

Advanced Customer Segmentation Project: A Synthesis of Key Findings and Strategies

Executive Summary

This briefing document synthesizes the findings of a comprehensive customer segmentation project that utilized RFM (Recency, Frequency, Monetary) analysis and HDBSCAN machine

learning to transform an undifferentiated e-commerce customer base into actionable, data-driven segments. The project successfully processed over 541,000 transaction records to identify seven distinct customer groups and a high-value outlier segment, enabling a shift

from inefficient, one-size-fits-all marketing to precise, high-ROI targeted strategies.

The model revealed that **80% of total revenue is generated by the top two segments**

(**"Champions"** and **"Potential Loyalists"**) , which constitute just 51% of the customer

base. Furthermore, a micro-segment of "VIP Whales," representing only 1% of customers,

accounts for a staggering 10% of all revenue. Key operational insights include the **"51-Day**

Rule," a critical recency threshold for customer intervention, and the **"November Effect,"**

a seasonal peak where transaction volume triples, demanding significant logistical and

financial planning. The analysis also established that purchase **Frequency is the single**

strongest predictor of customer monetary value , with a 0.81 correlation. Based on these

findings, a detailed strategic framework has been developed for each segment, projecting a

2.5:1 return on investment (ROI) from an annual budget of £1.76 million, expected to

generate a revenue impact of over \$4.33 million. The project is production-ready, supported

.by a complete MLOps lifecycle using MLflow for ongoing monitoring and retraining

Project Foundation and Methodology .1

Business Problem and Objectives .1.1

The project addressed the critical business challenge of operating without clear customer segmentation, which led to inefficient marketing, wasted resources, and missed revenue opportunities. The primary objective was to answer the question: *"Can we segment our*

customers into meaningful groups based on their actual behavior, so we can make

data-driven decisions about resource allocation, marketing strategies, and retention efforts

for each segment?"

The mission was to evolve from a generic customer approach to a personalized management system through deep behavioral analysis and unsupervised

.machine learning, creating actionable segments to maximize ROI

Technical Methodology .1.2

The project was executed in four phases using a combination of RFM analysis and advanced machine learning:| Phase | Activities || ----- | ----- || **1. Data Preparation** |

Cleaning and validating 541,909 raw transaction records, engineering RFM metrics, and normalizing distributions with log transformation. || **2. Exploratory Analysis** | Performing

correlation analysis, identifying behavioral patterns, and discovering seasonal and geographic trends. || **3. Machine Learning** | Applying HDBSCAN clustering to the prepared RFM data to automatically discover natural customer segments. || **4. Business Intelligence**

| Profiling and naming the discovered segments, developing strategic recommendations, and

| .projecting the financial impact

Data Processing and Model Development .2

Data Cleaning and Feature Engineering .2.1

The initial dataset of 541,909 transaction records underwent a rigorous cleaning pipeline. This involved removing cancelled orders, negative quantities, duplicate entries, and records with missing CustomerIDs. This process resulted in a final clean dataset of approximately **400,000 valid transactions from ~4,000 unique customers**, representing an acceptable

:25% data loss rate.RFM features were then engineered for each customer

.**Recency (R)**: Days since the last purchase •

.**Frequency (F)**: Total number of unique purchases •

.**Monetary (M)**: Total amount spent •

Distribution Normalization .2.2

Initial analysis revealed extreme right skewness in the Frequency (12.07) and Monetary (19.34) distributions, which violates the assumptions of distance-based clustering algorithms

and can lead to unreliable results. To address this, a **log transformation** was

applied.**Results After Log Transformation:**| Metric | Skewness Before | Skewness After | Status || ----- | ----- | ----- | ----- || **Recency** | 1.25 | -0.38 | Normalized || **Frequency** | 12.07 | 1.21 | Excellent || **Monetary** | 19.34 | 0.40 | Perfect

This crucial step successfully normalized the distributions while preserving the relative relationships between customers, ensuring the subsequent clustering was accurate and reliable. Outliers, identified as legitimate high-value "whale customers," were intentionally .kept as they represent strategic assets

Core Behavioral Insights .3

The analysis uncovered six critical insights into customer behavior that form the foundation .for strategic planning

Insight 1: The Power of Frequency

A correlation analysis revealed that purchase **frequency has a very strong positive correlation of 0.81 with monetary value** . This indicates that frequency is the most powerful predictor of revenue. Customers who purchase more often are overwhelmingly the ones who spend the most over their lifetime. This finding underscores that the primary business goal should be to build purchasing habits and increase repeat business, as this .yields a higher ROI than focusing solely on new customer acquisition

Insight 2: The 51-Day Recency Rule

The median recency across all customers is **51 days** . This means half of the customer base has not made a purchase in over 51 days, marking a critical threshold. The gap between the median (51 days) and the mean (92 days) highlights a long tail of dormant customers, indicating a significant retention problem. This "51-Day Rule" serves as a key .intervention point for "At Risk" customer segments

Insight 3: Two Types of Whales - The VIP Paradox

:Analysis of the top 10 customers revealed a critical distinction among high-value spenders

- **Type 1: Champions:** Customers with high monetary value, high frequency, and low recency (e.g., Customer 14646 with £280k spend over 201 orders). They are deeply loyal and provide a consistent revenue stream
- **Type 2: Lost Whales:** Customers with high monetary value from a single purchase but extremely high recency (e.g., Customer 12346 spent £77k in one order and then disappeared for 326 days). This represents a catastrophic failure in post-purchase engagement
- **Type 3: Potential B2B:** Customers with very high monetary value but very low frequency (e.g., £168k in only 2 purchases), suggesting a wholesale or reseller relationship

Insight 4: Seasonal Explosion - The November Effect

The business experiences extreme seasonality, with Q4 (September-December) accounting for **41% of the entire year's transactions**. November is the peak month, showing **3.15 times the activity of January**, driven by holiday shopping. This "November Effect" necessitates massive operational scaling for inventory, staffing, and logistics. For example, warehouse space requirements jump from a 10,000 sq ft baseline to 32,000 sq ft, and daily shipments scale from ~700 to ~2,100

Insight 5: Geographic Gold Mines

- :International revenue is heavily concentrated in two key markets
- **Ireland (EIRE):** Contributes an estimated **40-45%** of international revenue. The data (average frequency of ~148 purchases) strongly suggests these are B2B customers (retailers, wholesalers) rather than individual consumers
- **Netherlands:** Contributes an estimated **25-30%** of international revenue. This over-dependence (65-75% of international revenue from two markets) presents a critical business risk that requires dedicated logistics, account management, and a diversification strategy

Insight 6: Product Analysis - The RETROSPOT Phenomenon

Analysis of products purchased by "whale" customers shows that a specific design pattern, "**RETROSPOT**," is a powerful driver of loyalty. Three of the top 10 products feature this design. Whales are not buying single expensive items but rather many affordable items frequently. The likely customer profile includes small businesses like cafés, gift shops, and resellers. This insight points to strategic opportunities in product line expansion, bundling, and inventory management for these key items

Machine Learning Segmentation Results .4

HDBSCAN Model Performance .4.1

The HDBSCAN algorithm was chosen over traditional methods like K-Means because it automatically discovers the natural number of clusters, can identify non-spherical cluster shapes, and effectively handles outliers (noise). The model demonstrated strong performance, confirming its suitability for this business problem.

Metric	Score
Interpretation	-----
DBCV Score	0.0385
Silhouette Score	0.0936

| Good cluster separation and internal cohesion. || Reasonable segment boundaries. ||

Stability | 1.3347 | Exceptional; segments are robust and not fragile artifacts. || **Noise %** |
| .8.87% | Excellent (<10%); identifies genuine outliers for manual handling
Feature importance analysis confirmed that **Recency** is the primary driver of cluster
.separation, followed by Frequency and Monetary

Visual Confirmation: The Snake Plot .4.2

The RFM Snake Plot visually illustrates the distinct profiles of each discovered cluster. It plots the average scaled value for Recency, Frequency, and Monetary for each segment, allowing for quick comparison. The plot clearly shows the "crossing" patterns that define each group—for example, Cluster 5 (Champions) shows low Recency with high Frequency and Monetary, while Cluster 0 (Lost) exhibits the opposite pattern. This visualization confirms the data has a natural, non-random group structure that the HDBSCAN model successfully .captured

Customer Segment Profiles .5

The model discovered seven distinct customer segments plus a critical outlier group, each with unique behavioral characteristics and business value.| Cluster | Segment Name | Size (% of Base) | Recency | Frequency | Monetary | Total Revenue | % of Total Revenue || ----- | ----- | ----- | ----- | ----- | ----- | ----- || 5 | 🏆 **Champions** | 21% | ~23 days | ~9 orders | \$3,133 | \$2.66M | 45% || 4 | ☀️ **Potential Loyalists** | 30% | 37 days | 5 orders | \$1,761 | \$2.11M | 35% || 3 | ⚠️ **At Risk** | 15% | 52 days | 4 orders | \$1,405 | \$843K | 5% || 2 | 😴 **Hibernating** | 12% | 70 days | 3 orders | \$1,025 | \$512K | 3% || 1 | 💐 **New / Low Value** | 15% | 94 days | 2 orders | \$664 | \$398K | 1.5% || 0 | 💔 **Lost / Inactive** | 6% | ~157 days | 1 order | \$357 | \$89K | 0.5% || -1 | 💦 **Outliers (VIP Whales)** | 1% | 43 days | ~16 | orders | \$11,972 | \$599K | 10%

Strategic Recommendations and Business Impact .6

A detailed strategic framework was developed for each segment, with a focus on allocating .resources where they will generate the highest return

Summary of Segment Strategies
Segment, Objective, Key Strategies
Champions, Protect and Reward, "VIP treatment, exclusive benefits, loyalty escalators, 🏆 .dedicated account management
Potential Loyalists, Convert to Champions, "Personalized recommendations, subscription ☀️ .programs, educational content, threshold rewards
At Risk, Immediate Win-Back, "Automated early warning system (triggered at 45 days), ⚠️ .multi-phase ""We Miss You"" campaigns, personal outreach
Hibernating, Aggressive Win-Back, "A strong, one-time ""Last Chance"" mega-offer to 😴 .reactivate, followed by archiving if unresponsive
New / Low Value, Educate and Engage, "Robust onboarding email sequences, 💐 .first-purchase follow-ups, value demonstration through content
Lost / Inactive, Minimal Investment, "One final low-cost attempt, followed by database 💔 .cleanup to eliminate marketing waste
VIP Outliers, White-Glove Treatment, "Manual, individual investigation; B2B terms, 💦 .bespoke solutions, and executive-level relationship management

Projected Financial Impact

The implementation of these segmented strategies is projected to deliver significant financial returns. | Metric | Projected Value | ----- | ----- || **Total Annual Investment** | £1.76M || **Total Annual Revenue Impact** | +\$4.33M || **Net Profit Impact** (40% margin) | +\$1.73M || **Overall Program ROI** | 2.5:1 || **Payback Period** | < 5 months

Production Readiness and Next Steps .7

The project is complete and ready for production deployment. The entire MLOps lifecycle has been managed using **MLflow** for experiment tracking, model registry, and versioning.

An API endpoint has been built with **FastAPI** to serve real-time segment predictions. **Immediate next steps** involve launching the foundational programs for the "Champions" and "At Risk" segments, deploying the technology stack, and establishing baseline metrics. The long-term roadmap includes developing a predictive churn model, a next-best-action engine, and integrating omnichannel data to further enhance personalization and strategic precision.