# Clustering Algorithms
## Online Retail Case Study

Timo Sachs | AI2 – Clustering Assignment
Barcelona, 31.03.2024

Do Good. Do Better.

esade

# Smart customer segmentation provides opportunity to retain high-value client revenues and (re-)engage dormant customers

**Project Objective**

Optimize marketing and customer retention strategies using clustering algorithms to segment existing customer base and make actionable recommendations to improve campaign effectiveness for groups with similar characteristics and purchasing behaviors.
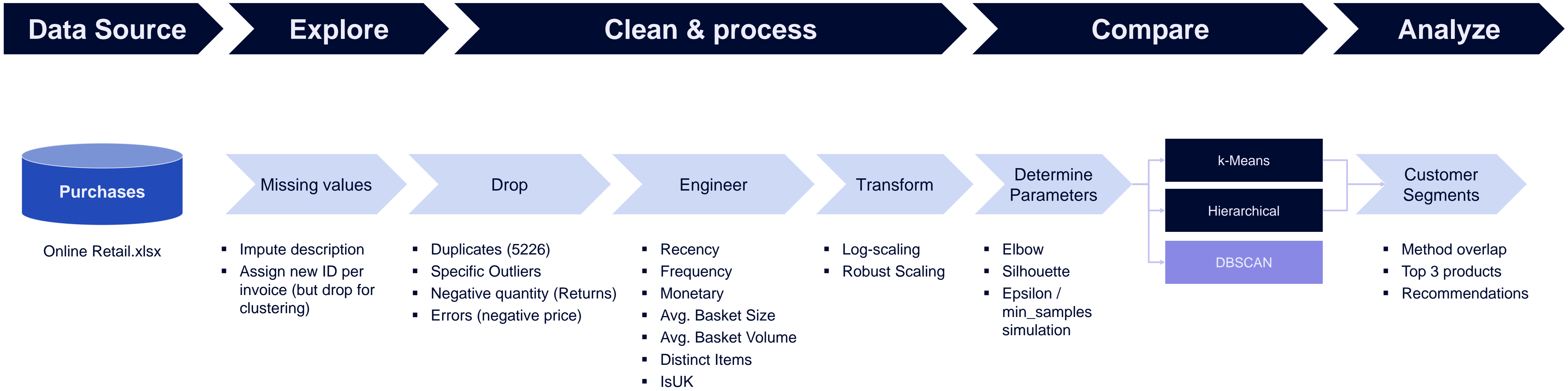
**Customer insights**

- **4 customer segments**: Dormant Majority (68%), Occasional Spenders (24%), VIP Customers (8%), and High-Volume Anomalies (<1%) with specific marketing strategies available
- Top **7.2% of customers** ('VIP Customers' segment) generate **47.4% of total revenue**
- Approximately 34.4% of customers generate 80% of revenue
- Overall >90% of customers from the UK, while profitable segments are more internationally diverse
- Heavy reliance on a single market presents risk and opportunity

**Model findings**

- **K-Means performed slightly better** than hierarchical clustering (Silhouette Score: 0.344 vs. 0.300)
- DBSCAN **not effective** with the given dataset and feature engineering approach (only 1-2 clusters)
- Log-transformation and robust scaling absolutely crucial for meaningful cluster results
- **Remaining anomalies (0.02-0.07%)** indicate potential B2B customers or data issue

esade

# Analysis process is based on **step-by-step processing pipeline** fed into 3 clustering algorithms to derive customer segments

| Data Source | Explore | Clean & process | Compare | Analyze |
|---|---|---|---|---|

**Purchases**

Online Retail.xlsx

**Missing values**
- Impute description
- Assign new ID per invoice (but drop for clustering)

**Drop**
- Duplicates (5226)
- Specific Outliers
- Negative quantity (Returns)
- Errors (negative price)

**Engineer**
- Recency
- Frequency
- Monetary
- Avg. Basket Size
- Avg. Basket Volume
- Distinct Items
- IsUK

**Transform**
- Log-scaling
- Robust Scaling

**Determine Parameters**
- Elbow
- Silhouette
- Epsilon / min_samples simulation

k-Means

Hierarchical

DBSCAN

**Customer Segments**
- Method overlap
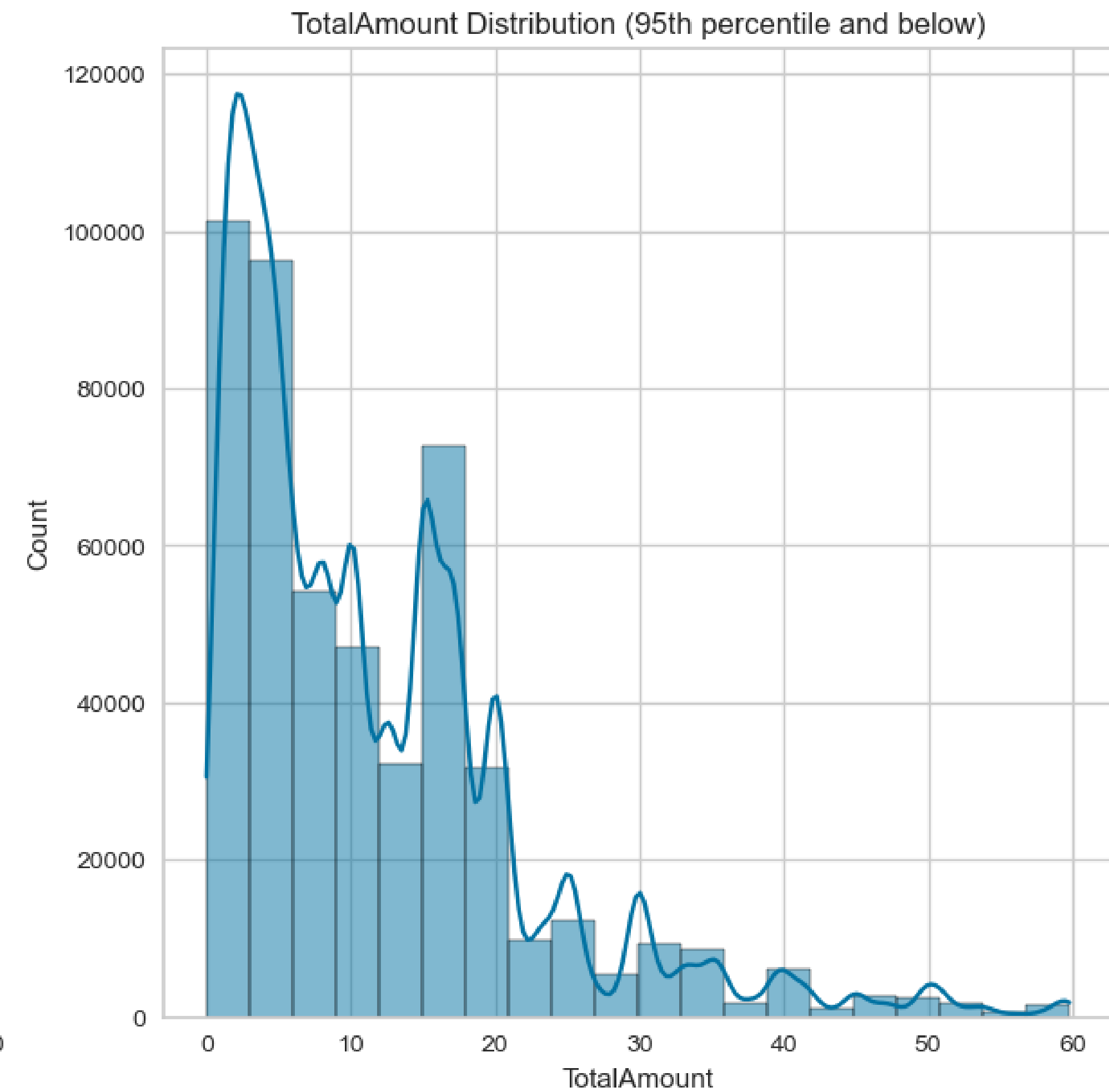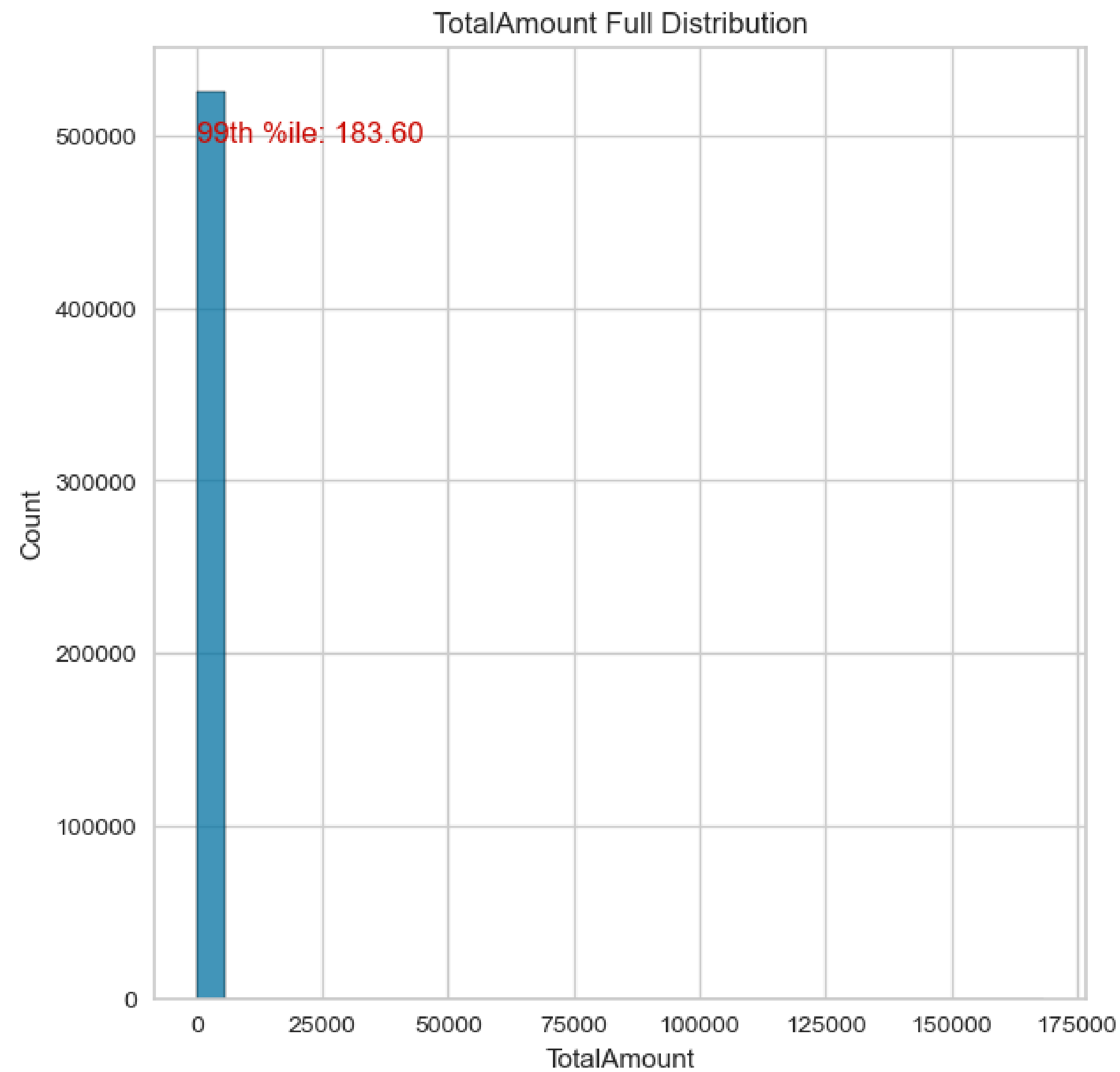- Top 3 products
- Recommendations

# 03

# Feature Cleaning, Selection & Processing

# **Feature selection and preprocessing** is based on a mix of exploration, imputation and dropping features

| | Contains | Missing values | Processing |
|---|---|---|---|
| **InvoiceNo** | Invoice number. Nominal, a 6-digit integral number uniquely assigned to each transaction. | No | Used for CustomerID imputation |
| **StockCode** | Product (item) code. Nominal, a 5-digit integral number uniquely assigned to each distinct product. | No | No |
| **Description** | Product (item) name. | Yes | Yes, impute with "unknown" |
| **Quantity** | The quantities of each product (item) per transaction. Numeric. | No | Calculate total amount |
| **InvoiceDate** | Invoice Date and time. Numeric, the day and time when each transaction was generated. | No | Derive monthly activity |
| **UnitPrice** | Unit price. Numeric, Product price per unit. | No | Calculate total amount |
| **CustomerID** | Customer number. Nominal, a 5-digit integral number uniquely assigned to each customer. | Yes | Group per customer id |
| **Country** | Country name. Nominal, the name of the country where each customer resides. | No | Create IsUK flag |

esade

# Total amounts (derived from quantity and unit price) are highly skewed and thus log-transformed after outlier removal



TotalAmount Full Distribution

99th %ile: 183.60



TotalAmount Distribution (95th percentile and below)

# **Outlier removal** is **not done based on percentiles**, but through in-depth analysis of top/bottom values with **hand-picked removals**
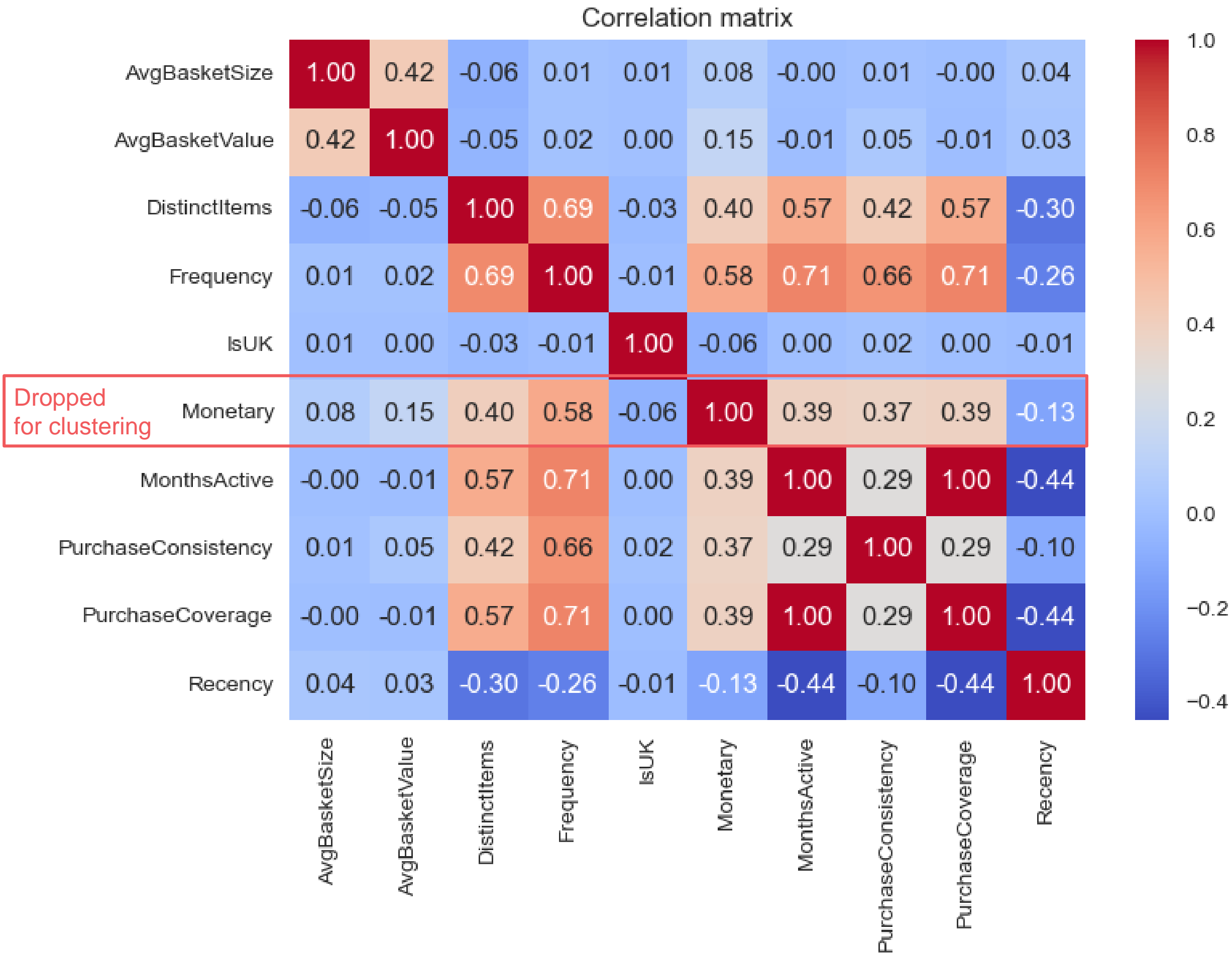
| | Stock Code / Invoice No. | Rows removed | Rationale |
|---|---|---|---|
| **Manual Bookings** | M | 316 | No insights / information can be derived; Might delute results |
| **Postage Fees** | DOT and POST | 706 + 1126 | Postage is not a product to be bought analyzed |
| **Padding** | PADS | 3 | Padding distorts unit prices due to ultra low price; Also not a real product |
| **Bank Charges** | BANK CHARGES | 12 | Bank charges are part of the profit calculation, but are not counted as product / item in this analysis |
| **Wrong Quantity** | 556444 | . | Fixed quantity based on comparable transaction (reset quantity to 1 instead of 60) |
| **Amazon Fee** | AMAZONFEE | 2 | Excluded since not a real product |

esade

# **Common RFM metrics** are enhanced with **additional basket metrics** and seasonality

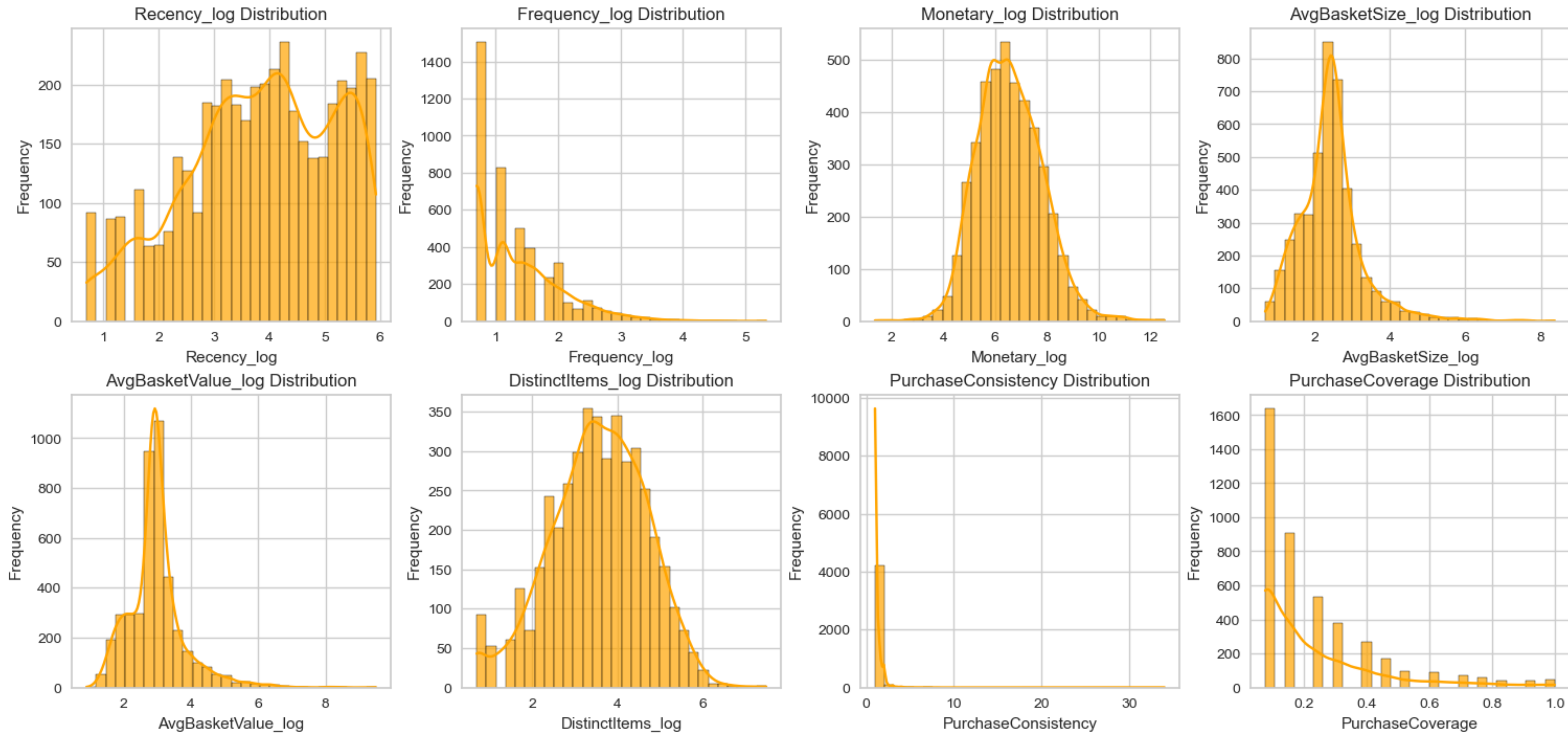| | Metric / Calculation | Insight goal |
|---|---|---|
| **Recency** | Time since last purchase | Active vs. passive customers |
| **Frequency** | Count of purchases | Count of purchases |
| **Monetary** | Sum of total spend per customer | Focus on highest-revenue customers |
| **Avg. Basket size** | Avg. total spend per invoice | Identify bulk buyers vs. small purchases |
| **Distinct items** | Number of unique products bought | Buying behavior (broad vs. specific) |
| **PurchaseConsistency** | Number of purchases within active month | Identify high engagement in active month |
| **Purchase Coverage** | Months with activity / all months | Identify consistent buyers over time |
| **IsUK** | IsUK | Region-specific analysis with less sparsity |

# Newly engineered **feature correlation** already provides **insights into behavioral trends** on global customer base level

- **Frequency** and **Distinct Items** have a strong positive correlation (+0.69), meaning that customers who buy frequently also tend to purchase a wider variety of items.

- **Frequency** and **Months Active** are highly correlated, which makes sense—customers who have been around longer tend to buy more often.

- **Purchase Coverage** and **Frequency** are also strongly linked, suggesting that frequent shoppers engage with multiple product categories.

- **Recency** has a negative correlation with **Months Active**, meaning that longer-term customers tend to have less recent purchases—potential sign of disengagement.

- **Monetary** value is most related to **Frequency**, confirming that the more often someone buys, the more they spend overall.

### Correlation matrix

| | AvgBasketSize | AvgBasketValue | DistinctItems | Frequency | IsUK | Monetary | MonthsActive | PurchaseConsistency | PurchaseCoverage | Recency |
|---|---|---|---|---|---|---|---|---|---|---|
| **AvgBasketSize** | 1.00 | 0.42 | -0.06 | 0.01 | 0.01 | 0.08 | -0.00 | 0.01 | -0.00 | 0.04 |
| **AvgBasketValue** | 0.42 | 1.00 | -0.05 | 0.02 | 0.00 | 0.15 | -0.01 | 0.05 | -0.01 | 0.03 |
| **DistinctItems** | -0.06 | -0.05 | 1.00 | 0.69 | -0.03 | 0.40 | 0.57 | 0.42 | 0.57 | -0.30 |
| **Frequency** | 0.01 | 0.02 | 0.69 | 1.00 | -0.01 | 0.58 | 0.71 | 0.66 | 0.71 | -0.26 |
| **IsUK** | 0.01 | 0.00 | -0.03 | -0.01 | 1.00 | -0.06 | 0.00 | 0.02 | 0.00 | -0.01 |
| **Monetary** | 0.08 | 0.15 | 0.40 | 0.58 | -0.06 | 1.00 | 0.39 | 0.37 | 0.39 | -0.13 |
| **MonthsActive** | -0.00 | -0.01 | 0.57 | 0.71 | 0.00 | 0.39 | 1.00 | 0.29 | 1.00 | -0.44 |
| **PurchaseConsistency** | 0.01 | 0.05 | 0.42 | 0.66 | 0.02 | 0.37 | 0.29 | 1.00 | 0.29 | -0.10 |
| **PurchaseCoverage** | -0.00 | -0.01 | 0.57 | 0.71 | 0.00 | 0.39 | 1.00 | 0.29 | 1.00 | -0.44 |
| **Recency** | 0.04 | 0.03 | -0.30 | -0.26 | -0.01 | -0.13 | -0.44 | -0.10 | -0.44 | 1.00 |

Dropped for clustering (Monetary row)

# Log-transformation leads to normalized distribution of quantitative features which can be fed into the clustering algorithms
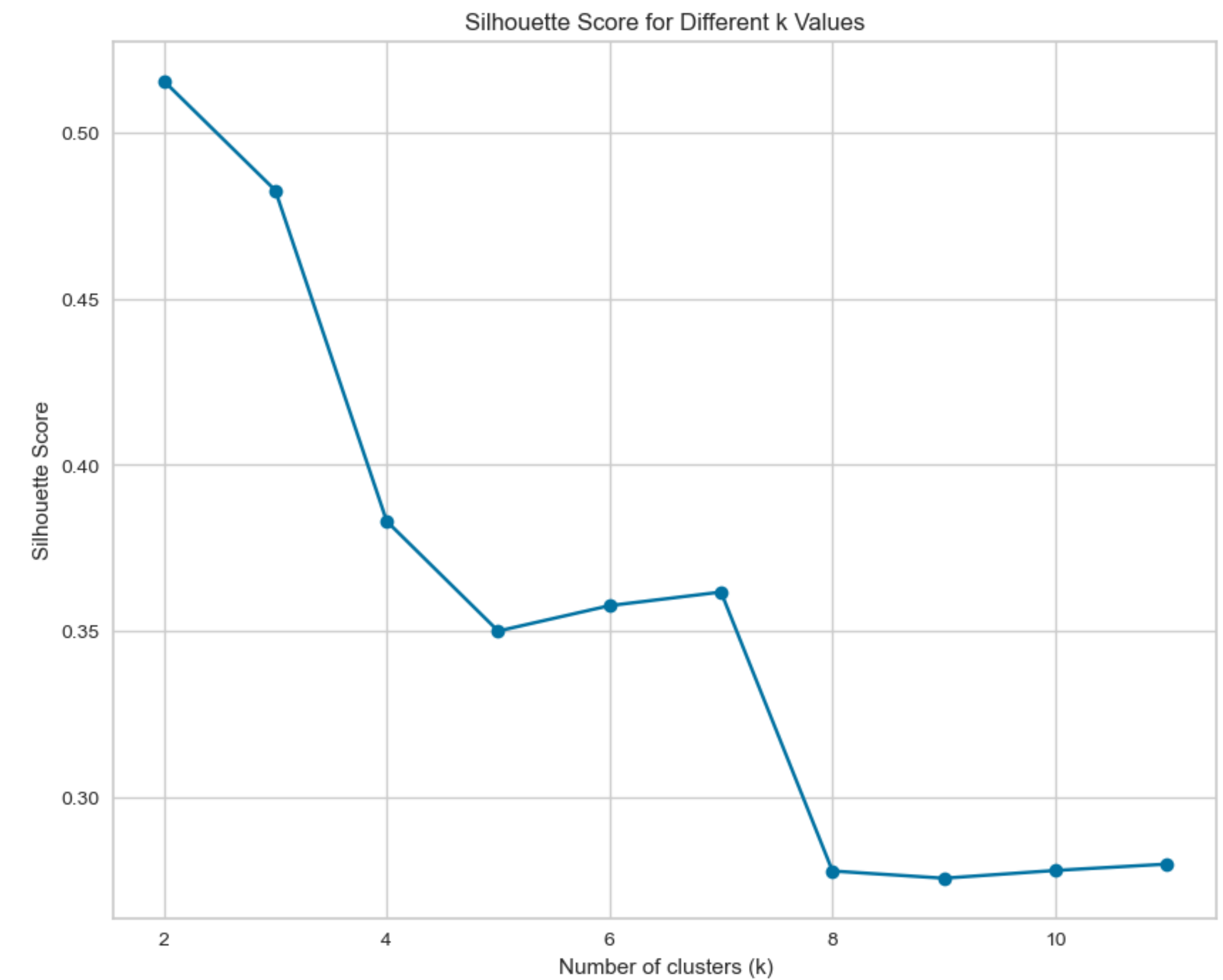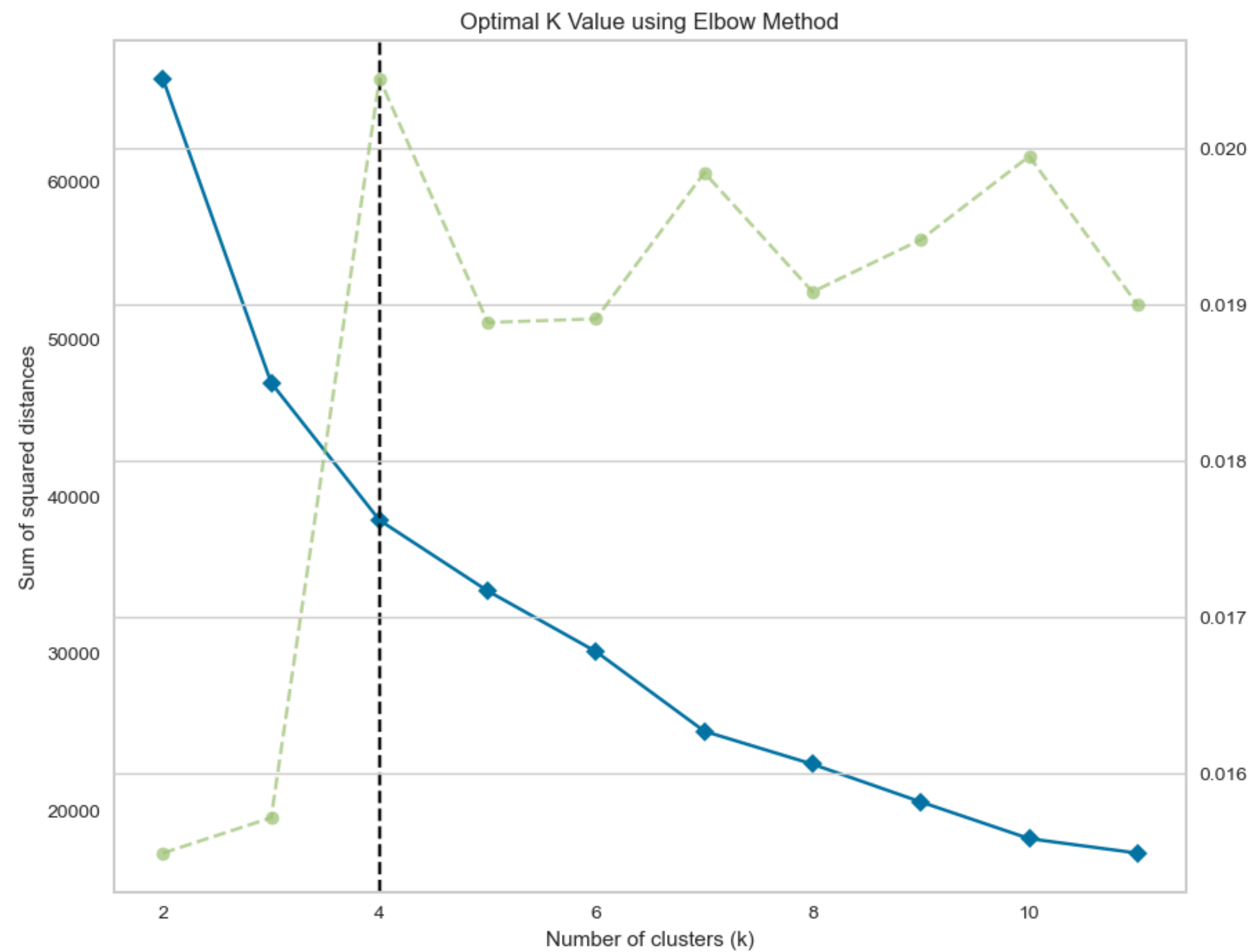
# 04
# Parameter selection

# Elbow method suggests an **optimal number of 4 clusters**; Silhouette score quite low due to complex difficult-to-separate customer groups

esade

# **Dendrogram** reveals clear customer segments despite heavy branching towards lower distance levels; Same cluster cut is applied as for k-Means
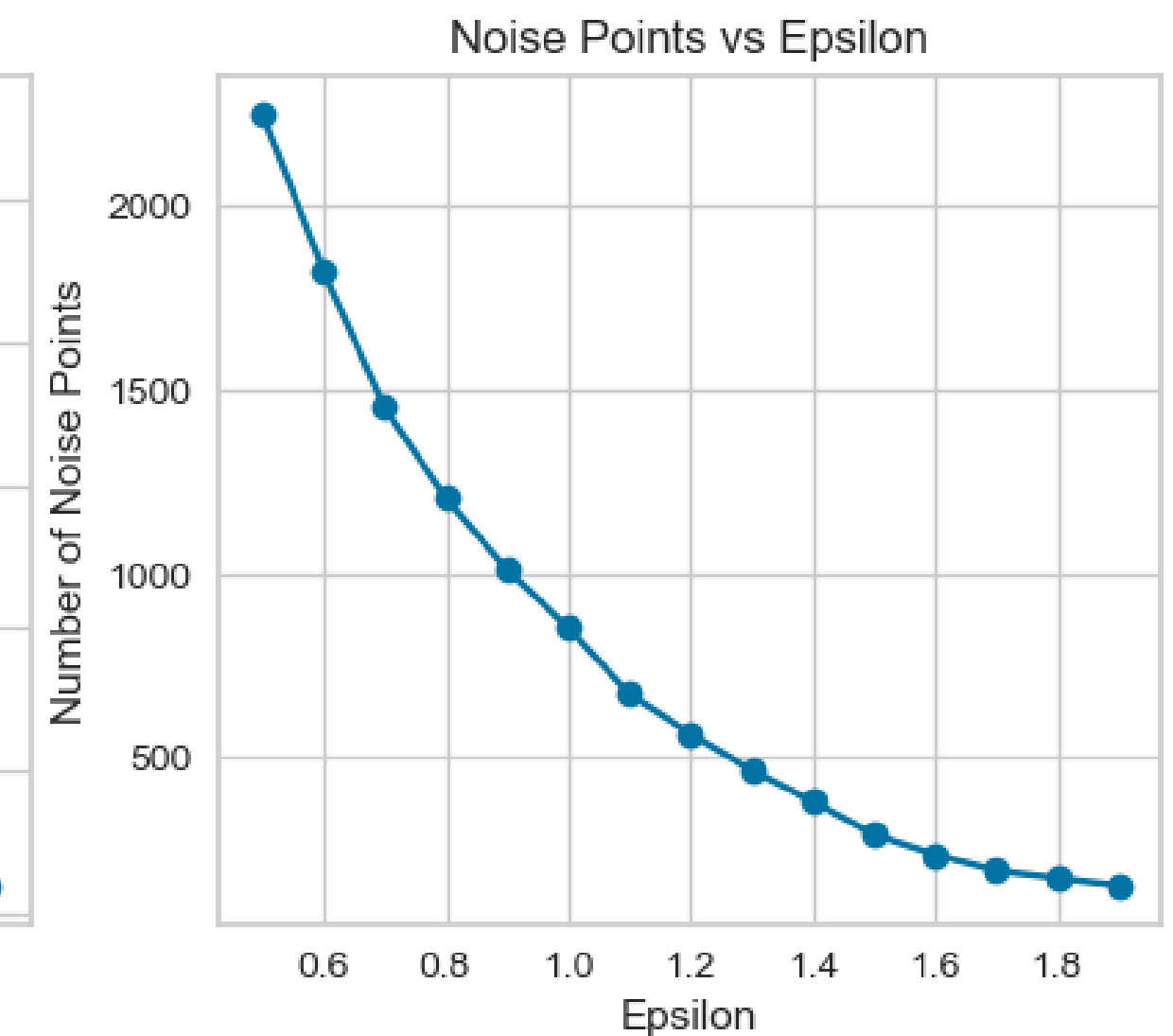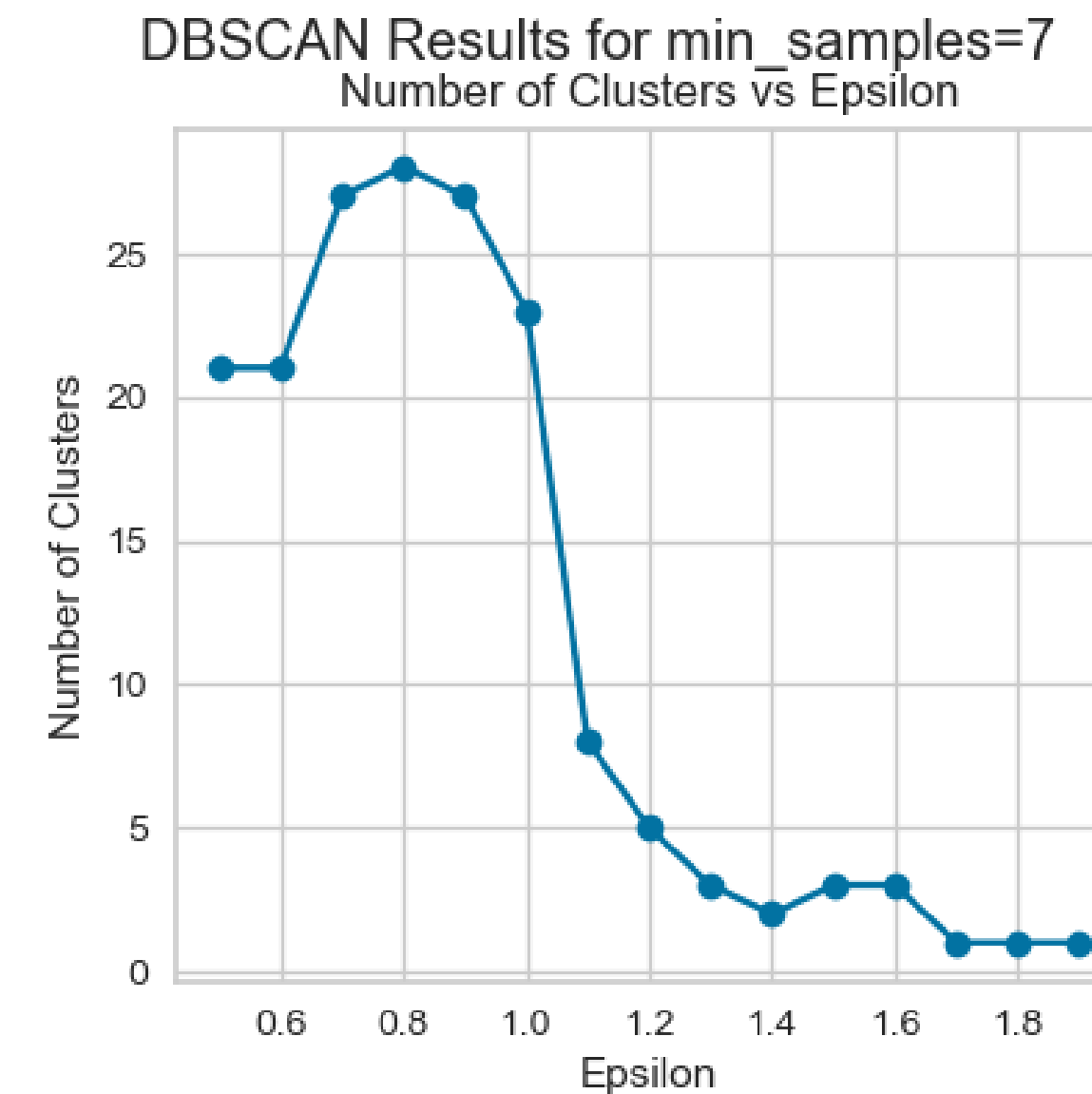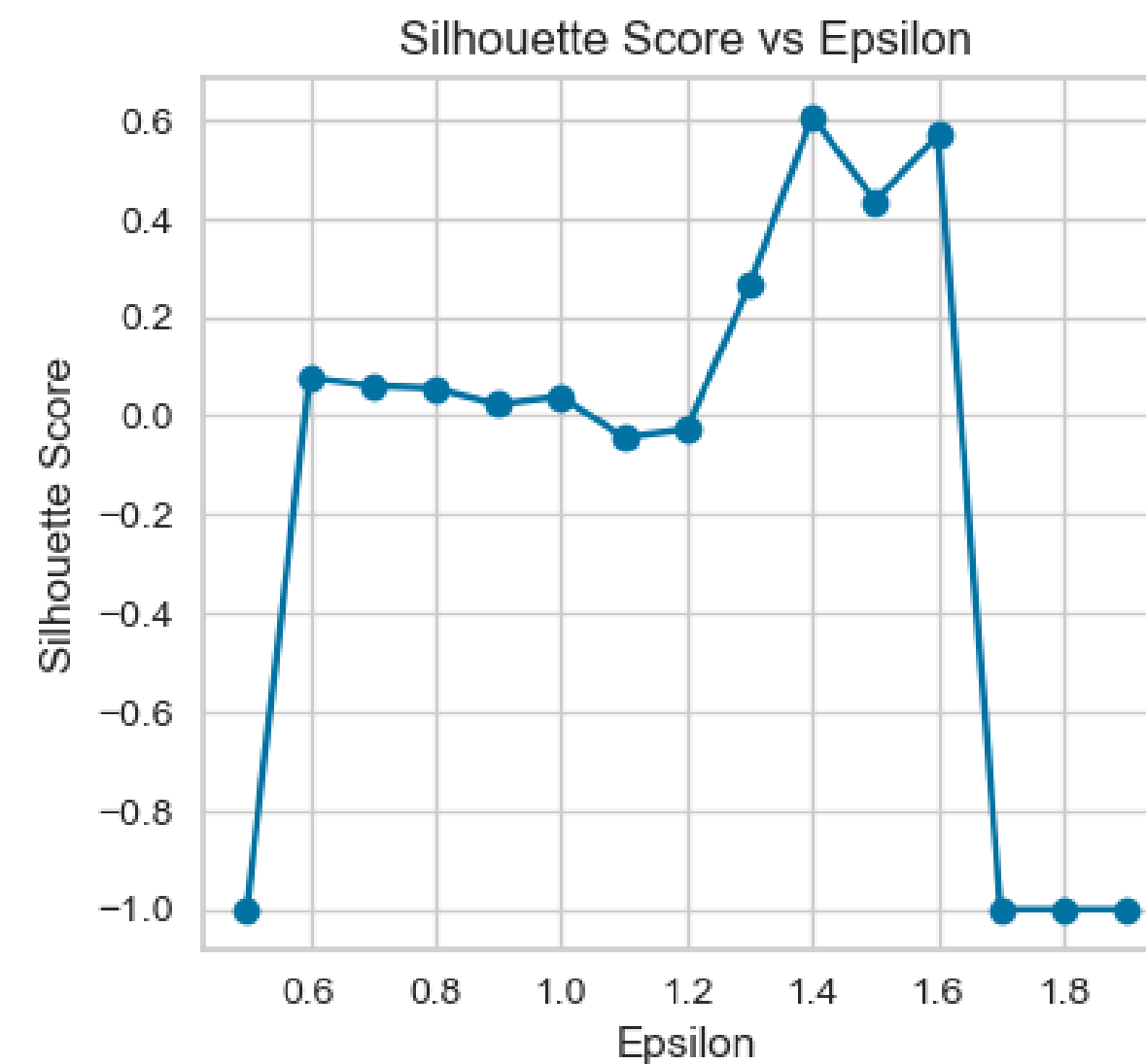


Customer Segmentation Hierarchical Clustering Dendrogram

esade

# Despite multiple parameter simulations, overall DBSCAN performance was insufficient to derive meaningful customer segments

**Best DBSCAN parameters**

- Epsilon: 1.4
- Min Samples: 7.0
- Number of clusters: 2.0
- Noise points: 381.0 (8.79%)
- Silhouette Score: 0.6071



DBSCAN Results for min_samples=7

# 05

# Clustering Results & Model Comparison

esade

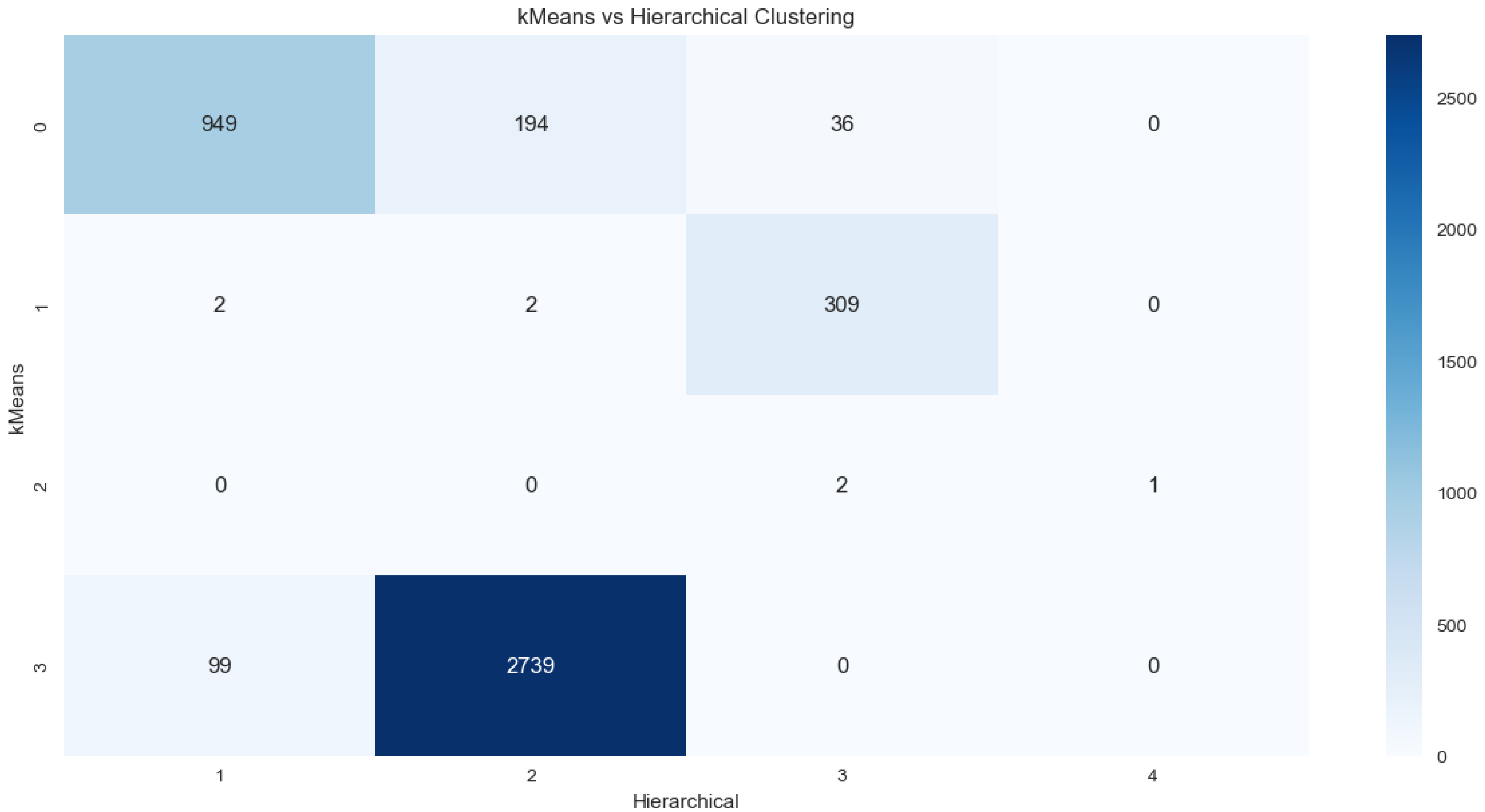# Clustering results show distinct customer segment characteristics with 3 customer groups and 1 outlier / anomaly cluster

| | Recency | Frequency | Monetary | Avg Basket Size | Avg Basket Value | Distinct Items | Purchase Consistency | Purchase Coverage | IsUK | # of Custom. | % of Custom. |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Dormant Majority** | 124.22 | 1.66 | 583.74 | 23.06 | 40.11 | 31.92 | 1.09 | 0.12 | 90.45 | 2838 | **65.50** |
| **Occasional Spenders** | 38.13 | 5.85 | 2238.66 | 18.82 | 27.15 | 98.96 | 1.29 | 0.36 | 90.25 | 1179 | **27.21** |
| **VIP Customers** | 12.06 | 20.27 | 12858.93 | 29.19 | 49.40 | 177.01 | 1.97 | 0.78 | 90.10 | 313 | **7.22** |
| **High-Value Anomalies** | 124.67 | 146.33 | 59126.27 | 8.56 | 16.70 | 1191.00 | 21.72 | 0.69 | 66.67 | 3 | **0.07** |

**Detailed distributions in appendix**

# k-Means and hierarchical clustering produced similar results; DBSCAN struggled to identify customer segments despite parameter tuning

| K-Means | Hierarchical | DBSCAN |
|---|---|---|

**+ Strengths**

**K-Means**
- Well-balanced customer segments (8%, 24%, 68%)
- Clear, actionable segments with distinct RFM profiles
- Good separation between high-/low-value customers
- Results immediately applicable to marketing strategy

**Hierarchical**
- Visualization (dendrogram) of customer relationships
- Flexibility in choosing segment count after analysis
- Less affected by outliers than K-means

**DBSCAN**
- Automatically detected outliers (high 13.5% of customers)
- Did not require pre-specifying number of segments
- Found clusters of varying shapes and densities

**— Weak points**

**K-Means**
- May not capture complex relationships
- Requires choosing k (=4)
- Sensitive to outliers (high-spent customers)

**Hierarchical**
- Computationally more intensive
- Some smaller segments with unclear business value
- More difficult to explain to business stakeholders

**DBSCAN**
- Separate parameter tuning
- Classified many potentially valuable customers as "noise"
- Did not identify meaningful clusters

**Compare results**

**Discard for this analysis**

# **Confusion matrix** of k-means and hierarchical clustering reveal similar customer attribution despite different labeling…



kMeans vs Hierarchical Clustering

esade

# ... which is supported by the overall similar comparison metrics across customer segments

**Robust Segmentation**

The **remarkable consistency between K-means and hierarchical clustering** validates the segmentation approach. Both methods identified similar customer groups with consistent behavioral patterns, increasing the confidence in these segments as a foundation for strategic planning.

**Segment Proportions**

Hierarchical clustering assigned slightly **more customers to the dormant segment** (67.74% vs. 65.50%) and fewer to the middle segment (24.23% vs. 27.21%), suggesting some borderline customers may display characteristics that could place them in either group.

**Different Outlier Detection**

The most significant difference appears in the outlier segment, where **hierarchical clustering (only) identified a single customer** with very different behavioral patterns than the three anomalies found in K-means. This suggests these methods might have different sensitivities to unusual purchasing patterns.
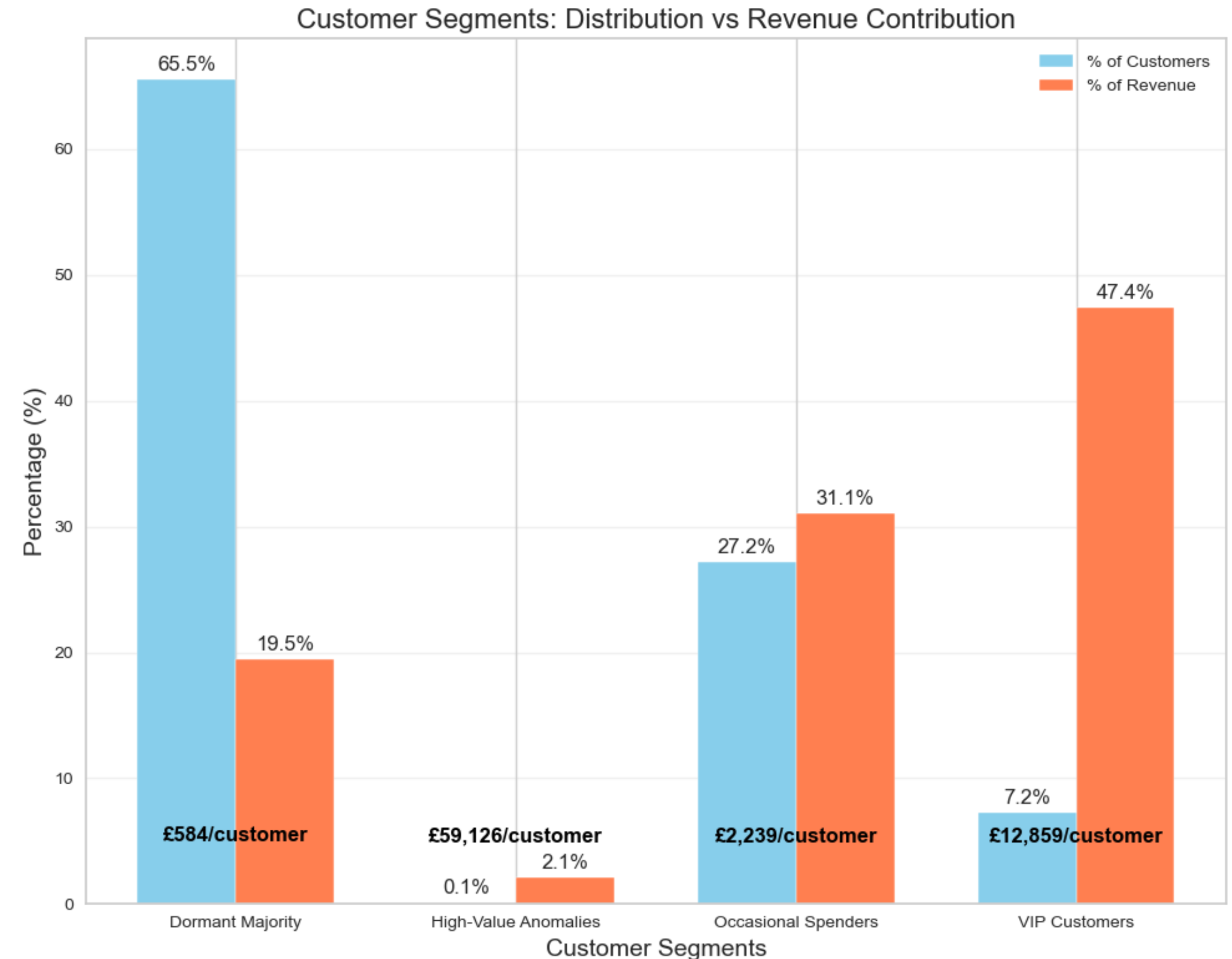
**Value Distribution**

The distribution of value across segments remains consistent, with a small percentage of customers (8.01% in hierarchical clustering) accounting for disproportionately high spending, frequency, and engagement, while the majority (67.74%) remain relatively disengaged with minimal spending.
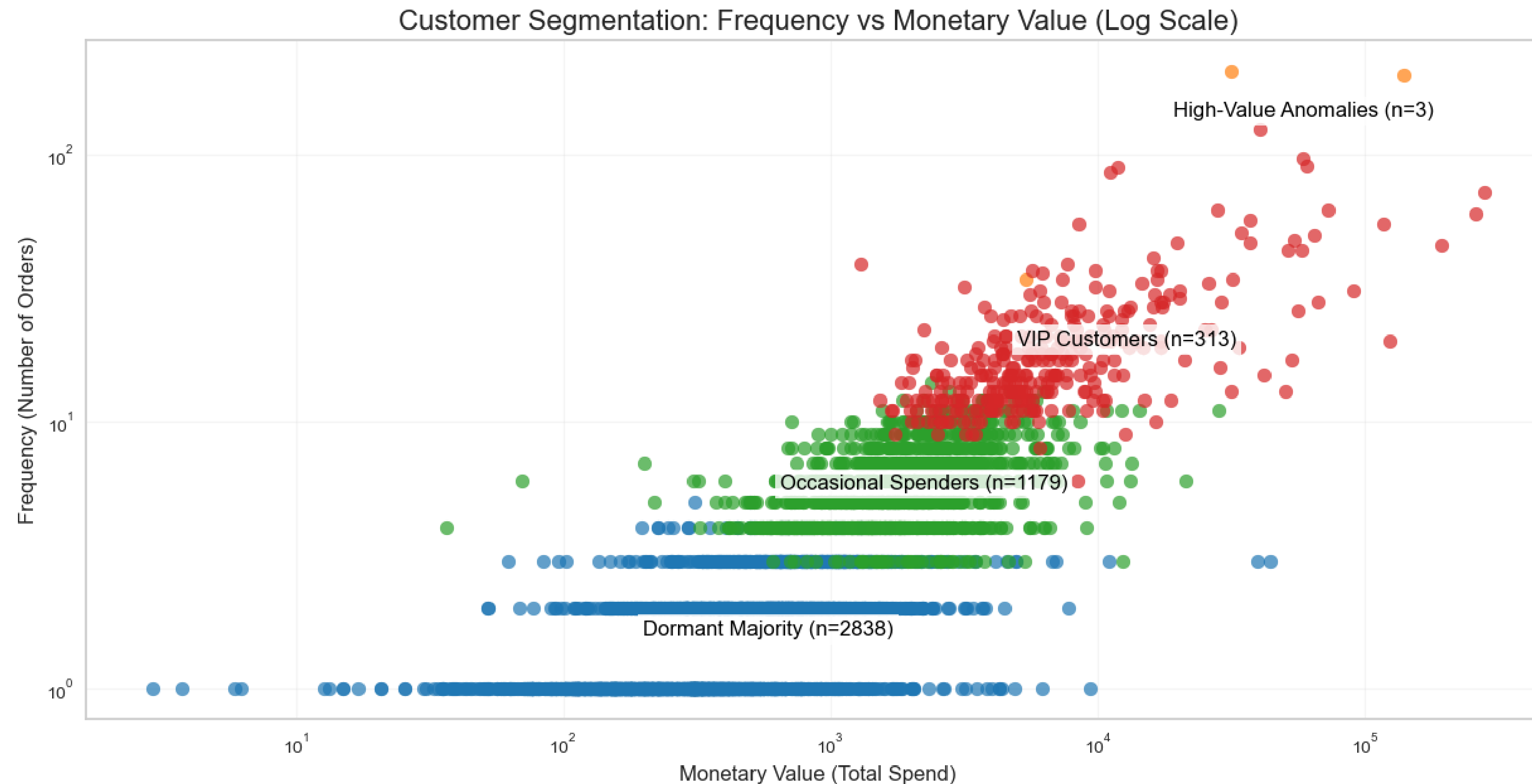
# 06

# Segment Deep-dive and Marketing Strategies

# Revenue contribution by customer segments reveals strong pareto distribution for VIP customers and dormant majority

- **Dormant majority** makes up 65% of customer based but only **19.5% of all revenues**

- **VIP customers** contribute **47.4% of revenues** (~2.5 of dormant but with only ~1/10 of customers demonstrating strong powerlaw distribution and need to retain those customers

- Occasional spenders might tip either way (dormant or VIP)

- High-value anomalies need to be revisited to confirm outlier removal approach or identify fraud / big B2B buyers



Customer Segments: Distribution vs Revenue Contribution

% of Customers
% of Revenue

Dormant Majority: 65.5% / 19.5% — £584/customer
High-Value Anomalies: 0.1% / 2.1% — £59,126/customer
Occasional Spenders: 27.2% / 31.1% — £2,239/customer
VIP Customers: 7.2% / 47.4% — £12,859/customer

# Customer segment clusters show relationship between number of orders and total spend, whereas broader base is separated from outliers



Customer Segmentation: Frequency vs Monetary Value (Log Scale)

High-Value Anomalies (n=3)

VIP Customers (n=313)

Occasional Spenders (n=1179)

Dormant Majority (n=2838)

Frequency (Number of Orders)

Monetary Value (Total Spend)

# The **dormant majority** should be targeted using **reactivation campaigns** to turn them into returning customers

| Segment | Dormant Majority \| 66 % |
|---|---|
| **Key Characteristics** | ▪ 124 days since last purchase<br>▪ Only 1.7 purchases on average<br>▪ $584 total spend per customer |
| **Top 3 Products** | ▪ World War 2 Gliders<br>▪ White Hanging Heart<br>▪ Fairy Cake Flannel |
| **Opportunity** | Reactivation could unlock significant revenue |

**Marketing Recommendation**

- **Reactivation Campaign:** "We Miss You" emails with personalized product recommendations

- **Win-Back Incentives:** First-time reorder discount or free shipping

- **Reminder Strategy:** Show previously purchased items with complementary suggestions

- **Seasonal Triggers:** Reach out during key seasonal moments based on past purchase timing

# Increase average spending of **occasional spenders** to turn them into VIP customers

| Segment | **Occasional Spenders** | 27 % |
|---|---|

| Key Characteristics | ▪ More recent activity (38 days)<br>▪ Moderate frequency (5.8 purchases)<br>▪ $2,239 lifetime value<br>▪ 98 distinct items vs. 31 for dormant |
|---|---|
| Top 3 Products | ▪ World War 2 Gliders<br>▪ Pack of 72 Retrospot<br>▪ Small Popcorn Holder |
| Opportunity | Most likely to convert to VIP status with right engagement |

## Marketing Recommendation

▪ **Frequency Program:** Reward increased purchase frequency with escalating benefits

▪ **Category Expansion:** Introduce related product categories based on past purchases

▪ **Mid-tier Loyalty:** Create "rising star" tier with visible pathway to premium benefits

▪ **Personalized Bundles:** Offer curated collections based on previous basket analysis

# Maintain **high engagement of VIP clients** by offering premium experience and exclusive benefits

| Segment | VIP Customers | 8 % |
|---|---|

| Key Characteristics | <ul><li>Very recent activity (12 days)</li><li>High frequency (20.3 purchases)</li><li>$12,858 lifetime value (22× dormant)</li><li>Highest basket value ($49.4)</li></ul> |
|---|---|
| Top 3 Products | <ul><li>Jumbo Bag Retros</li><li>Small Popcorn Holder</li><li>World War 2 Gliders</li></ul> |
| Opportunity | **Retention is critical** as they drive significant revenue |

## Marketing Recommendation

- **Premium Experience:** White-glove customer service and early access to new products

- **Relationship Building:** Personal shopping assistance and product customization

- **Exclusive Benefits:** VIP-only events and substantial loyalty rewards

- **Referral Program:** Incentivize bringing in similar high-value customers

esade

# Based on the cluster analysis, clear **next steps** can be derived to improve customer retention and increase overall revenues

**Customer management**

- **Develop VIP retention program:** Early product access, personalized service, exclusive benefits
- **Create conversion path for Occasional Spenders:** Incentivize increased purchase frequency
- **Design reactivation campaigns:** Target dormant customers with personalized recommendations

**Analytics Enhancement**

- Implement **real-time segment scoring**: Dynamically assign customers as behaviors change
- Conduct A/B testing: Compare segment-specific marketing approaches against generic campaigns
- Longitudinal analysis: Track customer movement between segments over time
- **Integrate customer service data:** Enhance segmentation with support interactions and satisfaction ratings

**Timeline**

- **Month 1:** Set up segment-based targeting in marketing platforms
- **Month 2-3:** Launch pilot campaigns for each segment with distinct messaging
- **Month 3-6:** Measure segment-specific conversion rates and adjust strategies
- **Month 6+:** Implement predictive analytics to identify at-risk customers and growth opportunities
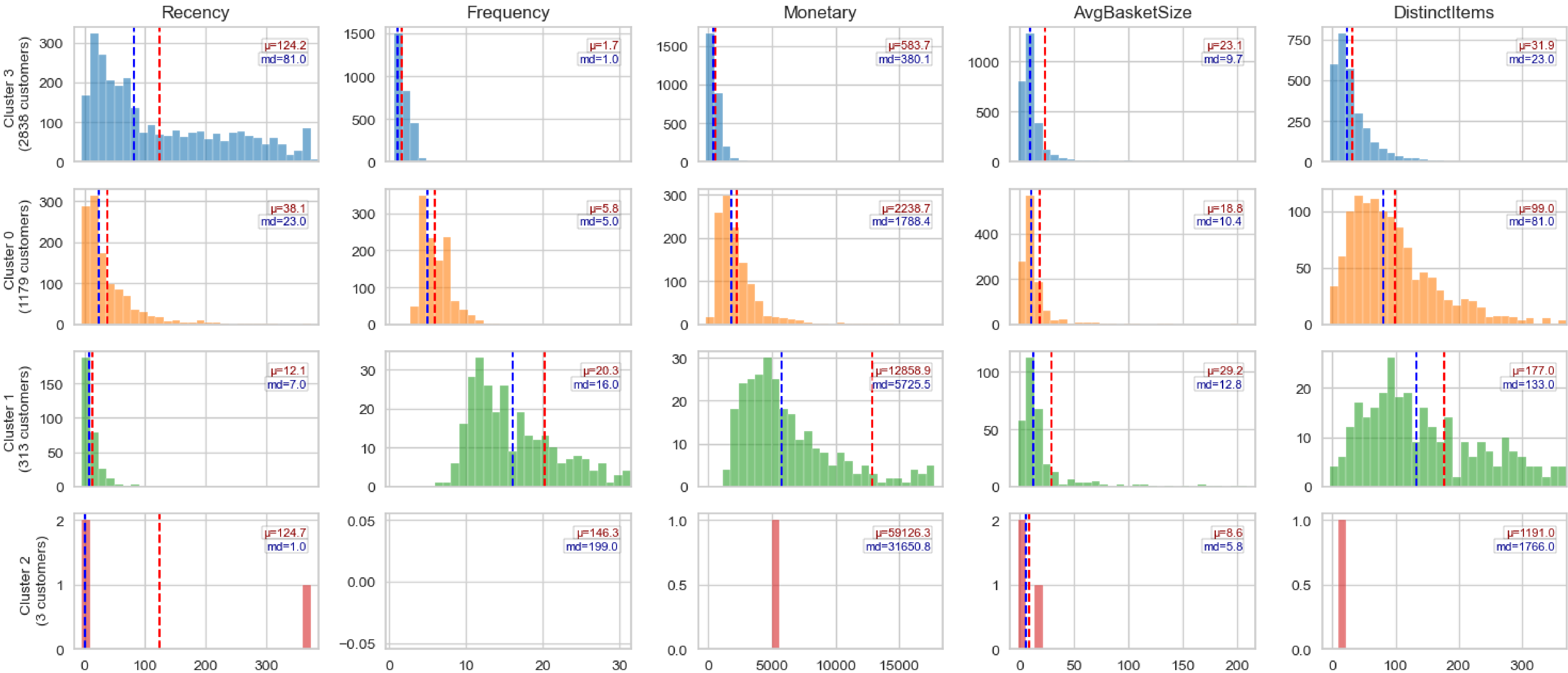
# Appendix

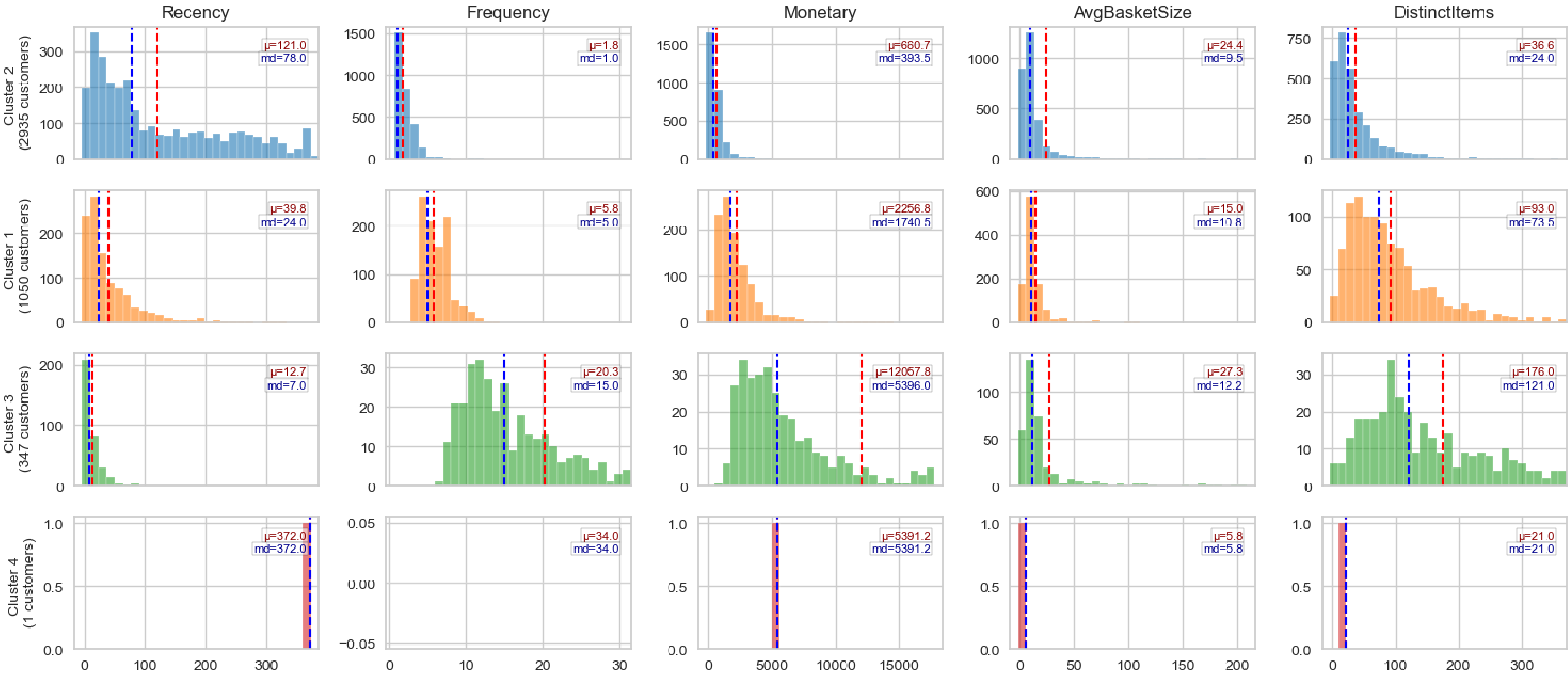# esade

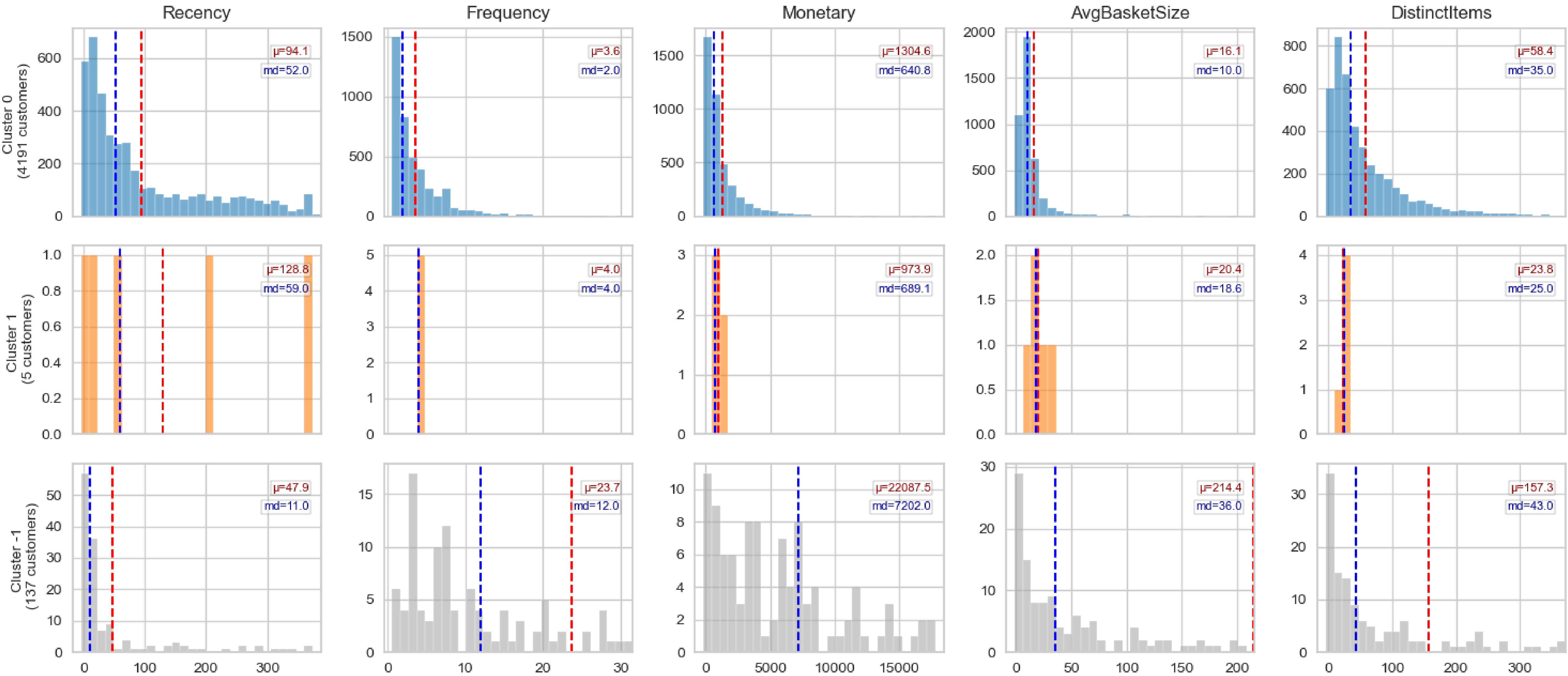## Feature Distributions by kM_Cluster (Sorted by Size)

# esade



Feature Distributions by Hierarchical_Cluster (Sorted by Size)

# esade

## Feature Distributions by DBSCAN_Cluster (Sorted by Size)

esade

# Analysis of top 3 product reveals commonalities across customer segments as well as deviations for VIP and anomalies