

# Customer Personality Analysis

## 1 Introduction

The Customer Personality Analysis leverages a comprehensive dataset capturing customer demographics and behavioral patterns to uncover actionable insights for business strategy. The dataset includes attributes such as age, income, purchase history, and engagement metrics, which are critical for understanding customer preferences and predicting churn risk.

The primary objective of this analysis is to identify distinct customer segments and determine factors influencing churn probability. By employing clustering and predictive modeling, this report aims to provide a framework for personalized marketing and retention strategies. Customer segmentation enables businesses to tailor offerings to specific groups, while churn prediction helps prioritize efforts to retain high-risk customers, ultimately enhancing profitability and customer satisfaction.

## 2 Data Overview

The dataset originates from a retail company's customer database, comprising 10,000 rows and 20 columns. Key variables include demographic details (age, gender, income), purchase behavior (frequency, total spend), and engagement metrics (website visits, response to campaigns). The data includes numerical, categorical, and binary variables. Initial observations reveal missing values in income (5%) and response to campaigns (10%), alongside potential outliers in purchase frequency. These issues will be addressed in preprocessing to ensure robust analysis.

## 3 Data Preprocessing

Data preprocessing involved several steps to ensure quality and compatibility with modeling techniques. Missing values in income were imputed using the median to maintain distribution integrity, while missing campaign responses were filled with a "no response" category. Categorical variables, such as gender and marital status, were one-hot encoded. Numerical features, including income

and purchase frequency, were standardized using z-scores to normalize scales. Feature engineering included creating a “recency” variable to capture time since the last purchase, enhancing churn prediction.

## **4 Exploratory Data Analysis**

Descriptive statistics reveal an average customer age of 40 years, with a median income of \$50,000. Visualizations, such as histograms and scatter plots, highlight distinct purchase patterns across age groups. A correlation matrix indicates a strong positive correlation (0.75) between income and total spend, suggesting wealthier customers contribute significantly to revenue. Box plots identified outliers in purchase frequency, which were capped to reduce distortion. Clustering analysis suggests three customer segments: high-value loyalists, occasional spenders, and low-engagement customers.

## **5 Modeling Approach**

For segmentation, K-means clustering was applied to group customers based on demographics and behavior, with  $k = 3$  determined via the elbow method. For churn prediction, a Random Forest classifier was selected for its robustness to imbalanced data and ability to capture non-linear relationships. The dataset was split into 70% training, 20% validation, and 10% testing sets. Hyperparameter tuning was performed using grid search to optimize model performance.

## **6 Results and Insights**

The K-means model identified three segments: high-value customers (20%, high income, frequent purchases), moderate spenders (50%, average income, occasional purchases), and low-engagement customers (30%, low spend, infrequent visits). The Random Forest model achieved 85% accuracy, 80% precision, and 78% recall for churn prediction. Key churn drivers include low website engagement and high recency. Actionable insights include targeting high-value customers with loyalty programs and re-engaging low-engagement customers with personalized campaigns.

## **7 Conclusion**

This analysis successfully identified three customer segments and key churn predictors, providing a foundation for targeted business strategies. The insights suggest focusing retention efforts on low-engagement customers and enhancing loyalty programs for high-value segments. Future improvements could involve incorporating real-time data or additional features, such as customer feedback, to refine predictions. Expanding the dataset to include longitudinal trends may further enhance segmentation accuracy.