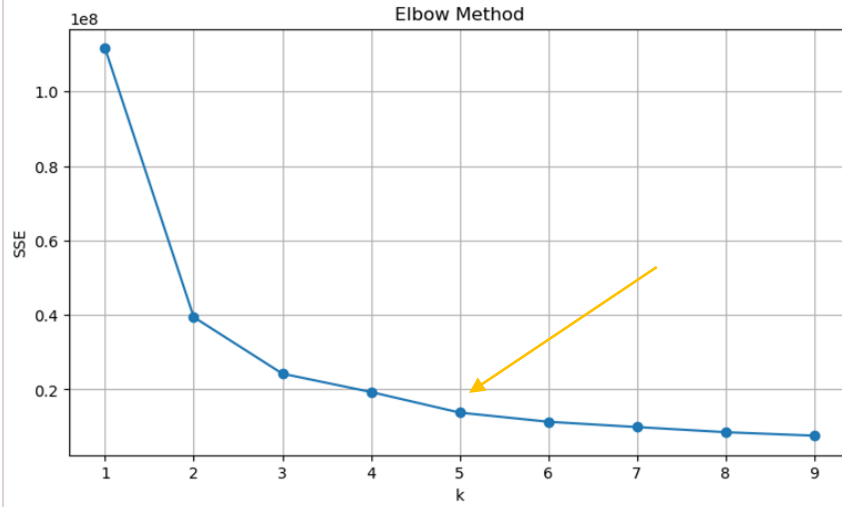
**Cluster 0**  
**396 days**  
**3 orders**  
**56 \$**  
**1153 customers (29.58%)**

# KMEANS

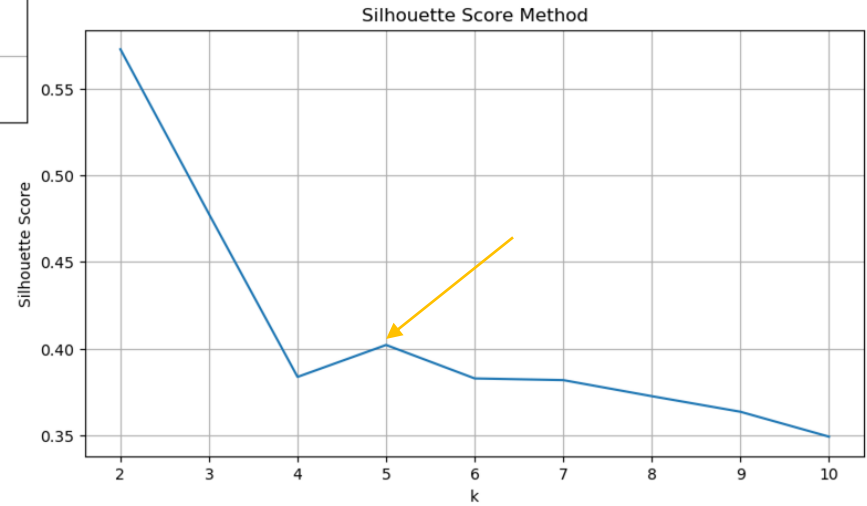




# KMEANS



k=5



# KMEANS

## Customers Segments

LOST

Cluster 1  
560 days  
2 orders  
37 \$  
301 customers (7.72%)

POTENTIAL

Cluster 0  
192 days  
4 orders  
78 \$  
1076 customers (27.6%)

AT-RISK

Cluster 4  
352 days  
3 orders  
61 \$  
715 customers (18.34%)

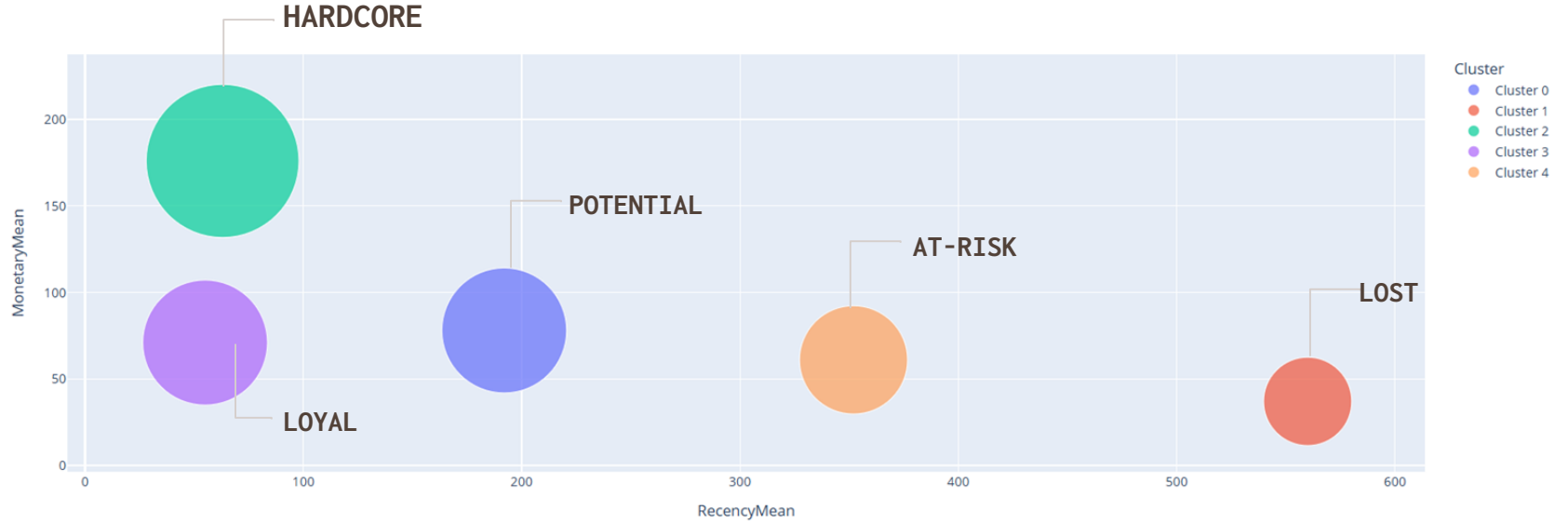
LOYAL

Cluster 3  
55 days  
4 orders  
71 \$  
1205 customers (30.91%)

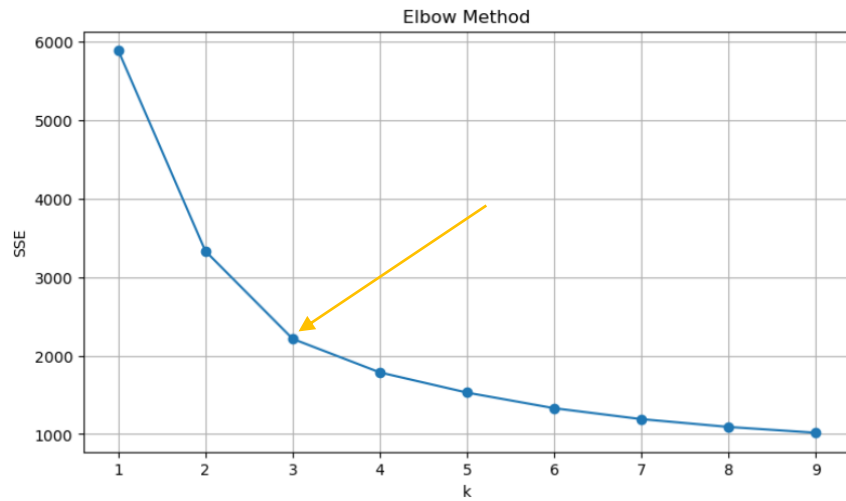
HARDCORE

Cluster 2  
63 days  
6 orders  
176 \$  
601 customers (15.42%)

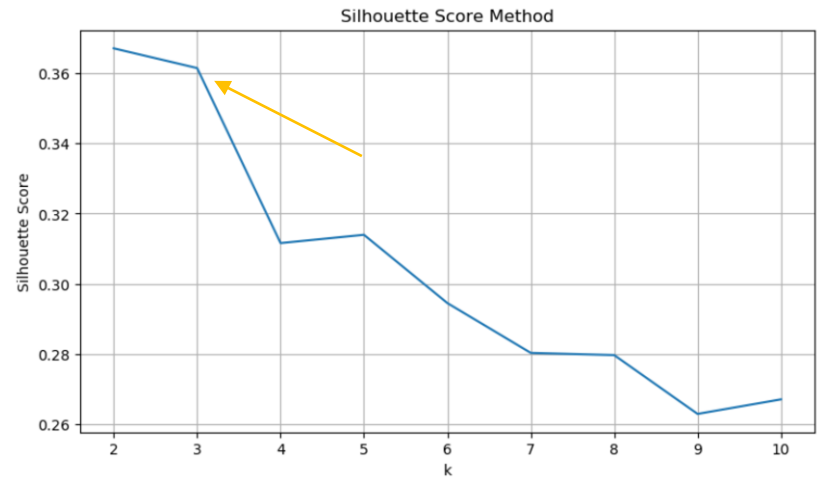
# KMEANS



# KMEANS (SCALE DATA)

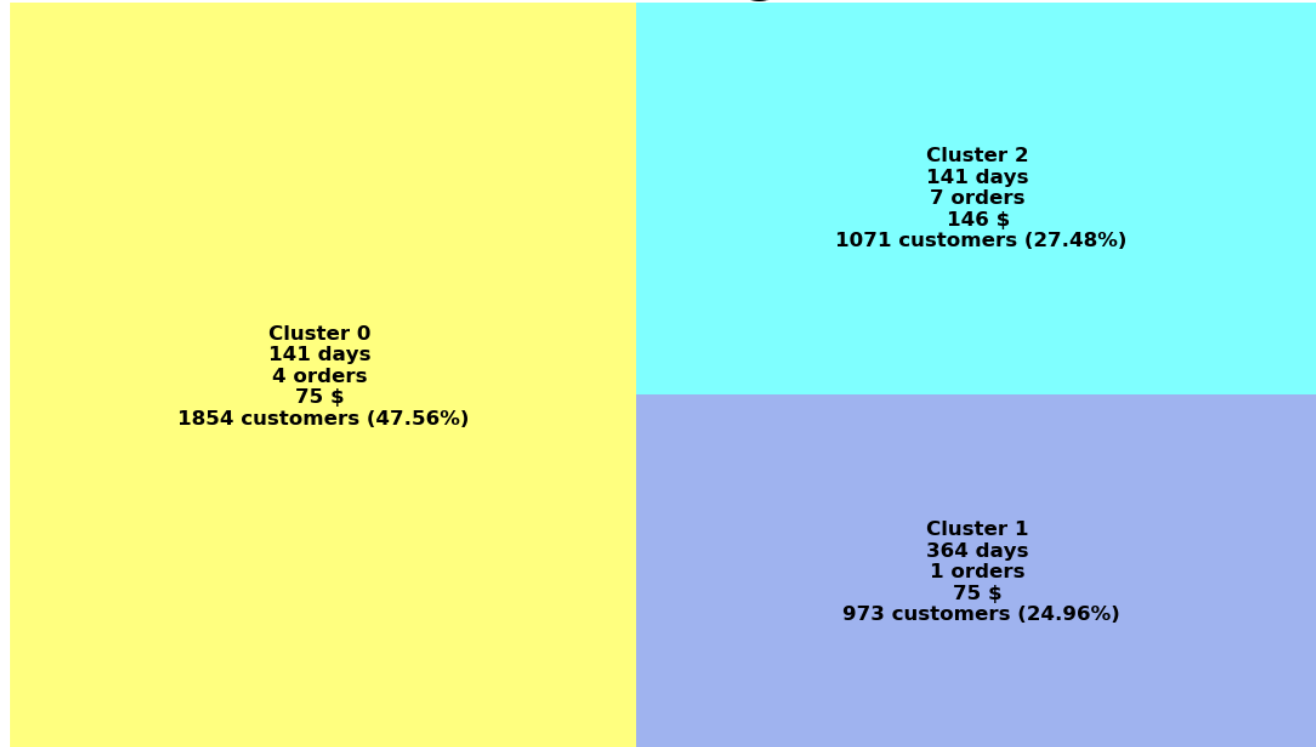


k=3

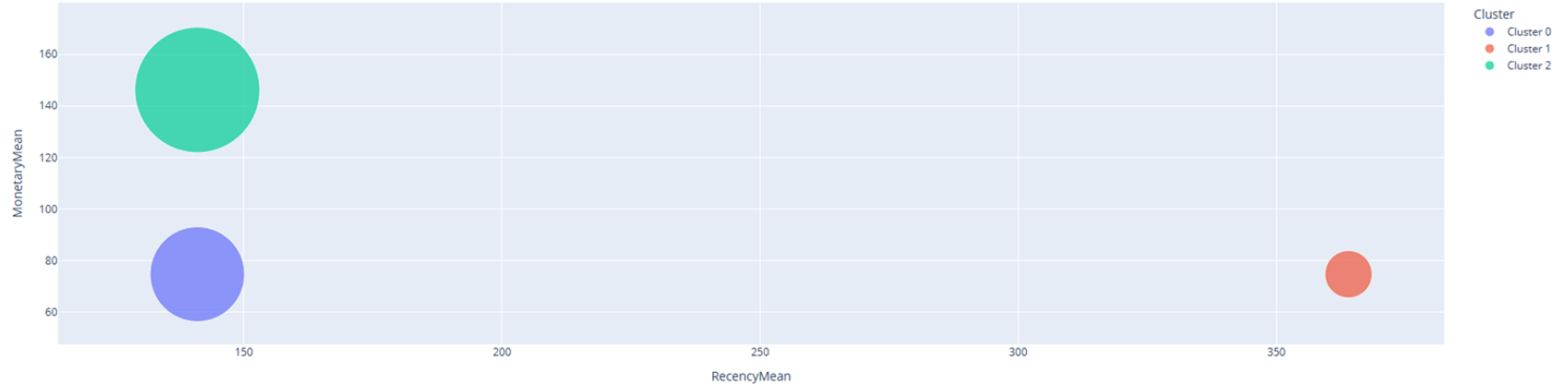


# KMEANS (SCALE DATA)

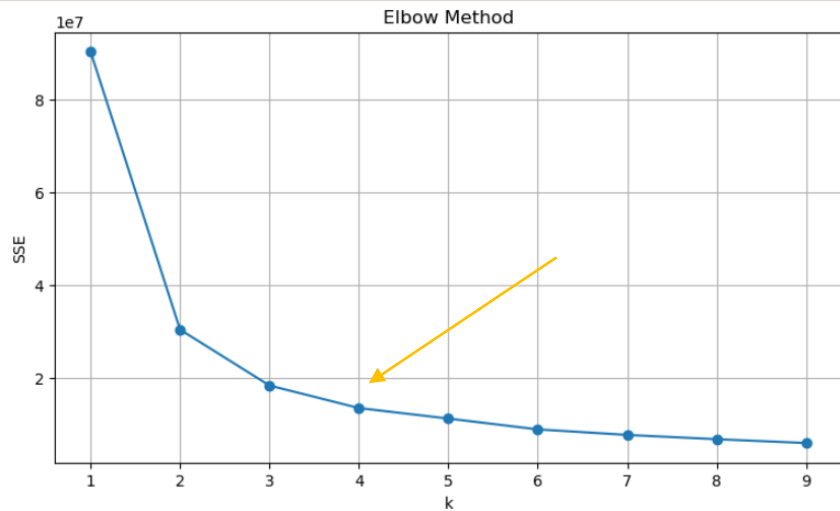
## Customers Segments



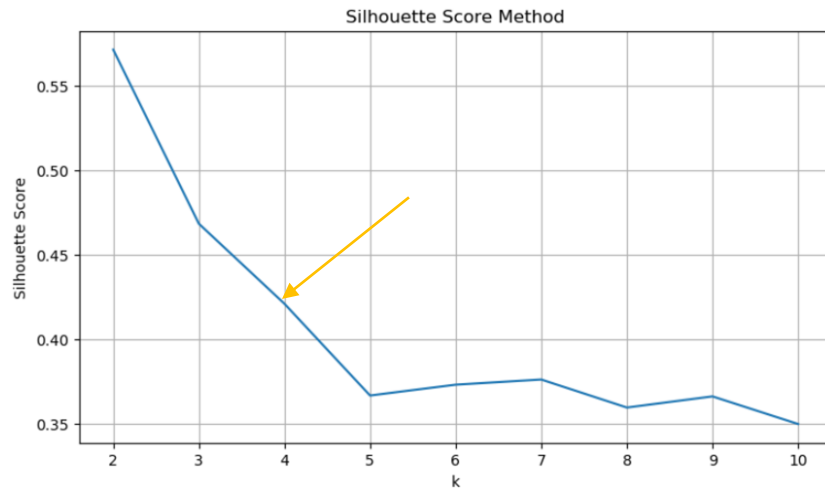
# KMEANS (SCALE DATA)



# KMEANS (REMOVE OUTLIER)

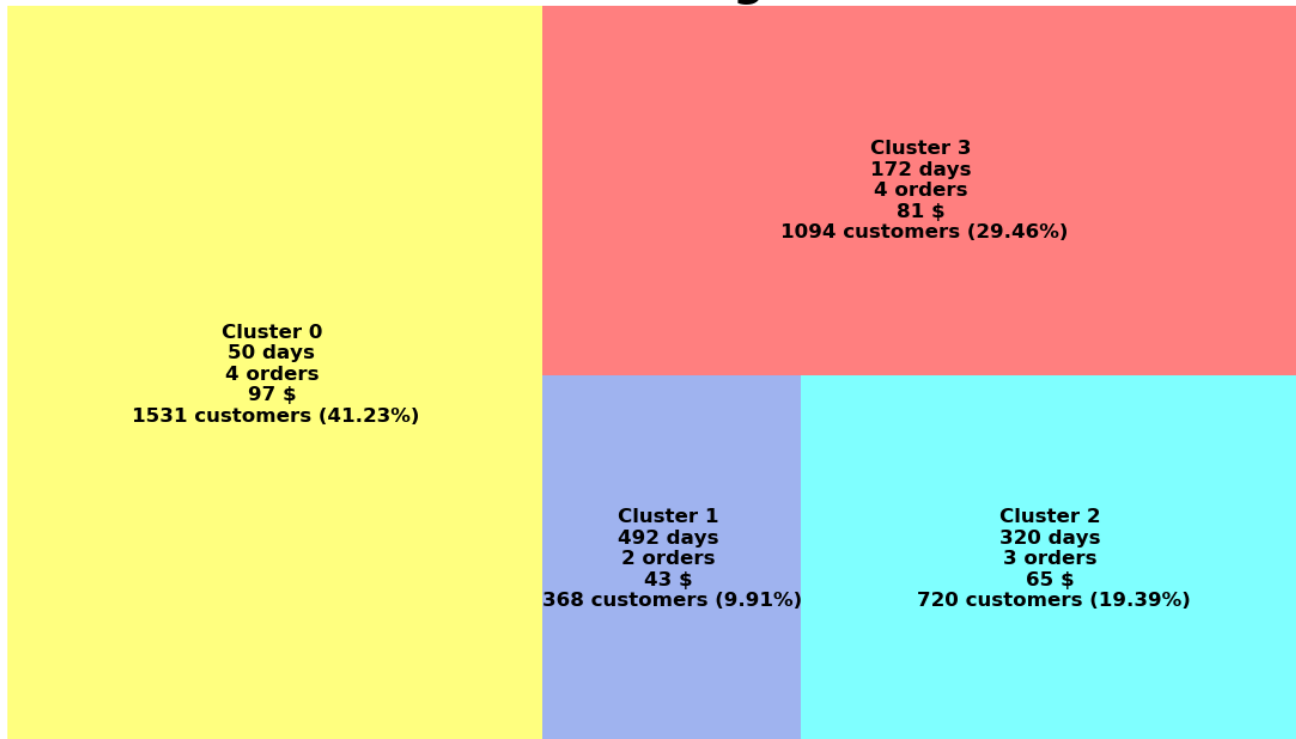


$k=4$



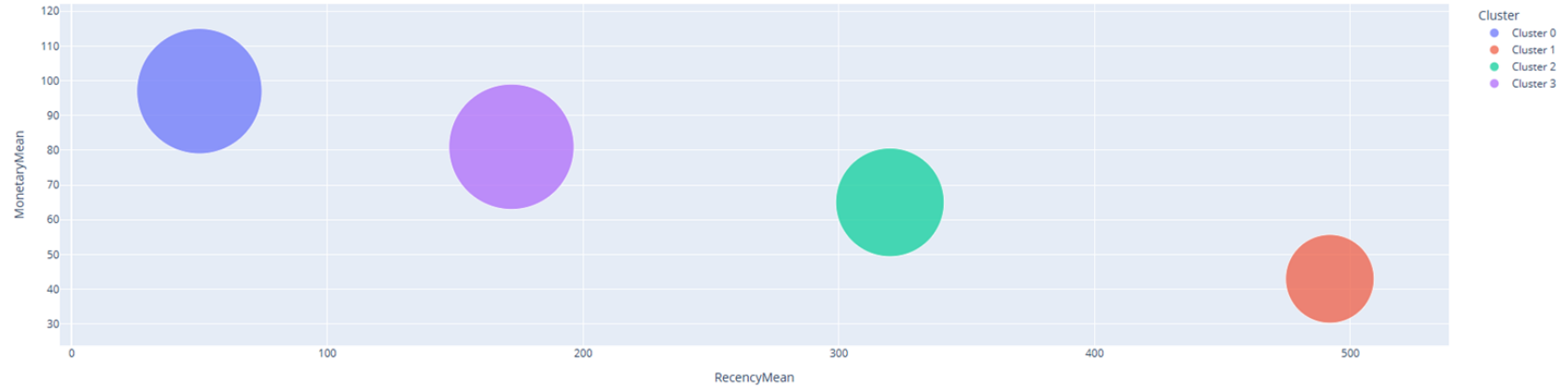
# KMEANS (REMOVE OUTLIER)

## Customers Segments

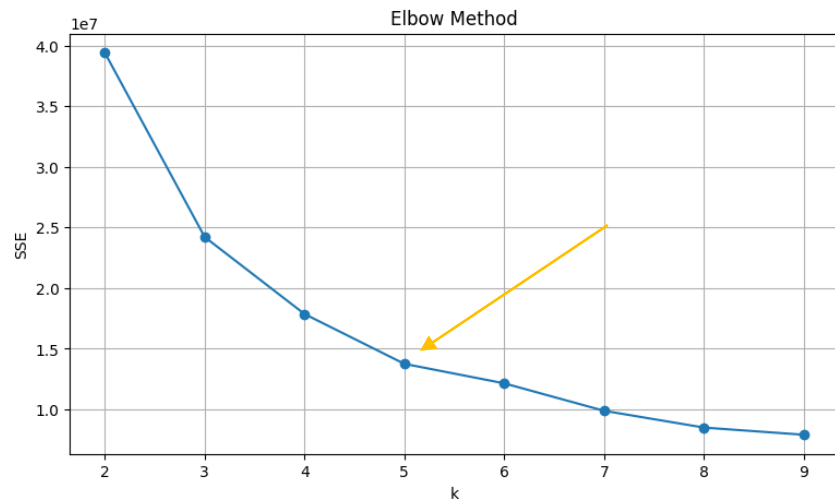




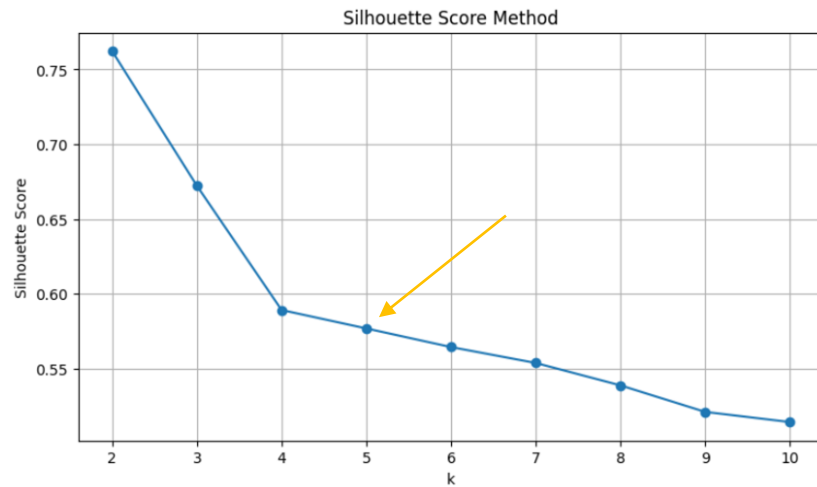
# KMEANS (REMOVE OUTLIER)



# KMEANS (PYSPARK)

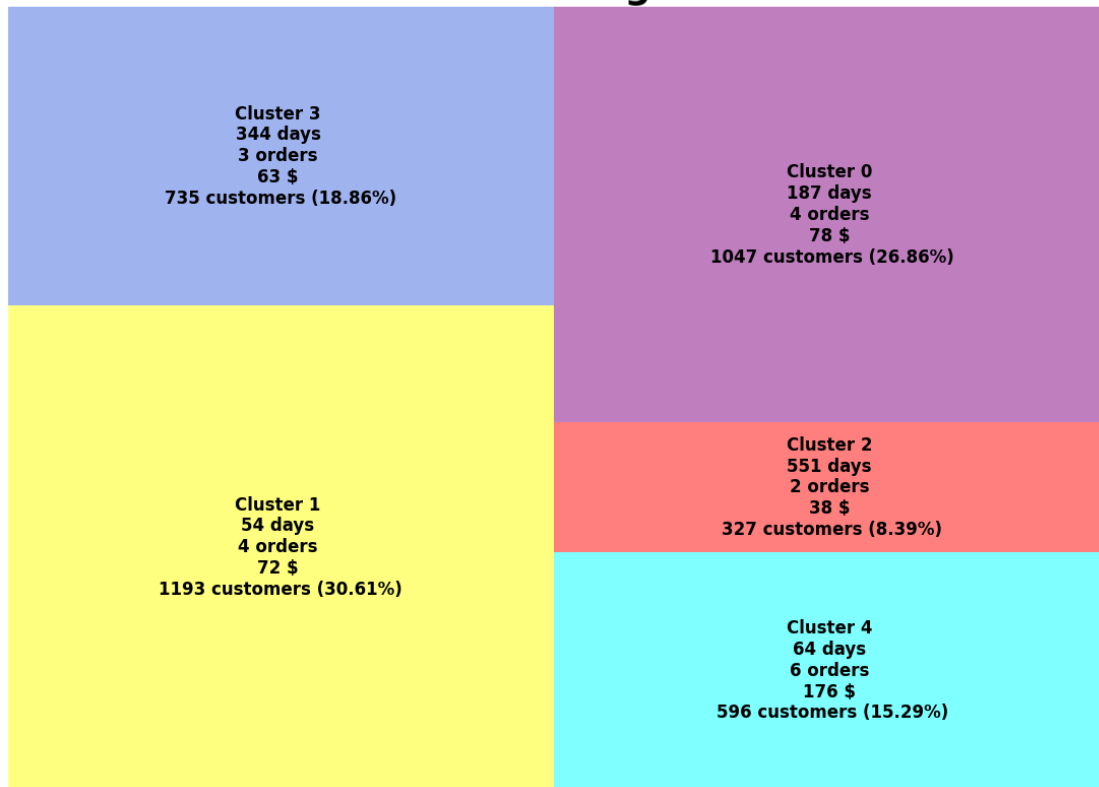


$k=5$

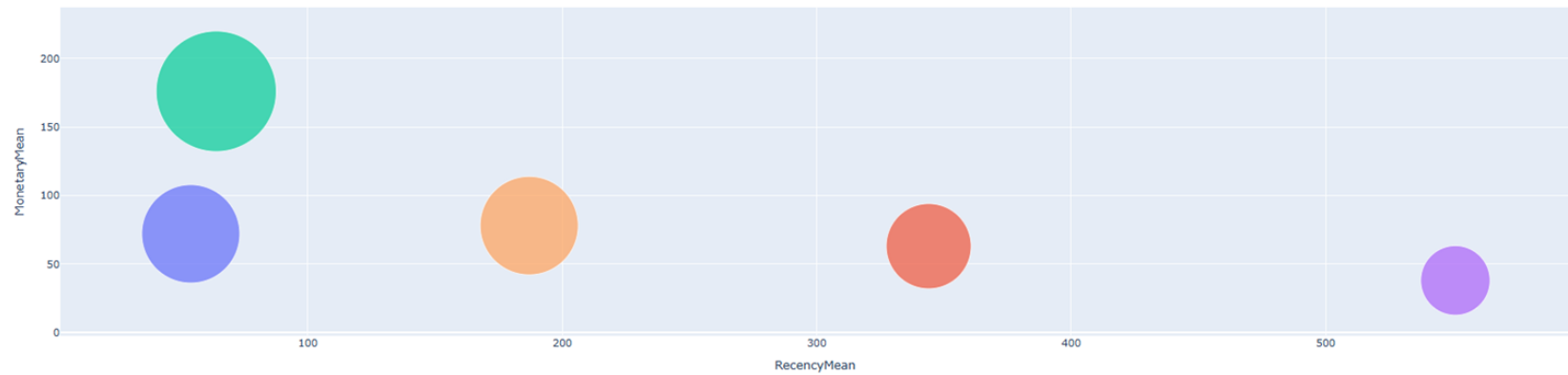


# KMEANS (PYSPARK)

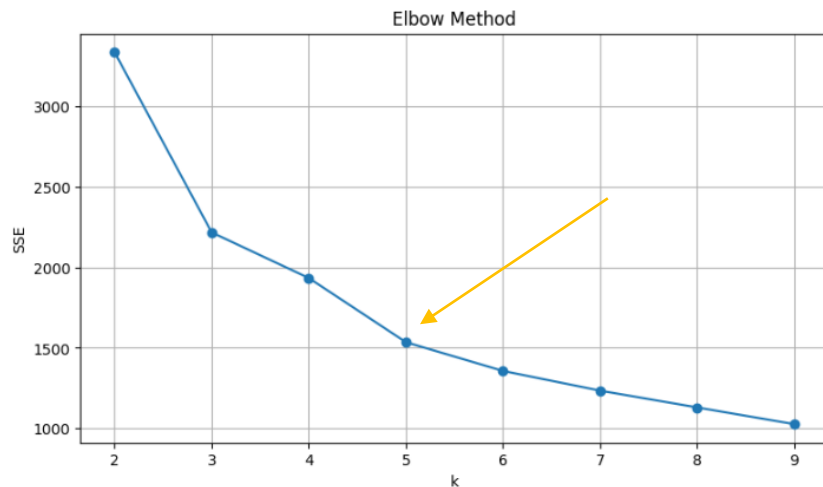
## Customers Segments



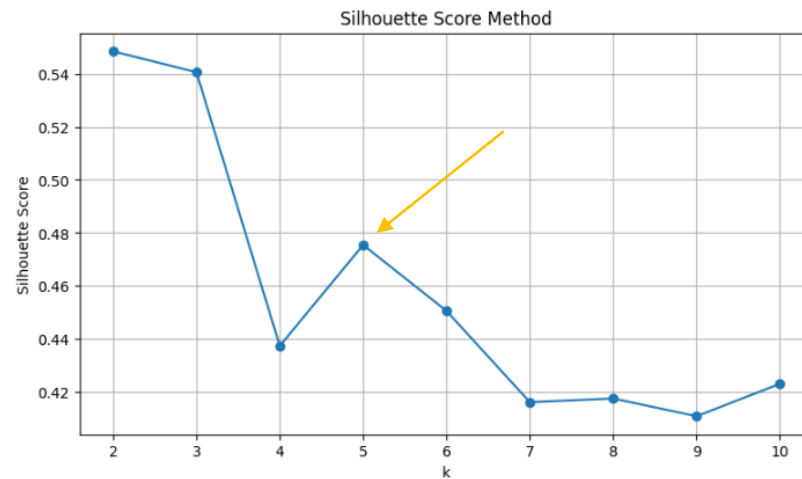
# KMEANS (PYSPARK)



# KMEANS (PYSPARK + SCALE DATA)

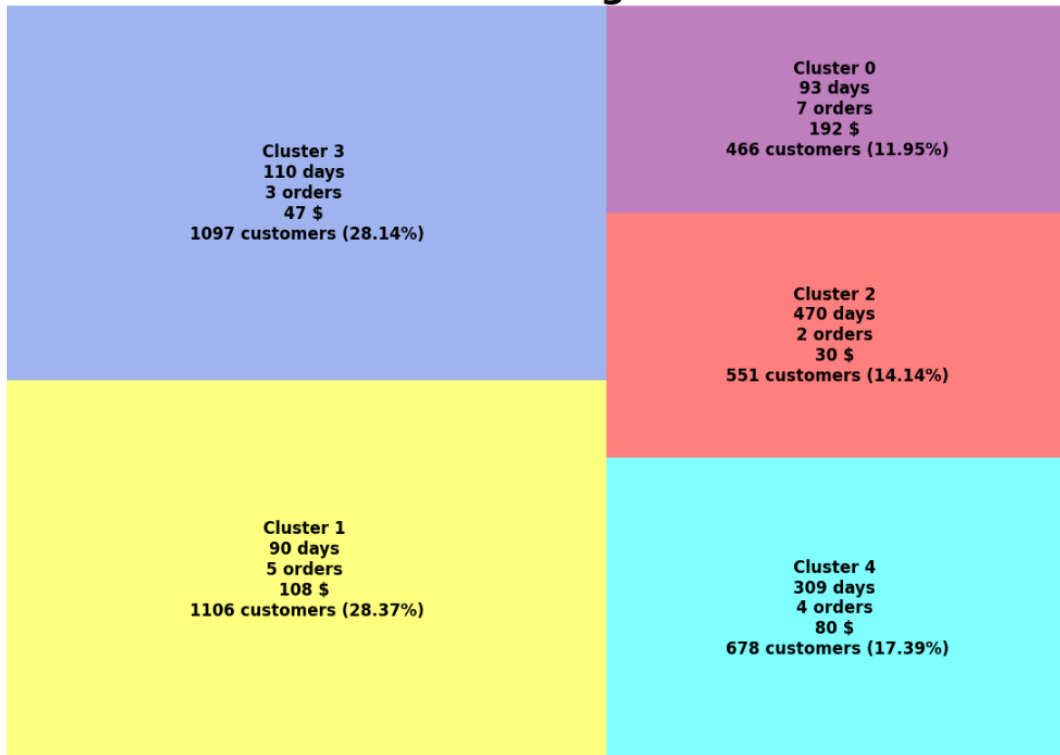


k=5

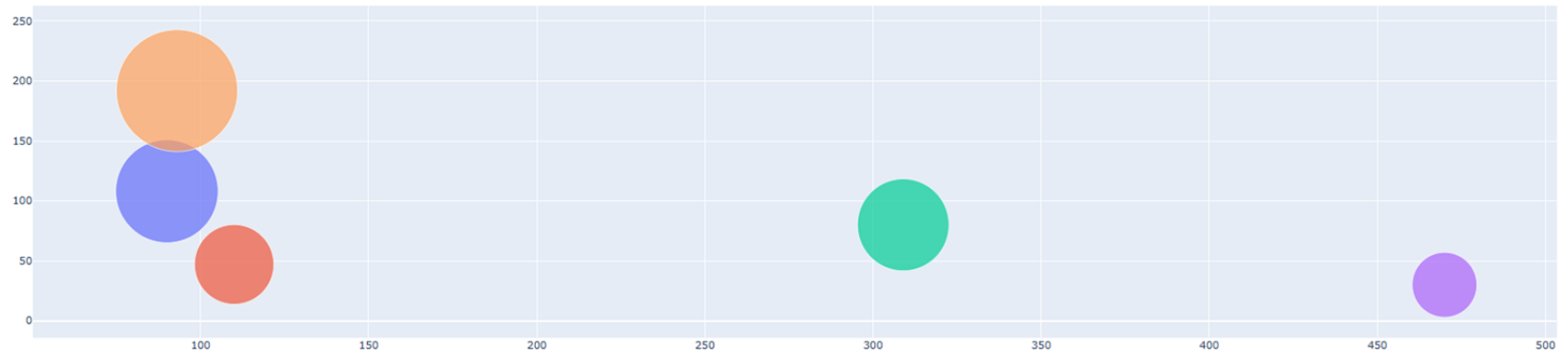


# KMEANS (PYSPARK + SCALE DATA)

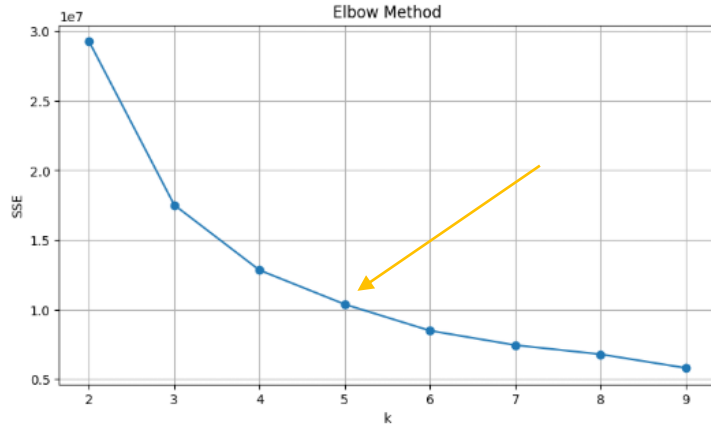
## Customers Segments



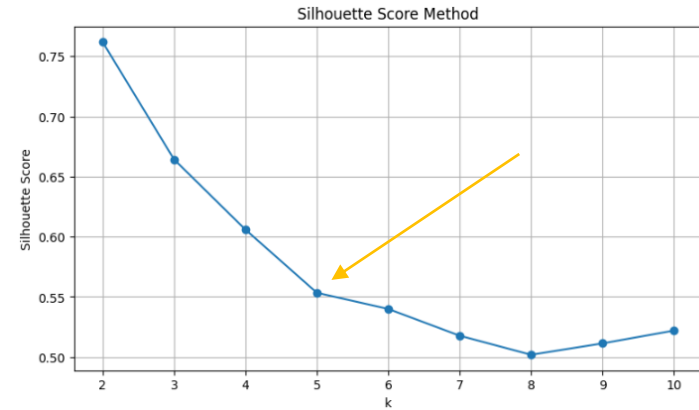
# KMEANS (PYSPARK + SCALE DATA)



# KMEANS (PYSPARK + REMOVE OUTLIERS)



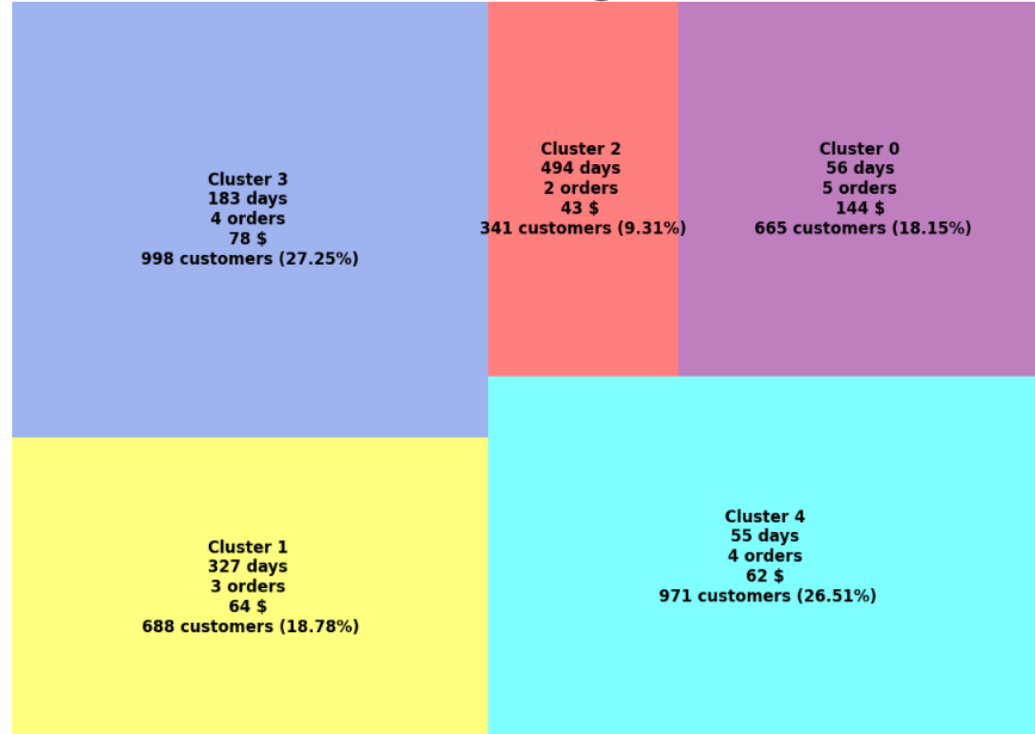
$k=5$



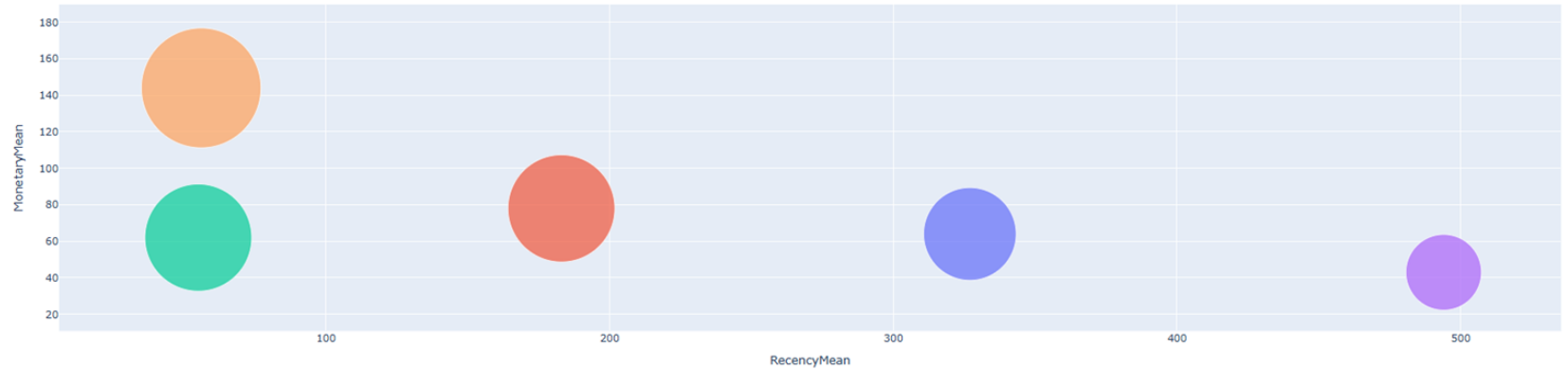


# KMEANS (PYSPARK + REMOVE OUTLIERS)

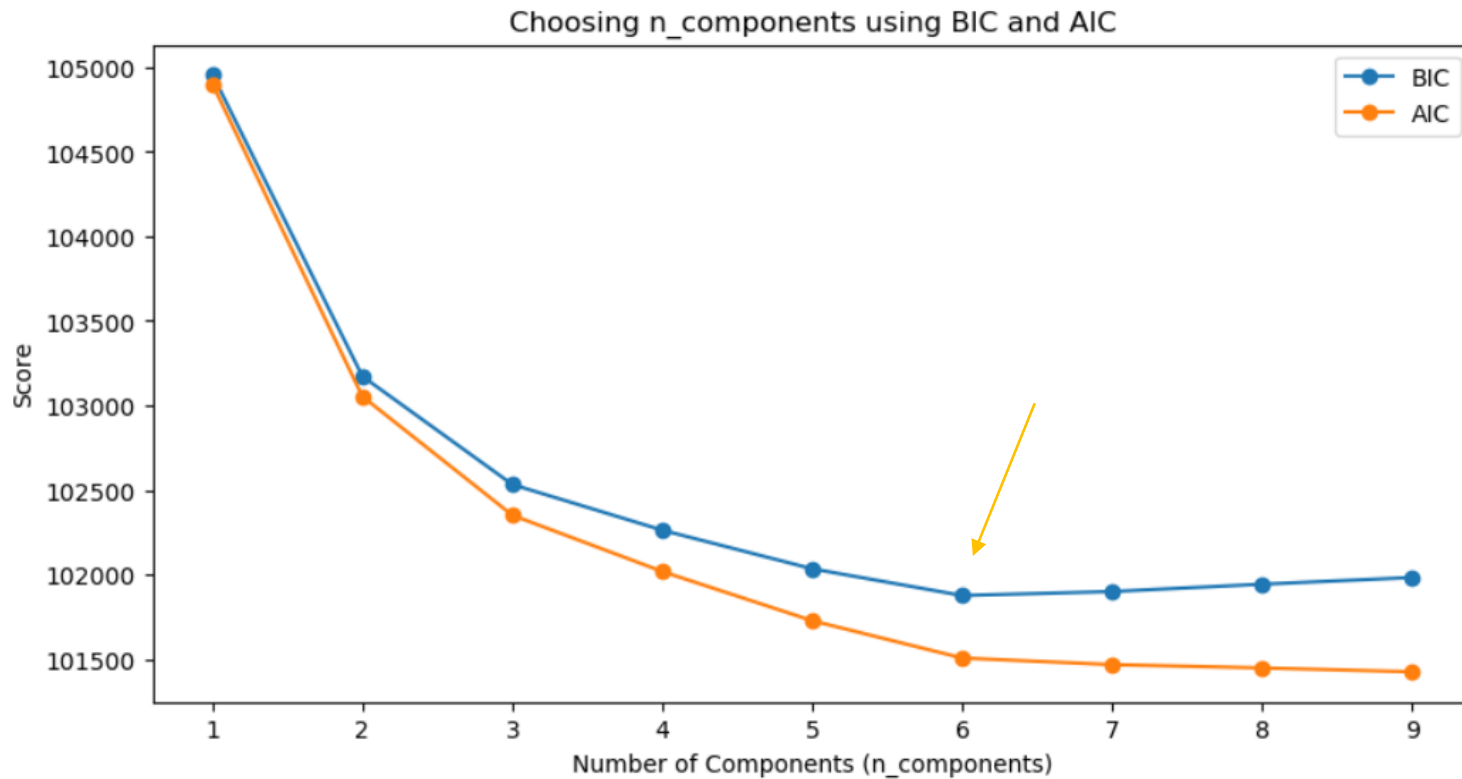
**Customers Segments**



# KMEANS (PYSPARK + REMOVE OUTLIERS)



# GMM



# GMM

## Customers Segments

LOYAL

Cluster 2  
32 days  
5 orders  
96 \$  
756 customers (19.38%)

AT-RISK LOYAL

Cluster 1  
263 days  
4 orders  
114 \$  
513 customers (13.17%)

HARDCORE

Cluster 0  
90 days  
6 orders  
155 \$  
611 customers (15.69%)

AT-RISK

Cluster 5  
351 days  
3 orders  
55 \$  
647 customers (16.61%)

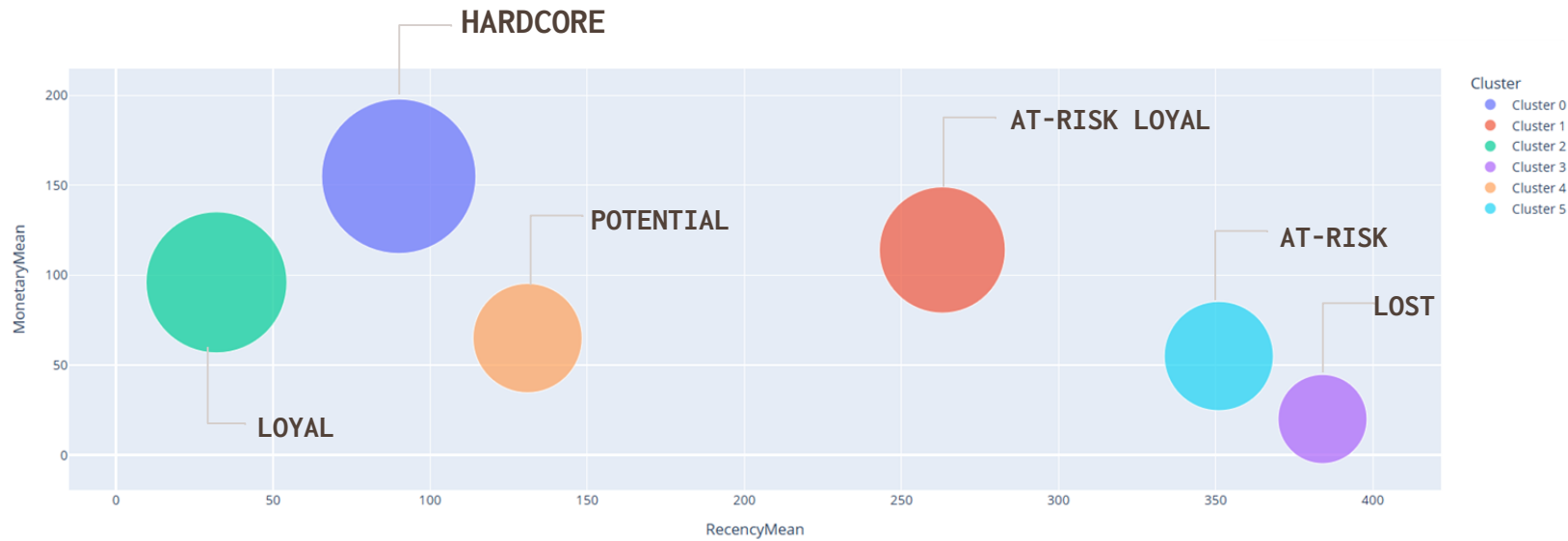
POTENTIAL

Cluster 4  
131 days  
3 orders  
65 \$  
932 customers (23.9%)

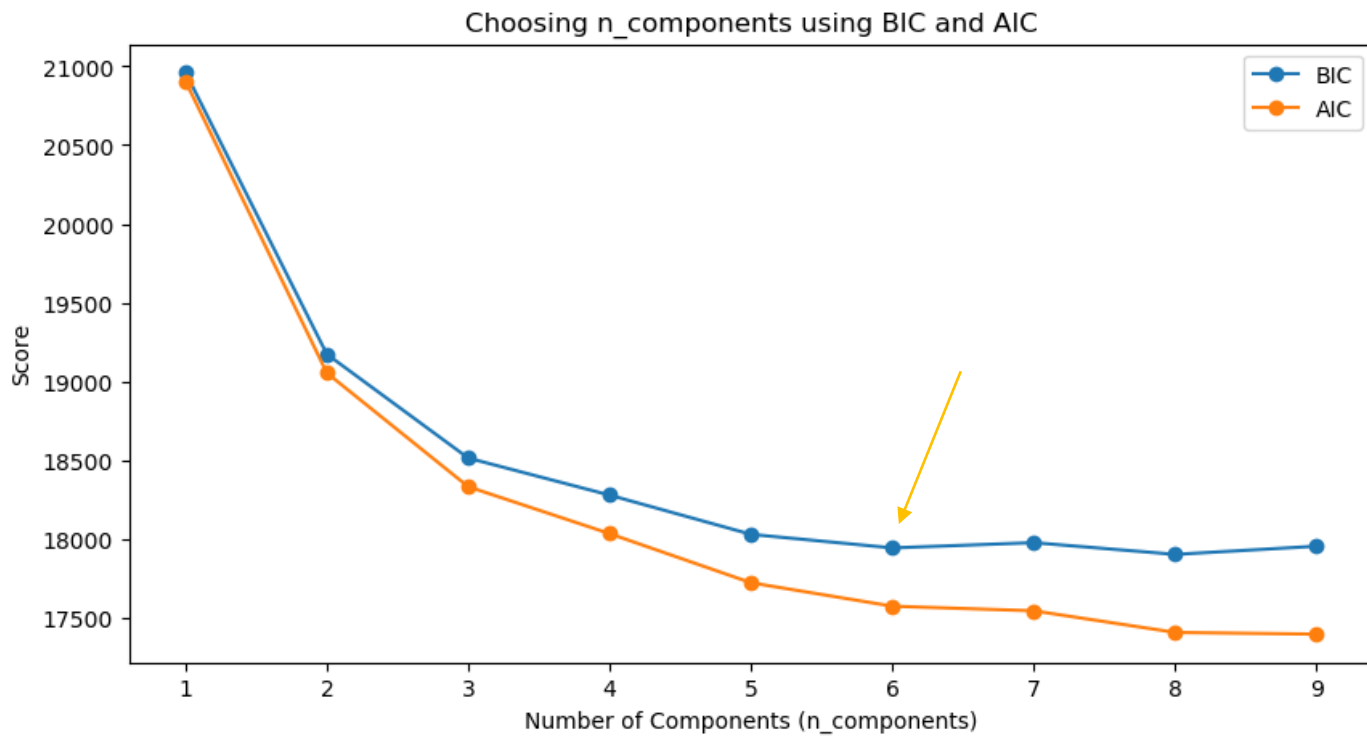
LOST

Cluster 3  
384 days  
2 orders  
20 \$  
438 customers (11.25%)

# GMM

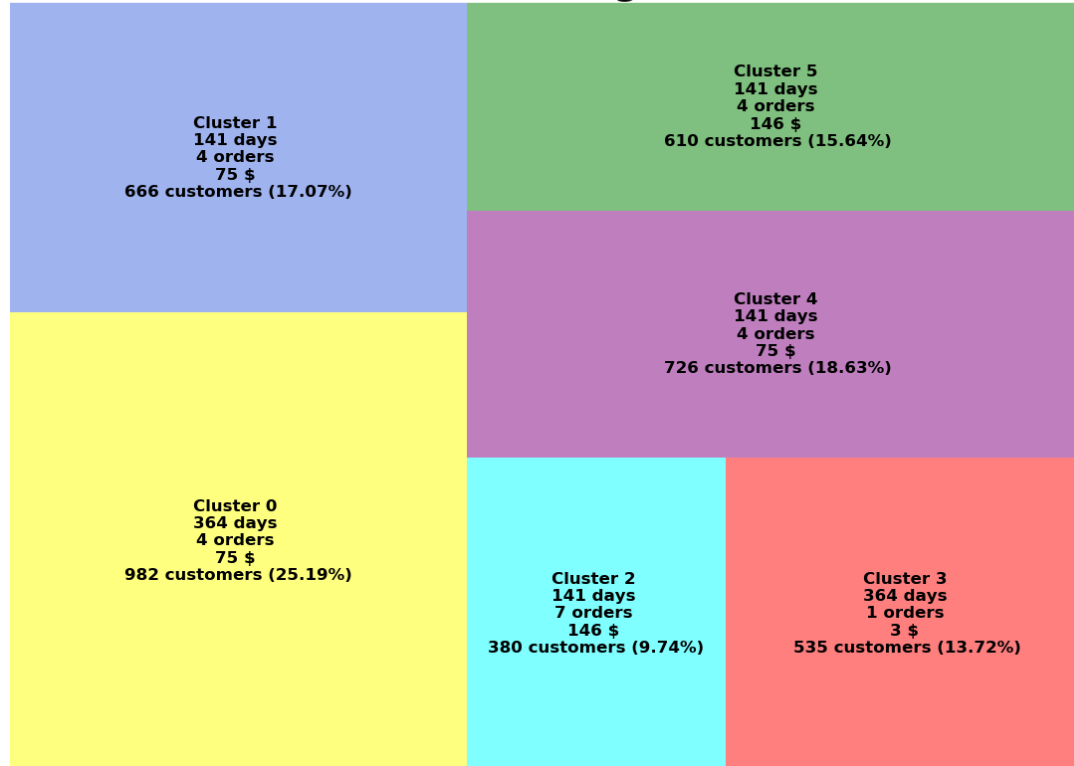


# GMM (SCALE DATA)

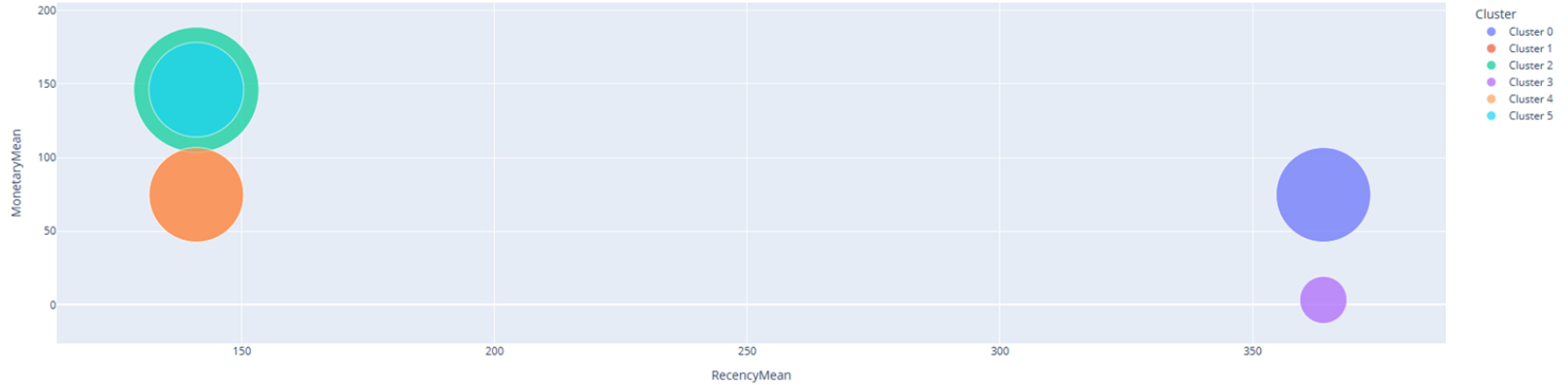


# GMM (SCALE DATA)

## Customers Segments

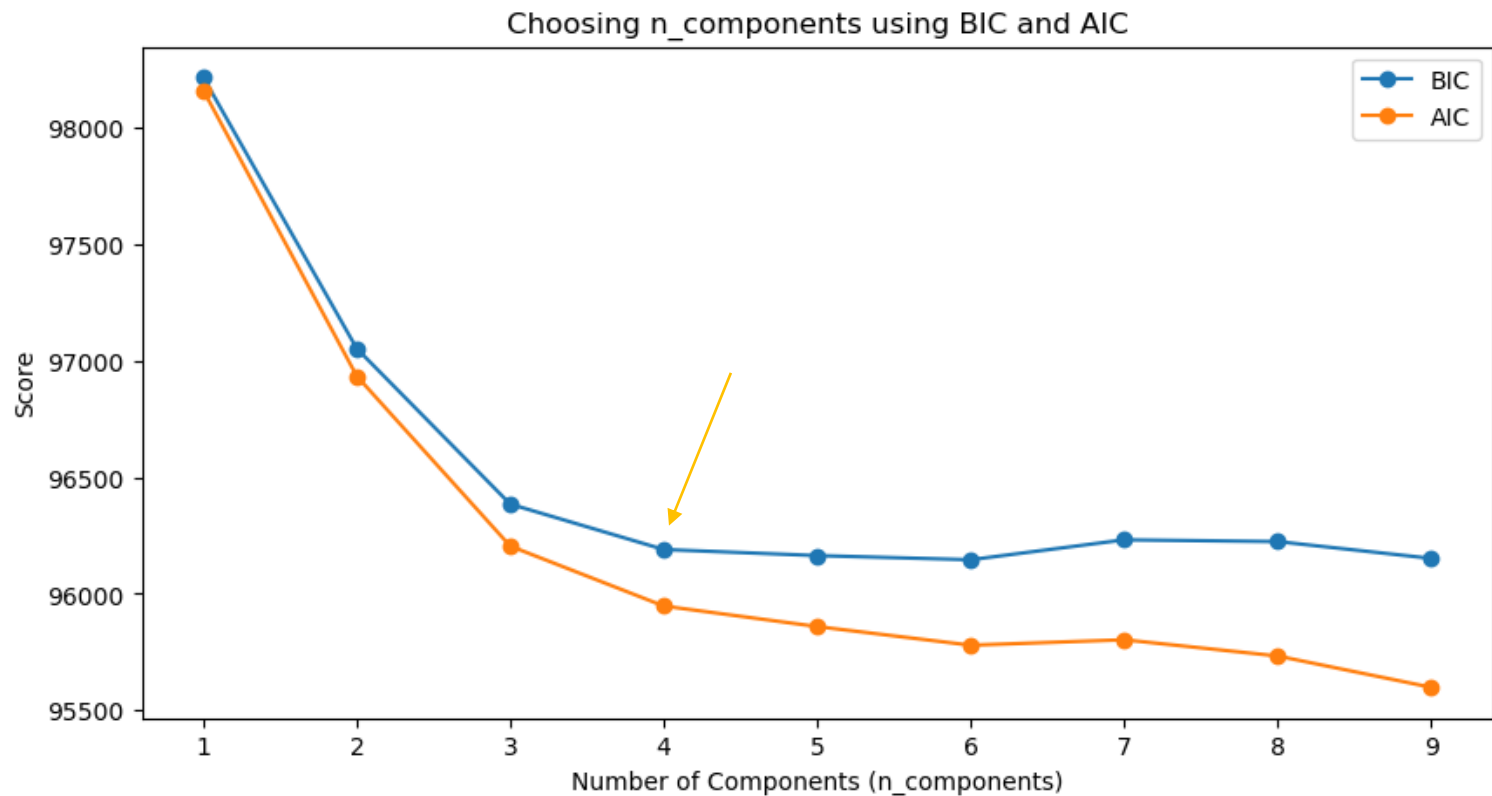


# GMM (SCALE DATA)



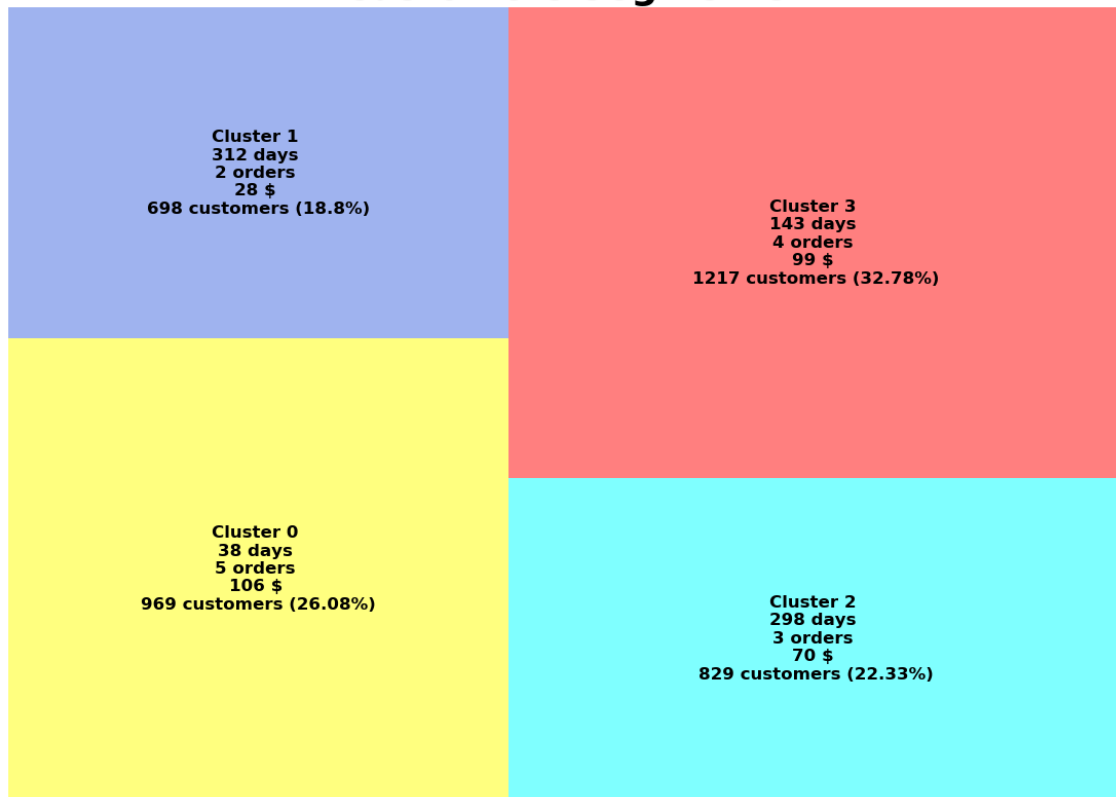


# GMM (REMOVE OUTLIER)

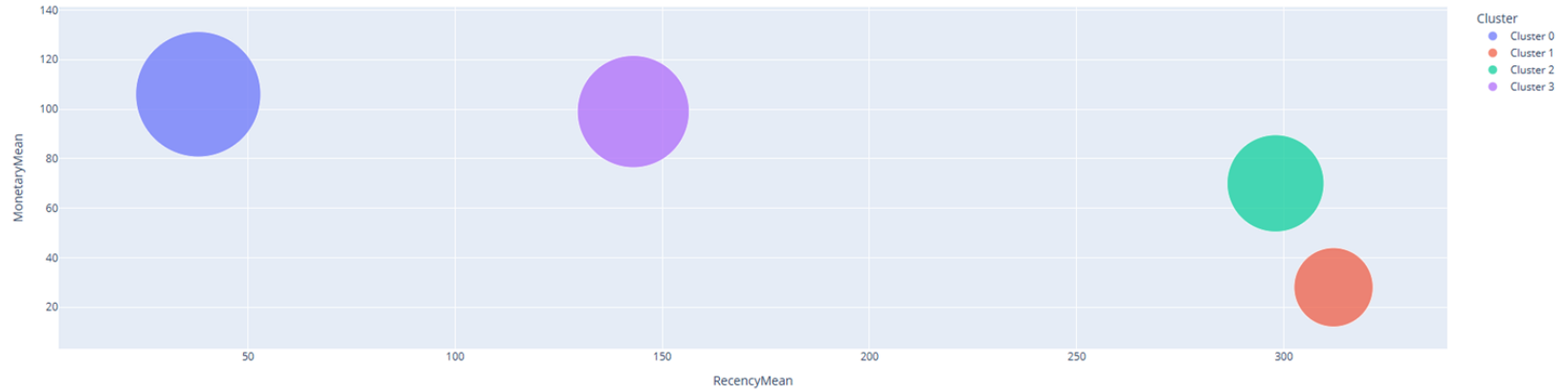


# GMM (REMOVE OUTLIER)

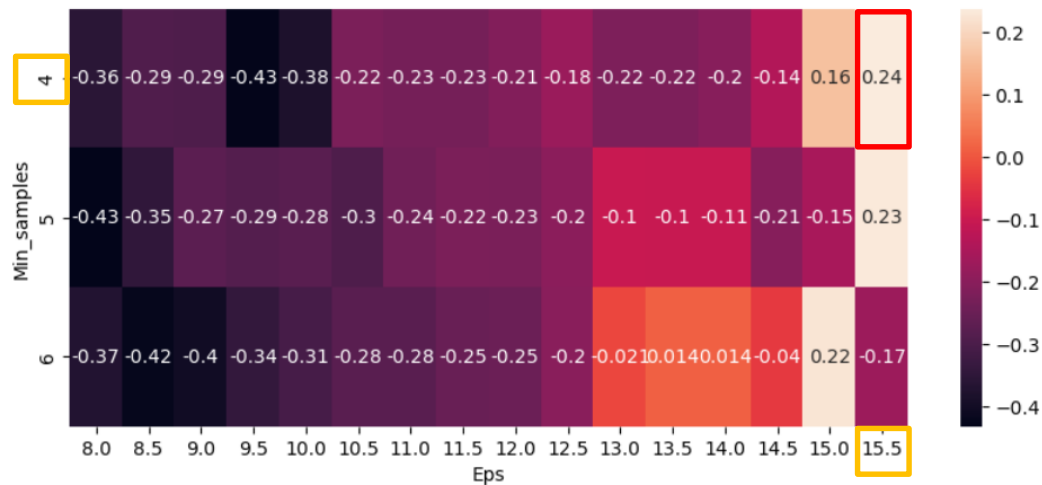
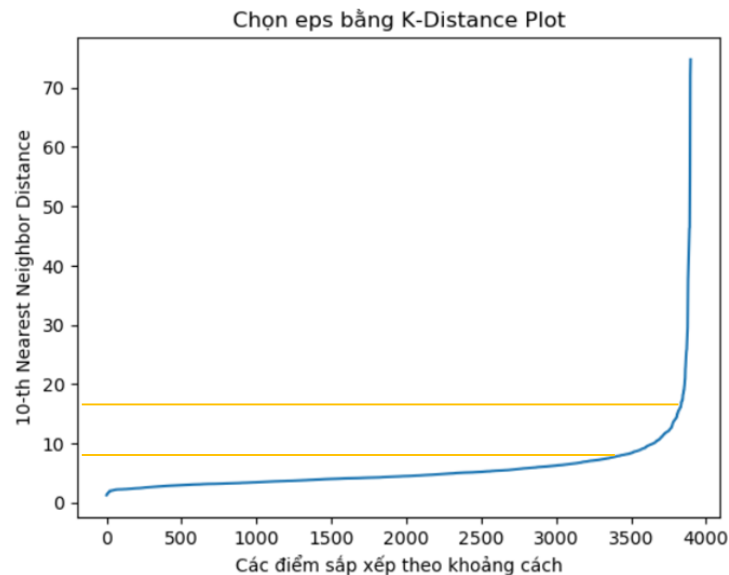
## Customers Segments



# GMM (REMOVE OUTLIER)



# DBSCAN

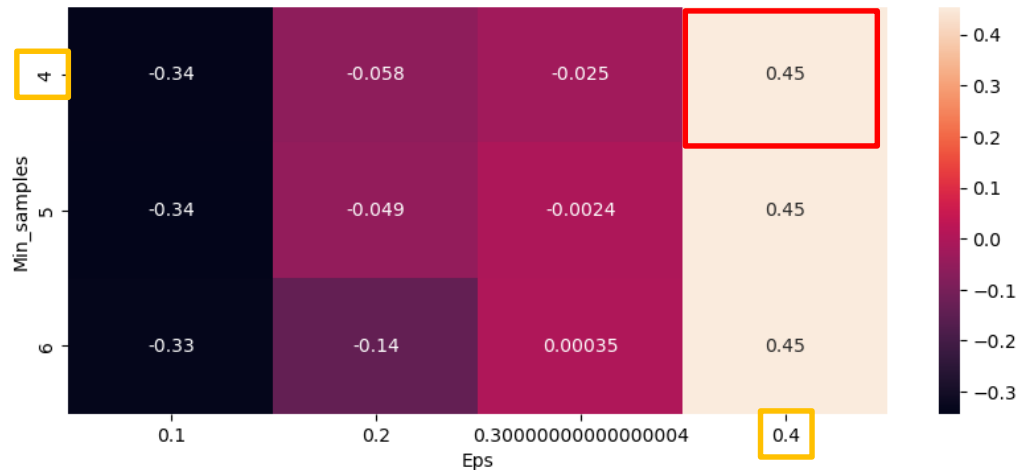
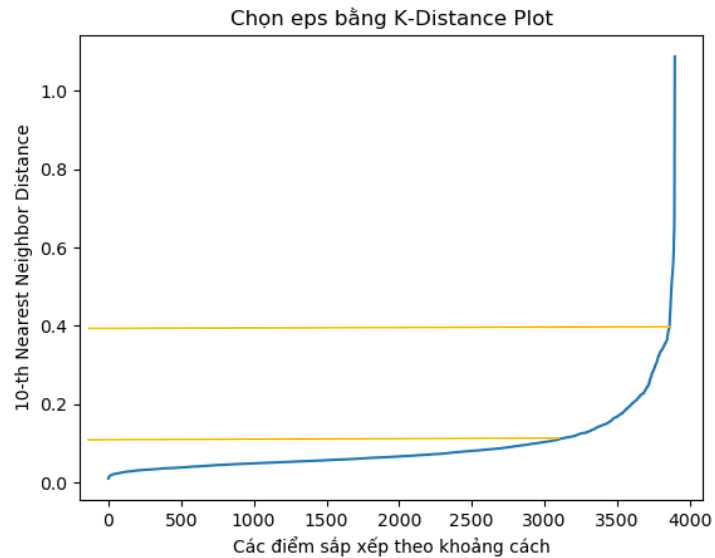


Nhiều (-1): 53 điểm

Cụm 0: 3839 điểm

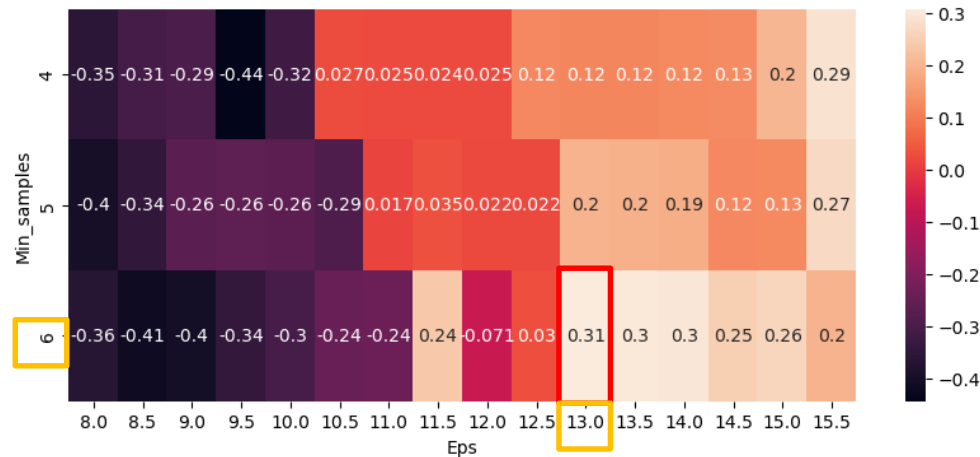
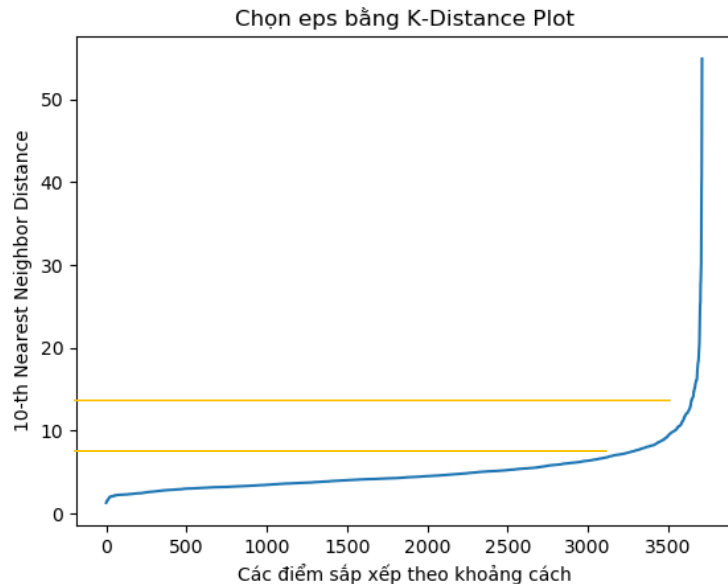
Cụm 1: 6 điểm

# DBSCAN (SCALE DATA)



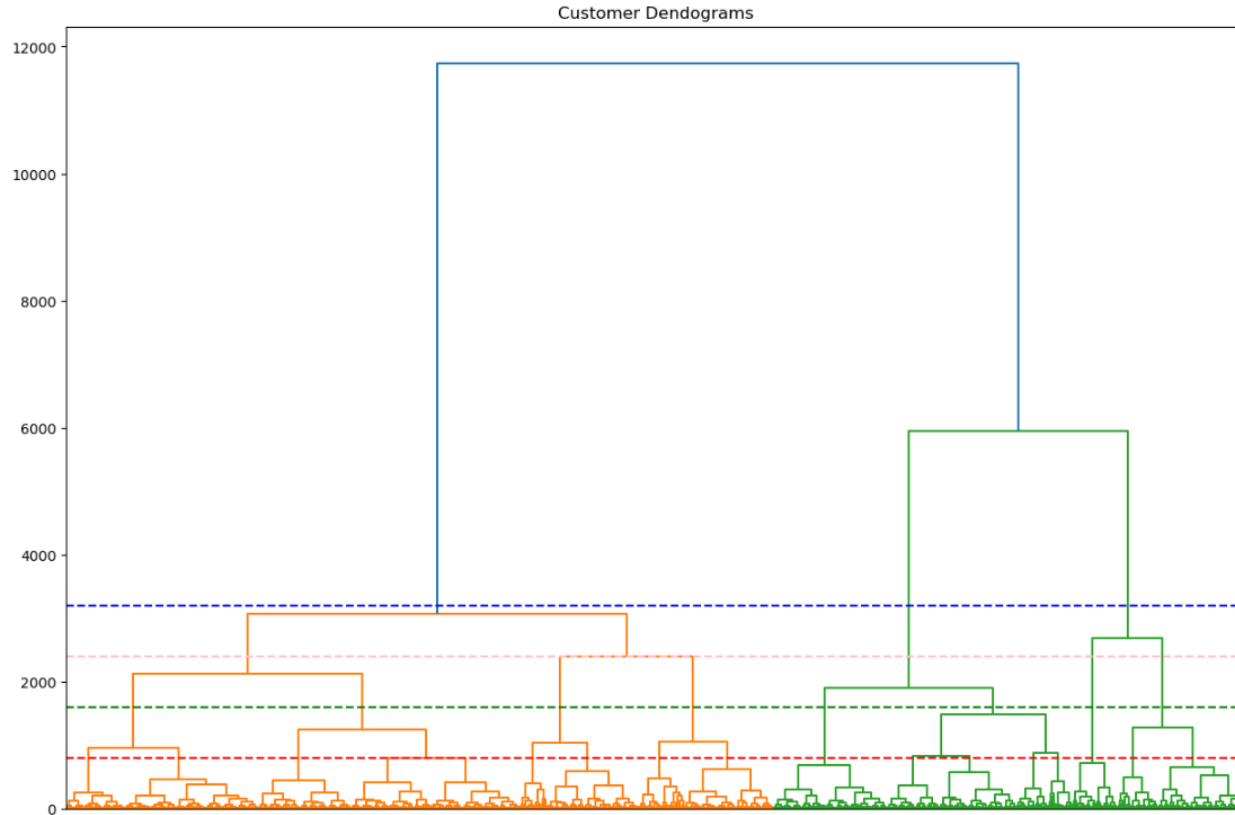
Nhiều (-1): 25 điểm  
Cụm 0: 3873 điểm

# DBSCAN (REMOVE OUTLIERS)



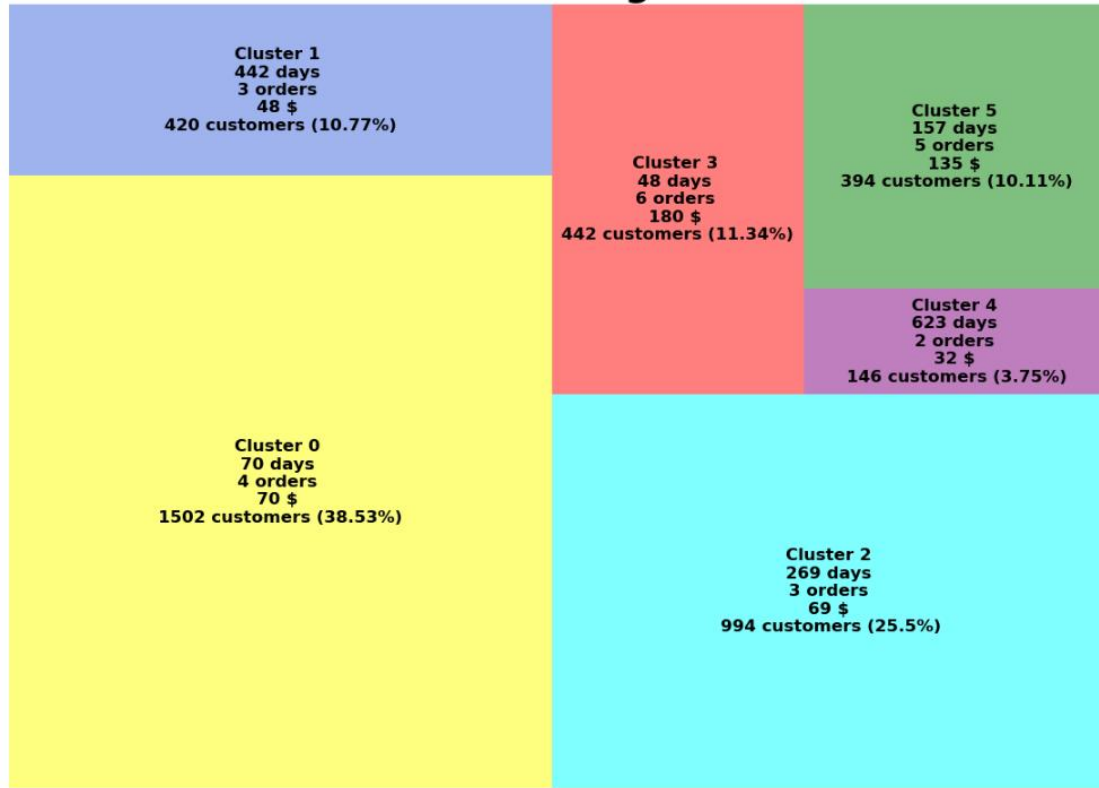
Nhiều (-1): 79 điểm  
Cụm 0: 3634 điểm

# HIERARCHICAL CLUSTERING



# HIERARCHICAL CLUSTERING

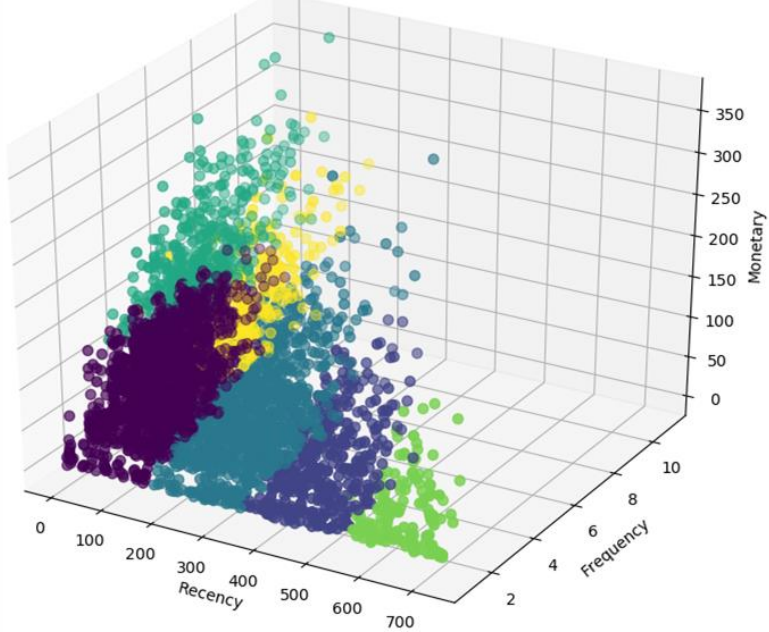
## Customers Segments



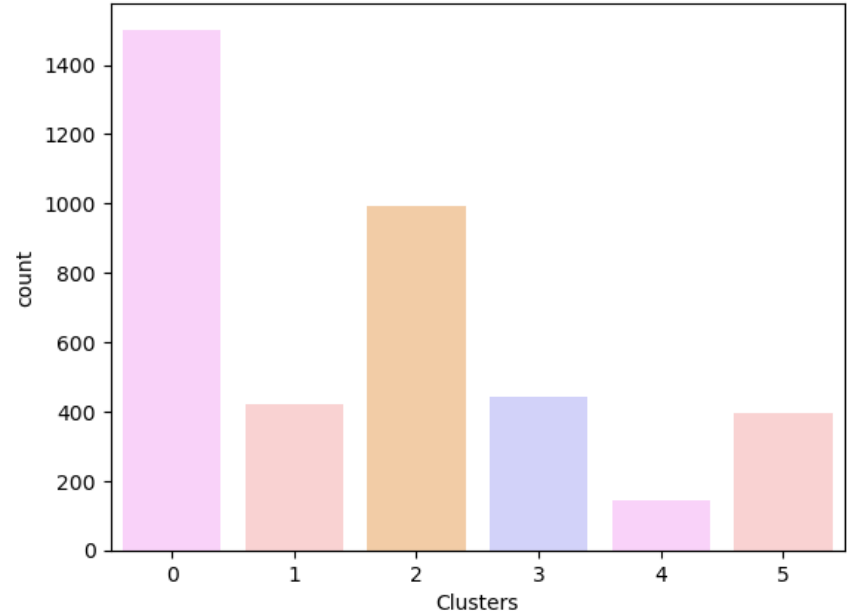


# HIERARCHICAL CLUSTERING

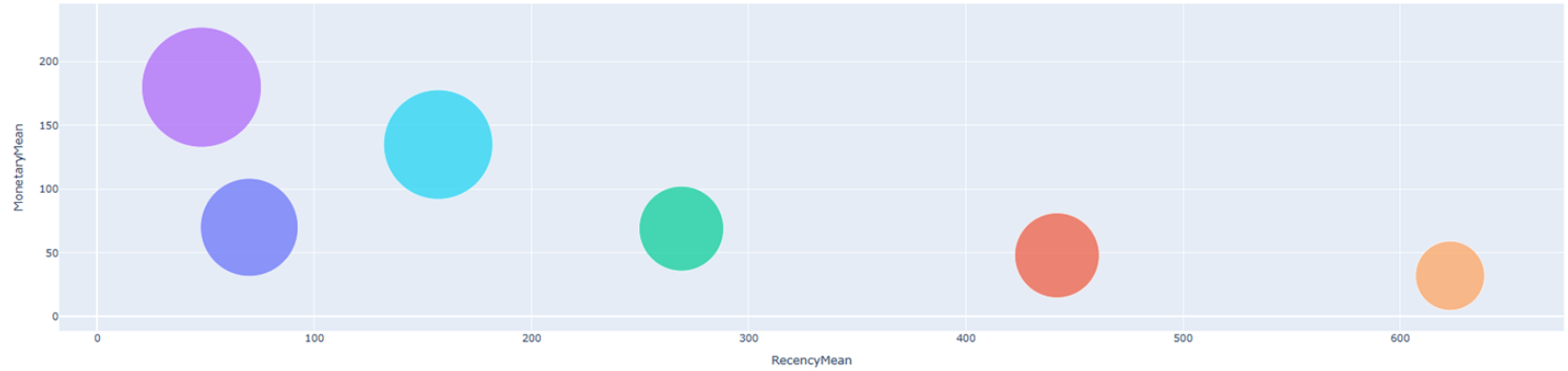
The Plot Of The Clusters



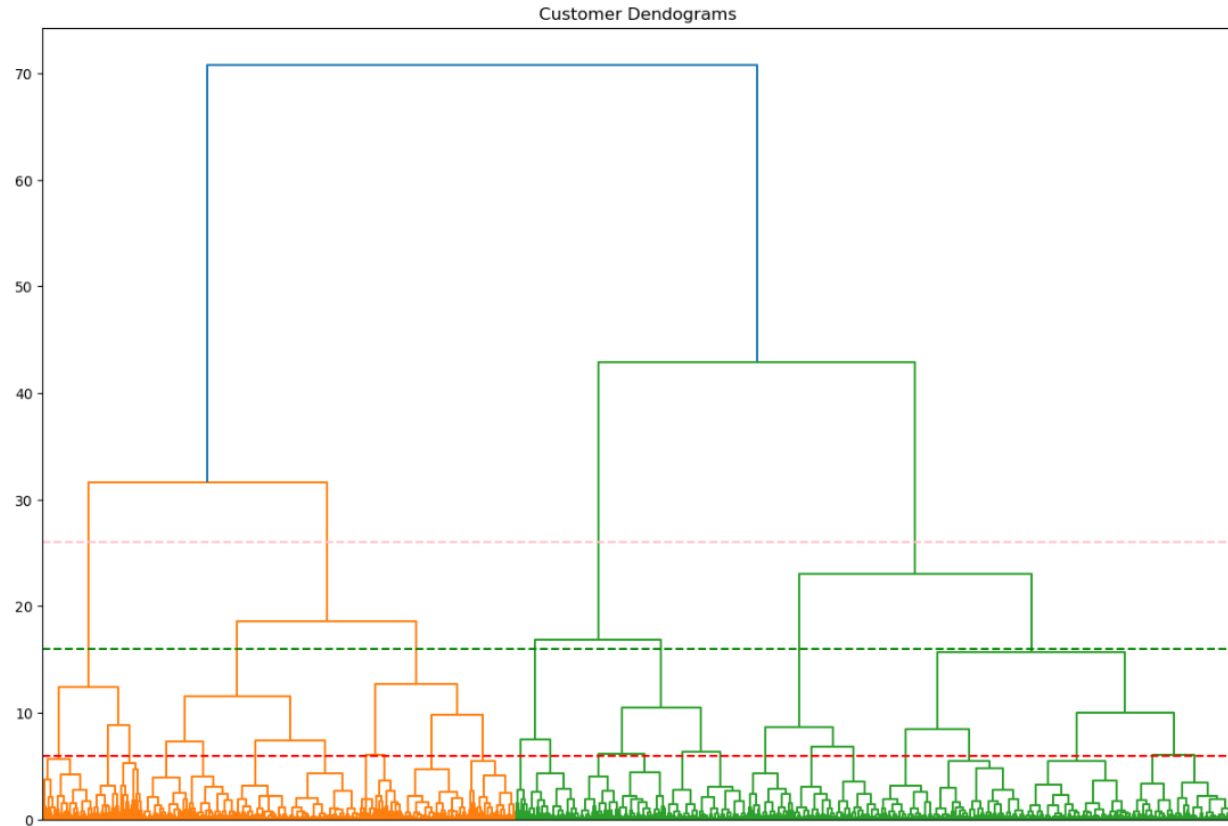
Distribution Of The Clusters



# HIERARCHICAL CLUSTERING

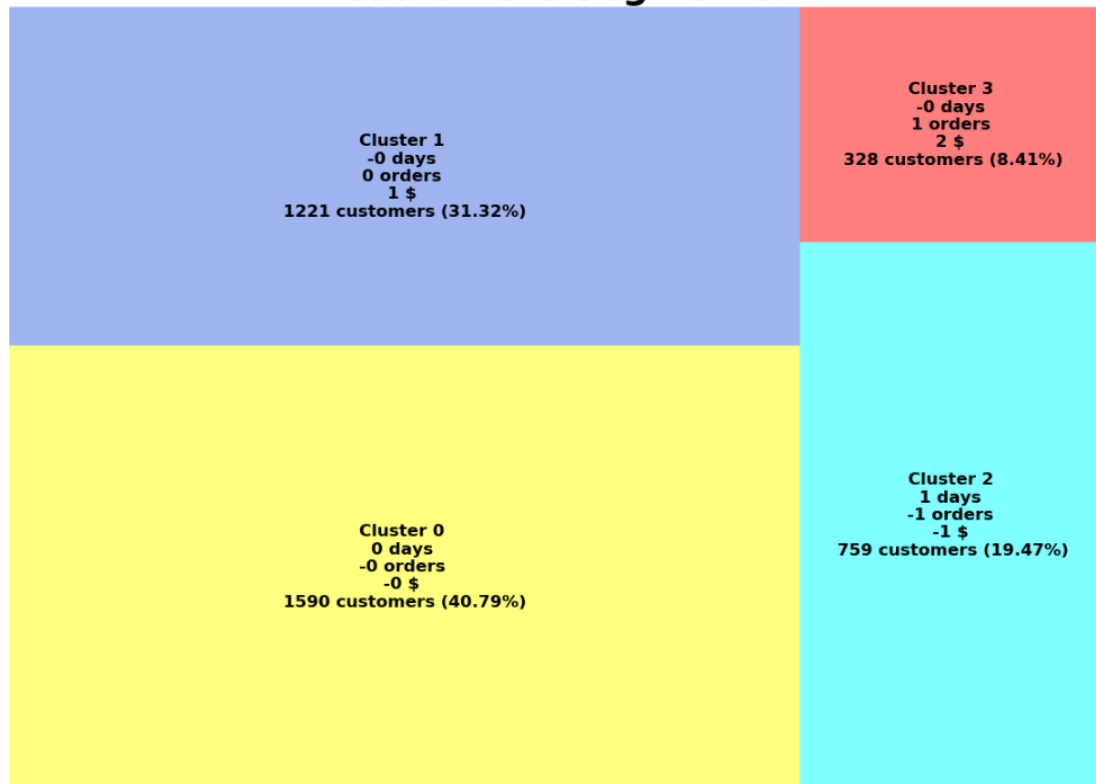


# HIERARCHICAL CLUSTERING (SCALE DATA)

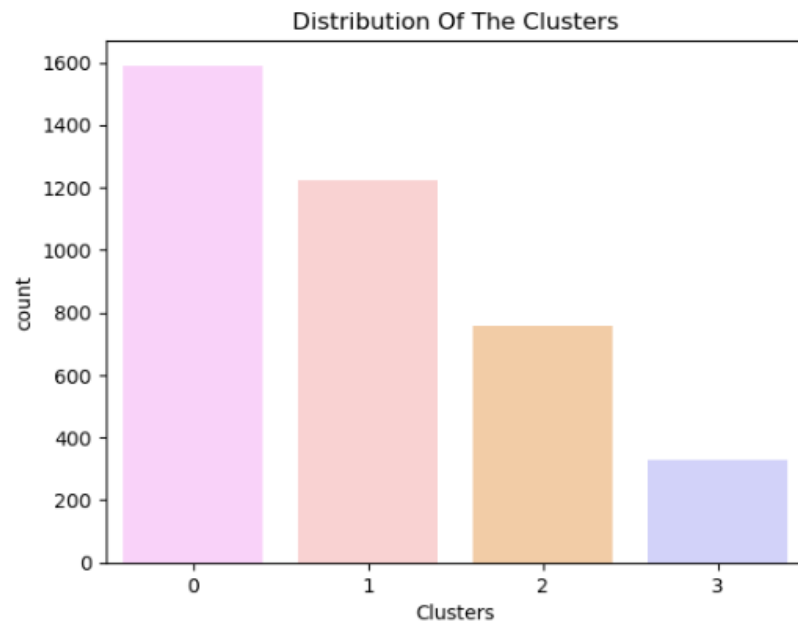
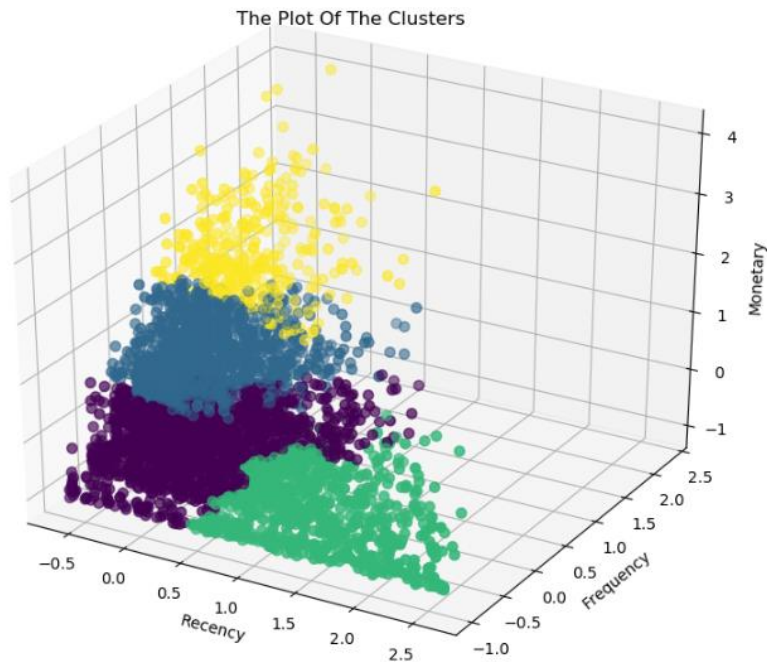


# HIERARCHICAL CLUSTERING (SCALE DATA)

## Customers Segments



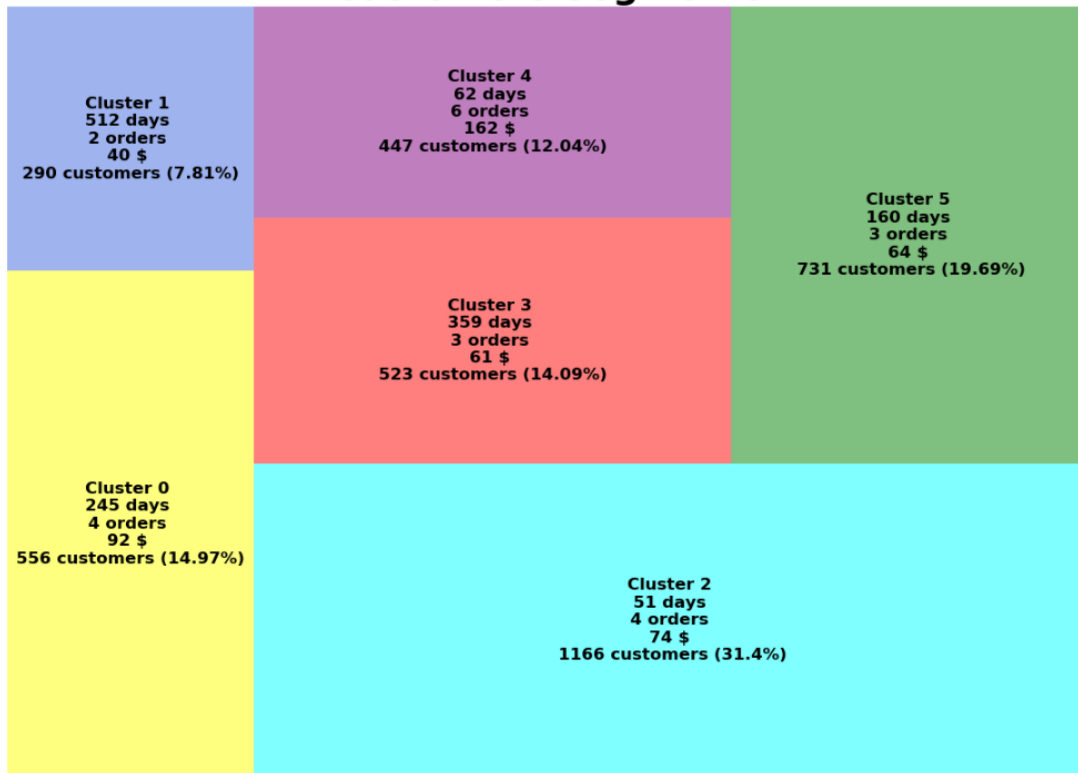
# HIERARCHICAL CLUSTERING (SCALE DATA)



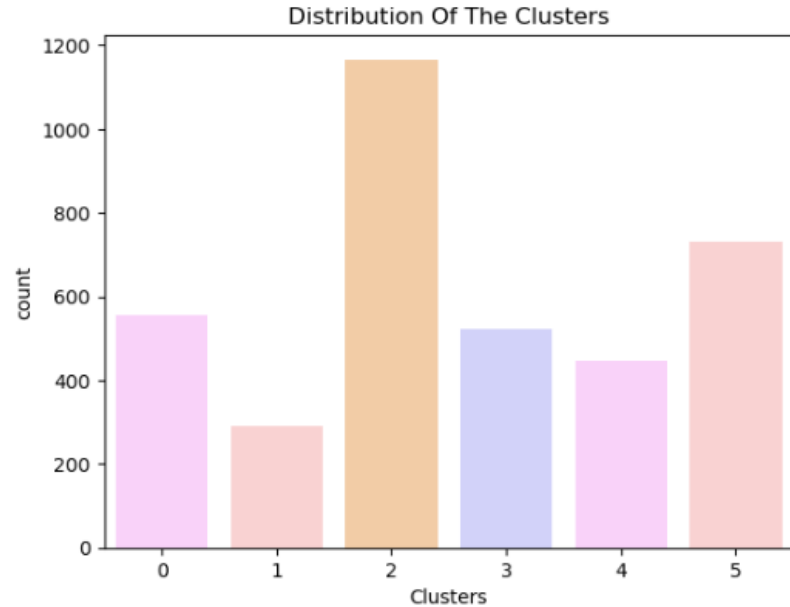
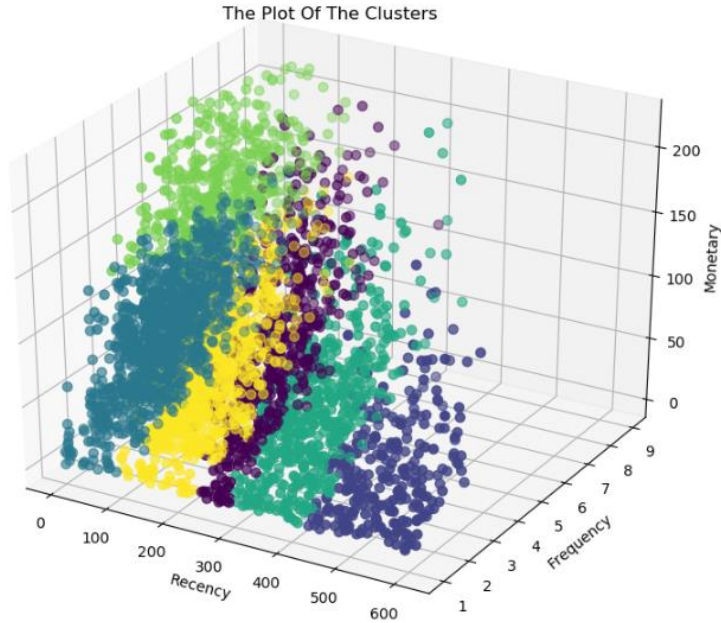


# HIERARCHICAL CLUSTERING (REMOVE OUTLIERS)

## Customers Segments

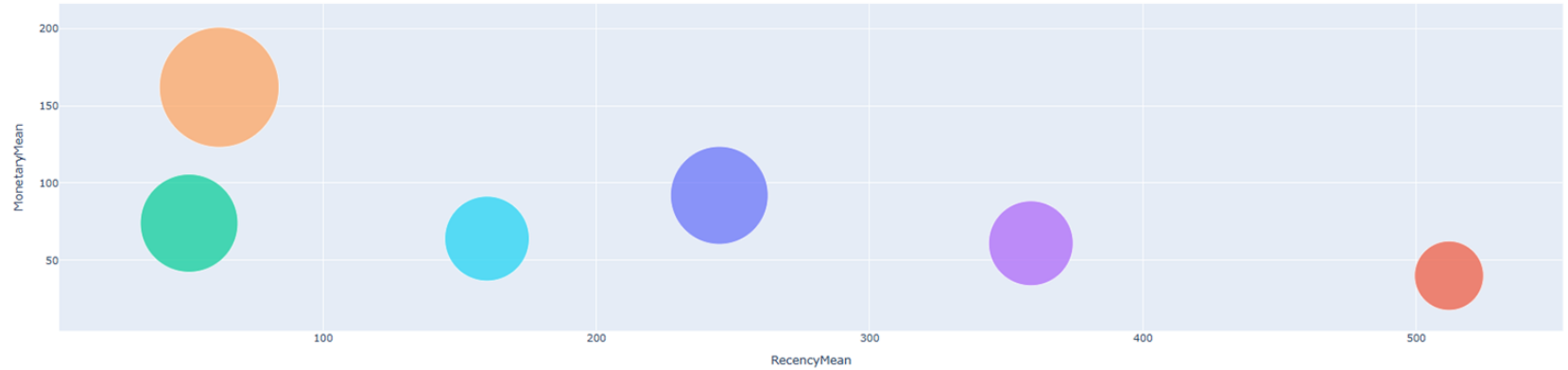


# HIERARCHICAL CLUSTERING (REMOVE OUTLIERS)





# HIERARCHICAL CLUSTERING (REMOVE OUTLIERS)



# EVALUATION

Phương pháp	Số cụm	Silhouette Score	Nhận xét
Manual Segmentation	6	-0.059	Thực hiện phân khúc thủ công dựa trên cảm tính, không phản ánh đúng cấu trúc dữ liệu
KMeans	2	0.57	Các cụm rõ ràng và tách biệt tốt, tuy nhiên 2 cụm không đủ để phản ánh sự đa dạng của khách hàng
<b>KMeans</b>	<b>5</b>	<b>0.4</b>	Cung cấp một sự phân chia chi tiết hơn k=2 mà vẫn giữ được tính chất phân tách tốt giữa các cụm
KMeans (PySpark)	5	0.58	Cùng chia 5 cụm nhưng kết quả Sihoutte tốt hơn machine learning truyền thống nhiều
GMM	6	0.2	Cách cụm khá gần nhau, phân tách chưa tốt
DBSCAN	2 + outliers	0.23	Phần lớn dữ liệu tập trung vào 1 cụm, không hiệu quả
Hierarchical Clustering	6	0.33	Cách cụm phân tách tốt hơn GMM, nhưng vẫn không tối ưu bằng KMeans

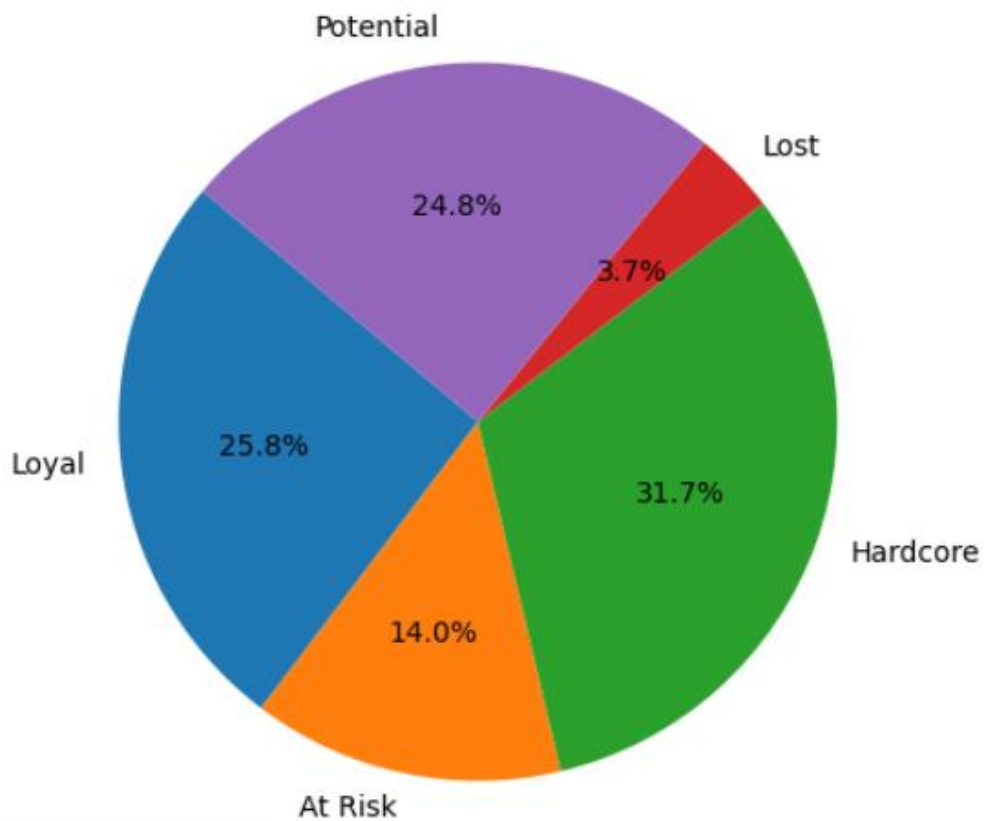
# EVALUATION (REMOVE OUTLIERS & SCALE DATA)

Phương pháp	Remove outliers		Scale data		Nhận xét
	Số cụm	Silhouette Score	Số cụm	Silhouette Score	
KMeans	4	0.42	3	0.36	Khi loại bỏ outliers thì các cụm phân tách tốt
KMeans (PySpark)	6	0.45	5	0.48	Các cụm chồng chéo lên nhau
GMM	4	0.2	6	0.13	Các cụm không được phân tách rõ ràng
DBSCAN	1 + outliers	0.3	1+ outliers	0.45	Không thực hiện phân cụm được
Hierarchical Clustering	6	0.23	4	0.25	Kết quả phân cụm không tốt

06

**SOLUTIONS**

## Tỷ trọng doanh thu theo nhóm khách hàng



Segment	Tỷ lệ	Đặc điểm			Doanh thu	Mục tiêu	Đề xuất phương pháp tiếp cận
		R	F	M			
Hardcore	15	Cao	Cao	Cao	31.7%	Giữ chân, tăng giá trị đơn	Chương trình khách hàng thân thiết cao cấp, dịch vụ cá nhân hóa, cross-sell
Loyal	19	Cao	Cao	TB	25.8%	Tăng giá trị đơn	Tích điểm thành viên, khuyến mãi combo, Ưu đãi sinh nhật & dịp đặc biệt
Potential	23	TB	Cao	TB	24.8%	Tăng tần suất mua	Ưu đãi cá nhân hóa diễn ra trong thời gian ngắn, gửi tin SMS/mail
At Risk	17	Thấp	TB	Thấp	14	Lôi kéo trở lại	Khảo sát và cải thiện dịch vụ, tương tác trên nhiều channel, giảm giá đặc biệt hoặc quà tặng
Lost	12	Thấp	Thấp	Thấp	3.7%	Cân nhắc nguồn lực, có thể bỏ qua nhóm này	Khảo sát lý do rời bỏ, vì nhóm chiếm tỷ lệ khá ít nên có thể bỏ qua



**THANK YOU!**  
**FOR WATCHING**

