

Customer Personality Analysis

Saba Madadi

March 2025

Abstract

In customer personality analysis, different factors such as customer details, business offers, and product types can have an impact. By examining how various features affect customer behavior, it is possible to improve sales strategies and increase business profits. In this study, after analyzing the dataset and taking necessary steps, important insights were gained regarding how these features influence sales. These findings can be useful for enhancing performance.

I. Introduction

Statistical tests were developed in early 20th as essential tools for analyzing data and extracting significant insights. They help researchers to make generalizations from smaller samples and evaluate their hypotheses. Customer personality analysis using these tests can help business owners improve their performance and increase profits. Analyzing a dataset containing customer information is crucial for boosting sales, gaining profits, and staying ahead of competitors. In this study, tests such as the Chi-square test, which examines the impact of two categorical variables, the Kruskal-Wallis H test for assessing the relationship between categorical and numerical variables, and the Friedman Test for analyzing multiple numerical variables against categorical ones, have been used. By reviewing the results of these tests, the report will provide insights into the effectiveness of campaigns created by the business and evaluate whether these campaigns capture people's attention. It will also suggest ways to improve them. In the second part, preprocessing of the dataset and its analysis is discussed to gain a deeper understanding of the data. The third section focuses on the statistical tests and their analysis, with the results presented in the fourth section. Finally, the conclusion part summarizes people's responses to the campaigns and offers several suggestions for improving performance.

II. Data Exploration

a. Dataset

The dataset has 2,240 rows and 29 columns. Among the features of the dataset, there are 10 attributes related to People, 6 attributes related to Products, 7 attributes related to Promotion, and 4 attributes related to Place.

Table 1: Related to People

Feature	Description
ID	Customer's unique identifier
Year_Birth	Customer's birth year
Education	Customer's education level
Marital_Status	Customer's marital status
Income	Customer's yearly household income
Kidhome	Number of children in customer's household
Teenhome	Number of teenagers in customer's household
Dt_Customer	Date of customer's enrollment with the company
Recency	Number of days since customer's last purchase
Complain	1 if the customer complained in the last 2 years, 0 otherwise

Table 2: Related to Products

Feature	Description
MntWines	Amount spent on wine in the last 2 years
MntFruits	Amount spent on fruits in the last 2 years
MntMeatProducts	Amount spent on meat in the last 2 years
MntFishProducts	Amount spent on fish in the last 2 years
MntSweetProducts	Amount spent on sweets in the last 2 years
MntGoldProds	Amount spent on gold in the last 2 years

Table 3: Related to Promotion

Feature	Description
NumDealsPurchases	Number of purchases made with a discount
AcceptedCmp1	1 if customer accepted the offer in the 1st campaign, 0 otherwise
AcceptedCmp2	1 if customer accepted the offer in the 2nd campaign, 0 otherwise
AcceptedCmp3	1 if customer accepted the offer in the 3rd campaign, 0 otherwise
AcceptedCmp4	1 if customer accepted the offer in the 4th campaign, 0 otherwise
AcceptedCmp5	1 if customer accepted the offer in the 5th campaign, 0 otherwise
Response	1 if customer accepted the offer in the last campaign, 0 otherwise

Table 4: Related to Place

Attribute	Description
NumWebPurchases	Number of purchases made through the company’s website
NumCatalogPurchases	Number of purchases made using a catalogue
NumStorePurchases	Number of purchases made directly in stores
NumWebVisitsMonth	Number of visits to the company’s website in the last month

b. EDA

Initially, by checking for duplicate data, it was found that dataset does not have any duplicates. By looking at a few rows from the beginning and the end, and selecting 5 random rows, a general understanding of dataset was obtained. After examining contents of 29 columns, ID column dropped because it was not useful for the analysis. Since Customer’s birth year does not provide much specific information, this column was replaced by a new column called ‘Age_Category’, which categorizes age into 5 groups: ‘Child’, ‘Teen’, ‘Young Adult’, ‘Adult’, and ‘Senior’. Consequently, the Year_Birth and Age columns were dropped since they were not needed anymore, and Age_Category remained in the dataset. Next, checking for missing values revealed that there were 24 missing values in the Income column, which were replaced with the mean value to handle them. To determine the data types for each column, the Education and Marital_Status columns, which are categorical in nature, were changed to categorical types. Since the decimal part of the Income column was 0 in all cases, its data type was changed to integer. For the Kidhome column, it was found that the values only fell into three categories (no children, 1 child, and 2 children), so it was converted to a categorical feature. The same process was repeated for the Teenhome column. Date of customer’s enrollment with the company did not provide much useful information in its initial form, so only the year was extracted. It was found that the data only includes the years 2012, 2013, and 2014, with the following counts for customer enrollments: 494 in 2012, 1189 in 2013, and 557

in 2014. By analyzing the feature 'Recency' and categorizing this column into three categories: less than a week, less than a month, and more than a month, it was found that the number of days since the customer's last purchase is highest for those in the 'more than a month' category, with only a small number (just 172 cases) being related to 'less than a week.'

Since the main goal of this research was to examine the impact of campaigns on customers, the content of the five columns related to whether customers accepted campaigns 1 to 5 provided interesting information that will be discussed further.

The 'Complain' column indicated whether a customer had made a complaint. Out of 2240 cases, only 21 customers complained, and these 21 cases were stored in a separate variable for potential future analysis.

Further analysis of 'Z_CostContact' and 'Z_Revenue' revealed that the values for these two columns were 11 and 3 for all rows, respectively, and due to the lack of variation, they were removed from the dataset. In following, 3 tables categorized by data type can be seen:

Table 5: Numerical Data

Feature
Income
MntWines
MntFruits
MntMeatProducts
MntFishProducts
MntSweetProducts
MntGoldProds
NumDealsPurchases
NumWebPurchases
NumCatalogPurchases
NumStorePurchases
NumWebVisitsMonth

Table 6: Categorical Data

Feature
Education
Marital_Status
Kidhome
Teenhome
Age_Category

Table 7: Boolean Data

Feature
AcceptedCmp1
AcceptedCmp2
AcceptedCmp3
AcceptedCmp4
AcceptedCmp5
Response

To achieve a proper analysis of the impact of campaigns on customers, it's important to first understand the features of the dataset. The variable 'Income' will be analyzed in detail in the next section. For now, aiming to examine the 'Amount spent in the last 2 years' on Wines, Fruits, Meat Products, Fish Products, Sweet Products, and Gold Products (Figure 1). On average, the highest spending is on Wine, followed by Meat, and then, with a significant difference, Gold. This analysis is important for understanding the influence of campaigns on customers, as it reveals which products customers prioritize for their purchases.

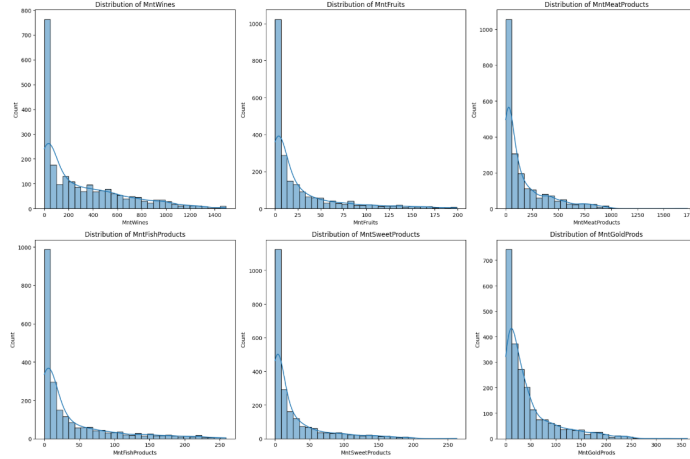


Figure 1: Amount spent in the last 2 years

Next, checking number of purchases made' through company's website, catalog, in-person shopping, and during discount periods show on average, most purchases are made in person, followed by purchases from store's website. There is a significant difference in purchases made through catalog, and least amount of purchases is made during discount periods. This initial analysis indicates that most purchases are made in person, and discounts do not seem to attract customers. (Figure 2)

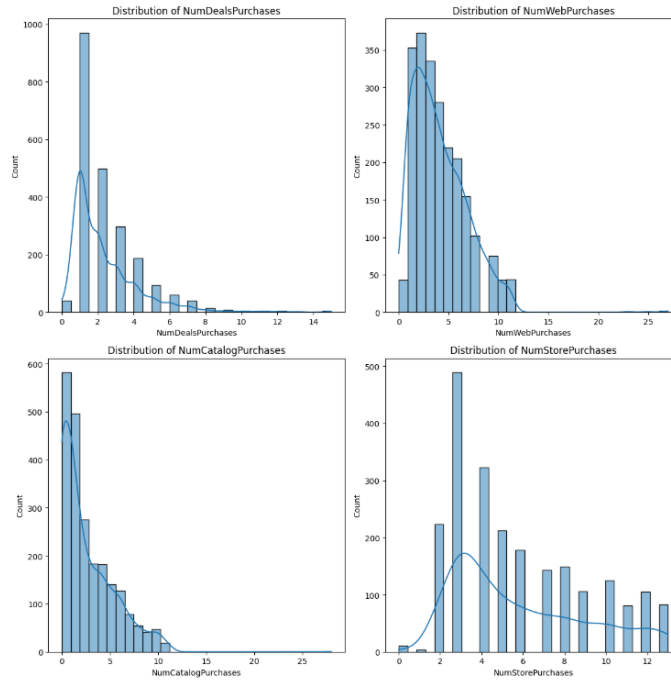


Figure 2: Number of purchases made

By looking at the plot below, one can observe the distribution of website visits per month. Considering that it was determined earlier that online shopping is the second priority for customers, approximately 7 visits seems reasonable. (Figure 3)

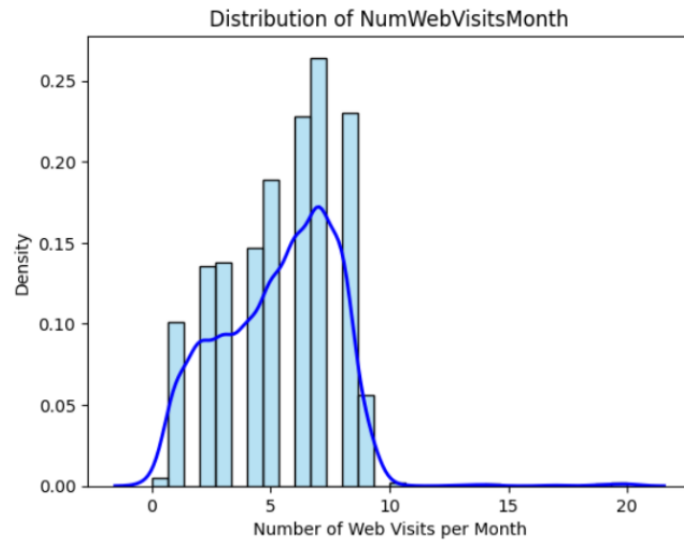


Figure 3: Web visiting

To have a better analysis of the categorical features, we will first examine 'Education.' According to the plot below, about half of the customers have a Graduation level, followed by those with a PhD at 21.6% and those with a Master's degree at 16.5%. (Figure 4)

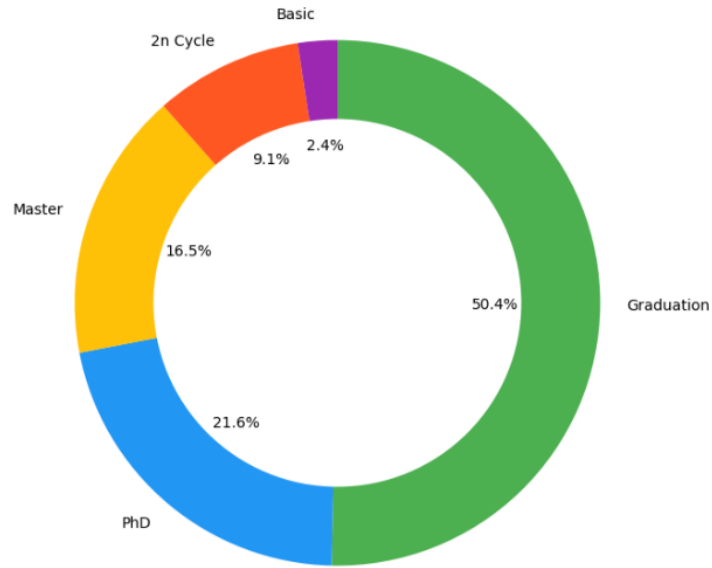


Figure 4: Education

Additionally, when examining marital status, we see that the highest number of cases are 'Married' with 861, followed by 'Together' with 756, and 'Single' with 480. (Figure 5)

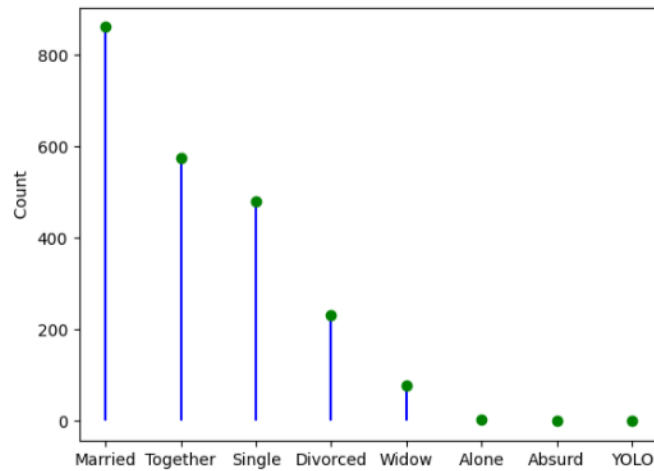


Figure 5: Marital Status

Table 8 shows information about three features: 'Number of children in the

customer's household,' 'Number of teenagers in the customer's household,' and 'Age Category.' This information will be useful for analyzing the acceptance of campaigns by customers.

Table 8: Customer Demographics

Age Category	Count
Adult	1571
Senior	618
Young Adult	43
Child	0
Teen	0

Teenagers Num	Count
0	1152
1	1028
2	52

Kids Num	Count
0	1288
1	896
2	48

Therefore, the dataset mainly contains information about adults who do not have children.

Since the main objective of this research is to examine the acceptance of campaigns by customers, a more detailed analysis will be conducted to determine if customers accepted the offer in 5 different campaigns. According to Table 9 and Figure 6, it can be observed that the second campaign has the lowest acceptance rate among customers, with a significant difference at 0.0134.

The remaining four campaigns have acceptance rates ranging from 0.0645 to 0.0748, with the fourth campaign leading the way.

Table 9: Campaign Acceptance

Campaign Num	Accept	Reject	Acceptance Rate
1	144	2088	0.0645
2	30	2202	0.0134
3	163	2069	0.0730
4	167	2065	0.0748
5	163	2069	0.0730

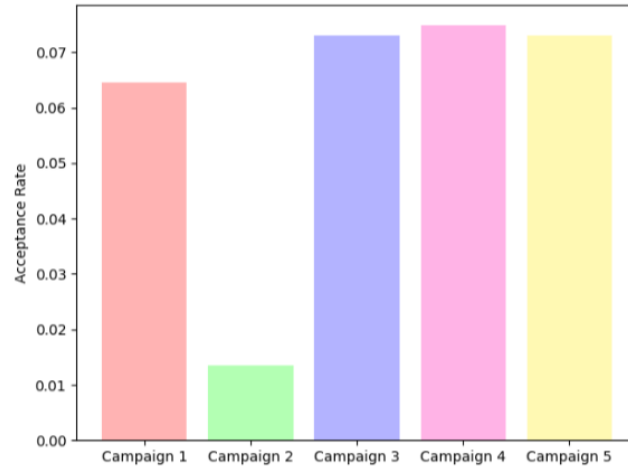


Figure 6: Campaign Acceptance

As the last point in this section, examining most important feature. Based on the analysis, out of 2,240 customers, only 334 responded positively to the offer in the last campaign, which corresponds to a 14.96% response rate. This is a very low number for a business.

Finally, to achieve better organization and easier implementation, used label encoding with LabelEncoder. To identify the corresponding labels, referring to the mapping provided in Table 10.

Table 10: Mapping

Code	Education Level
0	2n Cycle
1	Basic
2	Graduation
3	Master
4	PhD

Code	Marital Status
0	Absurd
1	Alone
2	Divorced
3	Married
4	Single
5	Together
6	Widow
7	YOLO

Code	Age Category
0	Adult
1	Senior
2	Young Adult

Code	Kidhome
0	0
1	1
2	2

Code	Teenhome
0	0
1	1
2	2

III. Statistical Analysis

The analyses conducted in this section aim to understand customer personality in order to design better campaigns to attract customers. As seen in the previous section, the response rate of customers to the campaigns is below 20%.

Therefore, by performing various statistical tests in this section, aim to identify the factors that influence this response rate.

First, analyzing Income to check the distribution of this feature. Using the Shapiro-Wilk test to see if it follows a normal distribution. This requires examining and removing any outliers from the dataset. By calculating the IQR and comparing the Income values with the lower and upper bounds, outliers have been handled. The Shapiro-Wilk test then shows that Income does not follow a normal distribution. (Figure 7)

Shapiro-Wilk Test:

$$W = \frac{(\sum_{i=1}^n a_i x_{(i)})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Where:

- W = Shapiro-Wilk test statistic
- n = number of observations in the sample
- $x_{(i)}$ = the i -th order statistic (the i -th smallest value in the sample)
- \bar{x} = sample mean
- a_i = constants derived from the expected values of the order statistics of a normal distribution

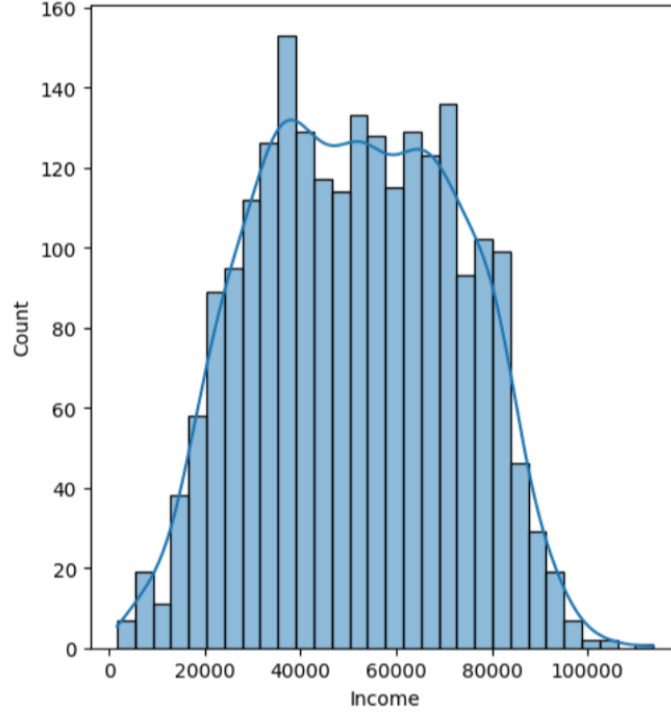


Figure 7: Income Distribution

Note: There is a common rule in statistics known as the Central Limit Theorem (CLT) which states that, given a sufficiently large sample size (usually n more or equal to 30 is considered sufficient), the sampling distribution of the sample mean will be approximately normally distributed, regardless of the shape of the population distribution.

As observed in the previous section, the majority of customers have education levels of Graduation, PhD, and Master (about 85% of the total dataset). A Kruskal-Wallis H test was conducted to check if customers from different education levels have different income levels. The results showed that there are significant differences between the income levels of different educational backgrounds.

Kruskal-Wallis H Test:

$$H = \frac{12}{N(N+1)} \sum_{i=1}^k \frac{R_i^2}{n_i} - 3(N+1)$$

Where:

- H = Kruskal-Wallis test statistic

- N = total number of observations across all groups
- k = number of groups
- R_i = sum of ranks for the i -th group
- n_i = number of observations in the i -th group

According to the violin plot shown in Figure 8, customers with a Basic education have lower incomes. However, since a substantial part of the dataset consists of the three groups: Graduation, PhD, and Master, and the plot shows that these three groups have similar income levels, it can be said that, overall, customers have an acceptable level of education and income. Therefore, analyzing other features may be more effective in understanding the lower acceptance of campaigns by customers.

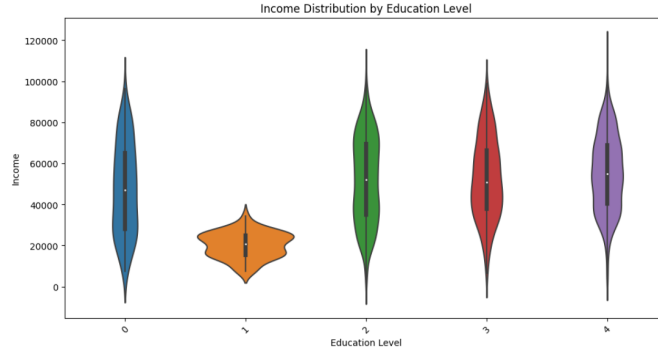


Figure 8: Education Level vs. Income

Another interesting point that can be observed from Figure 9 is that as income increases, spending on wine also increases. This makes sense because the dataset mainly consists of educated individuals with a good income. Since wine is the best-selling product, it shows a clear connection with income.

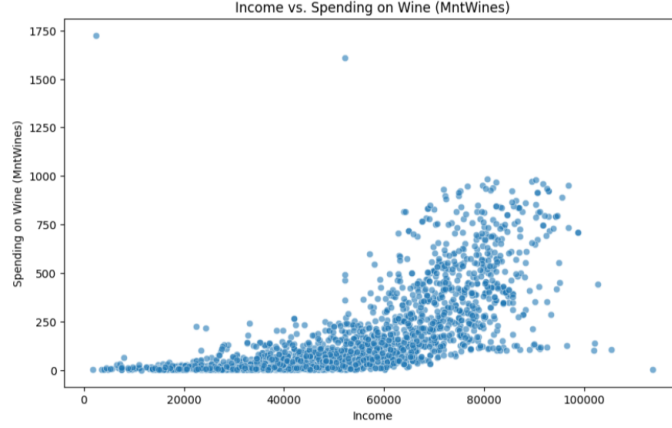


Figure 9: Spending on Wine vs. Income

To examine the impact of campaigns on spending, can analyze spending behavior across different categories. The dataset contains information on the amount spent on products like wine, fruits, meat, fish, sweets, and gold. By comparing spending distributions between those who accepted the campaign and those who did not, we can see where the campaign had an effect. Using the Independent Samples T-test, can look at the spending data for the six product categories based on the two groups: Accepted and Not Accepted. It is clear that there is a significant difference in spending behavior across all categories. For lower spending amounts, the distribution of those who did not accept the campaign is higher. However, as spending increases, those who accepted the campaign show higher spending levels. This suggests that customers who spend a lot are more likely to respond positively to the campaigns. Notably, wine had the highest sales, and the impact of the campaign on sales becomes visible at spending amounts above \$500. Therefore, we can conclude that campaigns have a positive effect on higher spending and popular products. (Figure 10)

Independent Samples T-test:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Where:

- \bar{X}_1 = Mean of group 1
- \bar{X}_2 = Mean of group 2
- s_1^2 = Variance of group 1
- s_2^2 = Variance of group 2

- n_1 = Sample size of group 1
- n_2 = Sample size of group 2

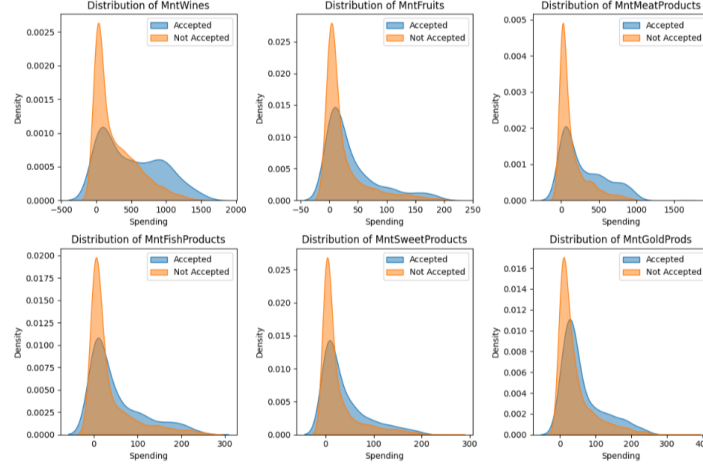


Figure 10: Spending in Accepted Group vs. Not Accepted

Considering that 71.6% of customers have children (kid or teen), will now examine whether customers with children spend differently than those without children. To do this, using the Mann-Whitney U Test to compare the spending distributions of the independent groups (with children and without children). Figure 11 shows that spending distributions in these two groups are significantly different, indicating that having children affects spending. Individuals without children have a strong spending distribution in the range of approximately 50 to 150, whereas the distribution for individuals with children is flatter.

Mann-Whitney U Test:

$$U = R_1 - \frac{n_1(n_1 + 1)}{2}$$

Where:

- R_1 = Sum of ranks for the first group
- n_1 = Number of observations in the first group

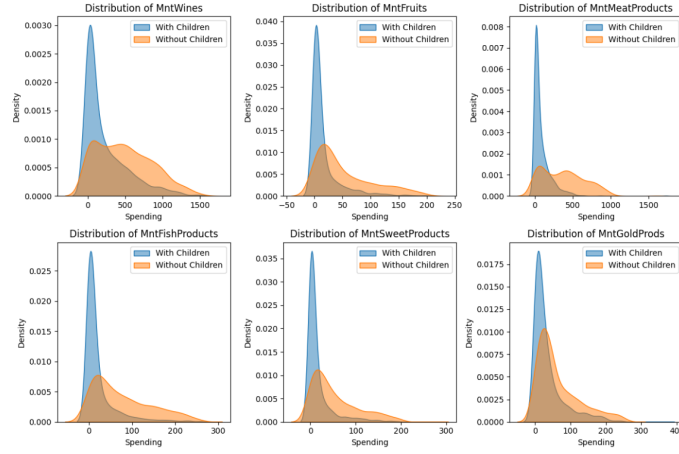


Figure 11: Spending in with Children vs. without

o analyze the impact of having children (either teenagers or young children) on campaign acceptance, a Chi-Squared Test was conducted. The results indicate a significant difference in campaign acceptance between customers with children and those without. According to Figure 12, the acceptance of campaigns is higher among individuals without children. However, overall, the majority of customers, whether they have children or not, show very low engagement with the campaigns.

Chi-Squared Test:

$$\chi^2 = \sum \frac{(O - E)^2}{E} \quad (1)$$

Where:

- χ^2 is the Chi-Squared statistic,
- O is the observed frequency,
- E is the expected frequency.

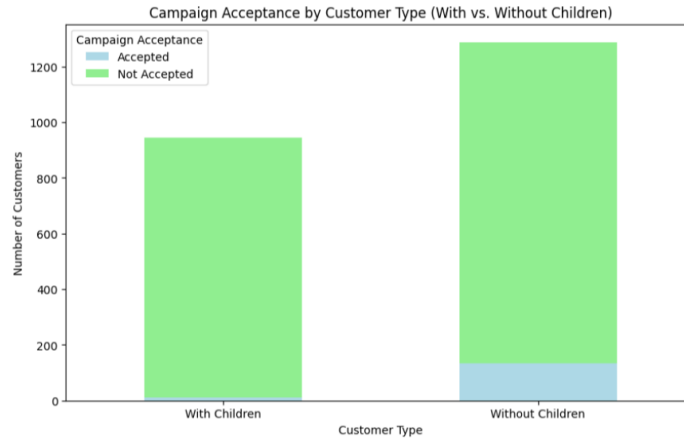


Figure 12: Acceptance in with Children vs. without

It has been determined that the acceptance of campaigns is lower among individuals with children. Since the goal of this research is to help businesses increase customer engagement with campaigns, Figure 13 shows that individuals with teenagers are more receptive to campaigns than those with young children. Therefore, changes should be made in the design of campaigns targeted at individuals with young children to capture their attention. This hypothesis is tested using the Chi-Squared Test, similar to the previous analysis.

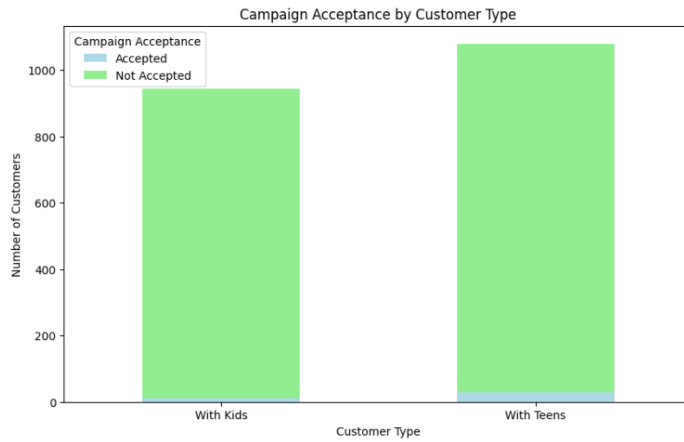


Figure 13: Acceptance in with Kid vs. with Teen

Further tests can be conducted to examine the impact of other categorical features on campaign acceptance. Figures 14 to 16 explore the effects of education level, age category, and marital status on the acceptance of promotional campaigns. According to Figure 14, although the overall acceptance of

campaigns is very low, different education levels show varying acceptance rates. Graduation, PhD, and Master's degree holders rank the highest, indicating that there is a significant relationship between education level and acceptance of promotional campaigns.

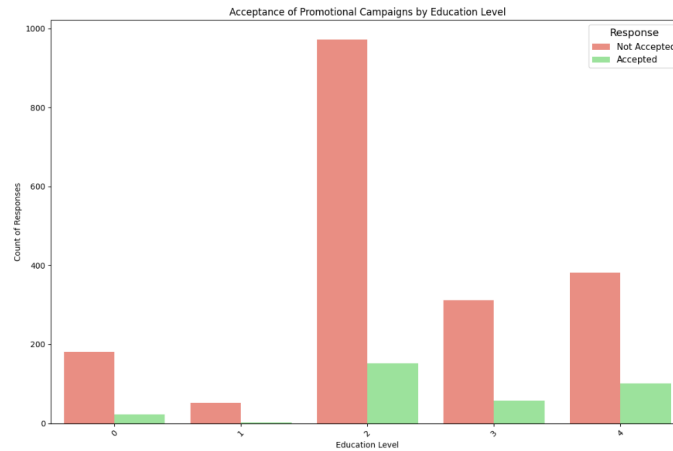


Figure 14: Education Level

In Figure 15, it is observed that there is no significant association between age category and response. The acceptance rates across different age groups do not vary much, suggesting that age is not a major factor influencing campaign effectiveness.

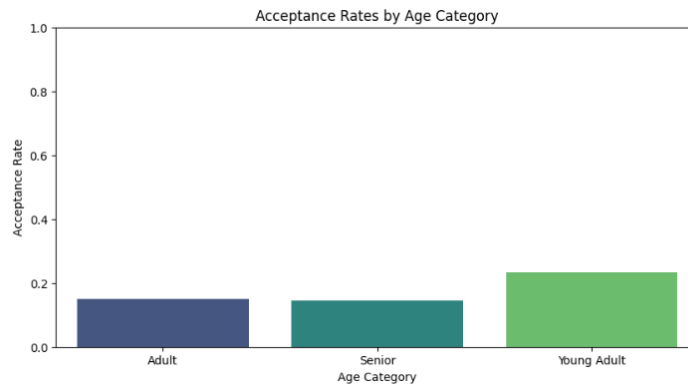


Figure 15: Age Category

Figure 16 reveals a significant association between marital status and response. Individuals with the marital statuses of YOLO, Absurd, and Alone have accepted campaigns more frequently, while those who are Together or Married

show the lowest levels of acceptance.

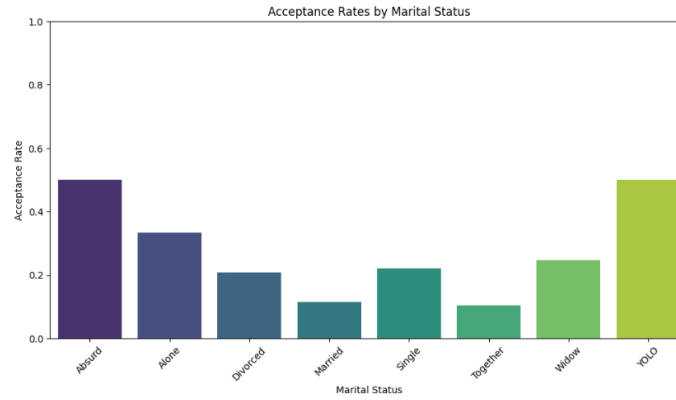


Figure 16: Marital Status

Next, to determine which product category to focus on, using the Friedman Test to compare customer spending across multiple product categories. The results show that the spending in different categories is significantly different. According to Figure 17, the top three product categories are wine, meat, and gold, in that order.

Friedman Test:

$$\chi^2 = \frac{12}{n \cdot k \cdot (k + 1)} \sum_{j=1}^k R_j^2 - 3n(k + 1)$$

Where:

- n = Number of subjects
- k = Number of treatment conditions (or product categories)
- R_j = Sum of ranks for the j -th treatment

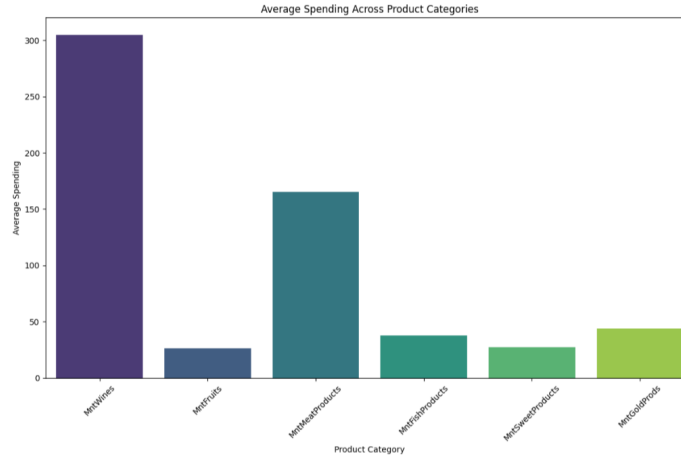


Figure 17: Spending on Products

IV. Results

In this section, the results of the tests conducted earlier will be reviewed and analyzed. The findings will be used to evaluate the impact of campaigns and the level of acceptance by customers. Before that, the details related to Figures 1 and 2 will be examined.

Table 11: Amount of products spent in the last 2 years

Category	Wine	Fruit	Meat	Fish	Sweet	Gold
Mean	304.93	26.38	165.32	37.64	27.16	44.17
Standard Deviation	336.78	39.82	219.40	54.69	41.33	52.20
Minimum Value	0.00	0.00	0.00	0.00	0.00	0.00
Maximum Value	1493.00	199.00	1725.00	259.00	263.00	362.00

Table 12: Summary of Purchase Statistics

Category	Deals Purchases	Web Purchases	Catalog Purchases	Store Purchases
Mean	2.32	4.10	2.64	5.81
Standard Deviation	1.89	2.77	2.80	3.24
Minimum Value	0.00	0.00	0.00	0.00
Maximum Value	15.00	27.00	28.00	13.00

Analyzing Tables 11 and 12 shows that the highest sales for this business

are, on average, related to wine, meat, and gold. Additionally, the most popular purchasing methods, in order, are in-store purchases, online purchases, and catalog purchases, with deals being the least popular. To improve campaign performance, more focus can be placed on the best-selling products and the more popular sales methods.

Also, Table 13 indicates statistics related to income and as seen mean income is about 52000 and based in p-value which is less than 0.05 Shapiro-Wilk Test shows Income does not have normal distribution. The null hypothesis (data is normally distributed) is rejected. (Related to Figure 7: Income Distribution)

Table 13: Income Statistics

Mean	Std	Min	Max
52247.248661	25037.797168	1730.000000	666666.000000

Test Name	Statistics	P-value
Shapiro-Wilk Test	0.986	3.26×10^{-14}

The analysis of the impact of education level on income using the Kruskal-Wallis H test shows that with a p-value less than 0.05, the null hypothesis is rejected. This indicates that there are significant differences in income levels among different education backgrounds. (Related to Figure 8: Education Level vs. Income)

Table 14: Kruskal-Wallis Test Results

Test Name	Statistic	P-value
Kruskal-Wallis H test	140.85	1.85×10^{-29}

The examination of how the marketing campaign influences spending behavior across customer groups was conducted using the Independent Samples T-test. The results showed that with a p-value less than 0.05 for all products, the null hypothesis is rejected. This means there is a significant difference in spending behavior for the products. (Related to Figure 10: Spending in Accepted Group vs. Not Accepted)

Table 15: Independent Samples T-test Results

Test Name	Category	Statistics	P-value
Independent Samples T-test	Wine	12.01	3.19×10^{-32}
Independent Samples T-test	Fruit	5.93	3.52×10^{-09}
Independent Samples T-test	Meat	12.02	2.57×10^{-32}
Independent Samples T-test	Fish	5.25	1.65×10^{-07}
Independent Samples T-test	Sweet	5.54	3.35×10^{-08}
Independent Samples T-test	Gold	6.62	4.40×10^{-11}

The impact of having children on spending was analyzed using Mann-Whitney U test. With a p-value less than 0.05, null hypothesis is rejected, indicating a significant difference in spending behavior for products. Furthermore, to examine whether customers with children accept campaigns as much as those without, Chi-Squared Test was applied. The results revealed that with a p-value less than 0.05, there is a significant difference in campaign acceptance between customers with children and those without children. Additionally, a more detailed analysis using the same test on the group with children showed significant differences in campaign acceptance between customers with children and those with teenagers. Related to Figures 11, 12, 13)

Table 16: Mann-Whitney U Test Results

Test Name	Category	Statistics	P-value
Mann-Whitney U test	Wine	293975.5	4.49×10^{-54}
Mann-Whitney U test	Fruit	221158.0	1.37×10^{-96}
Mann-Whitney U test	Meat	189243.5	3.52×10^{-118}
Mann-Whitney U test	Fish	201201.0	3.26×10^{-110}
Mann-Whitney U test	Sweet	228886.5	1.33×10^{-91}
Mann-Whitney U test	Gold	322735.5	6.88×10^{-41}

Table 17: Chi-Squared Test Results for with vs. without children

Test Name	Statistic	P-value
Chi-Squared Test	74.24	6.93×10^{-18}

Table 18: Chi-Squared Test Results for kid vs. teen

Test Name	Statistic	P-value
Chi-Squared Test	5.25	0.022

Next, Friedman Test used to check if there are any significant differences in spending on different product categories. P-value less than 0.05, so there is a significant difference in how much people spend across different categories. (Related to Figure 17)

Table 19: Friedman Test Results

Test Name	Statistic	P-value
Friedman Test	5949.41	0.0

Then, to see if there's a relation between customer education level and how respond to promotional campaigns, Chi-Squared Test indicated that there is a significant relationship. (Related to Figure 14)

Table 20: Chi-Square Test Results for relationship between customer education level and acceptance of promotional campaigns

Test Name	Statistic	P-value	Degrees of Freedom
Chi-Square Test	23.59	9.64×10^{-5}	4

However, when looking at relationship between age categories and acceptance of promotional campaigns, Chi-Squared Test found no significant connection. P-value was higher than alpha level of 0.05, meaning it wasn't statistically significant. (Related to Figure 15)

Table 21: Chi-Square Test Results for relationship between customer age category level and acceptance of promotional campaigns

Test Name	Statistic	P-value
Chi-Square Test	2.48	0.29

Finally, regarding marital status and acceptance of promotional campaigns, Chi-Squared Test showed a significant association. This means marital status

does seem to play a role in how individuals respond to promotions. (Related to Figure 16)

Table 22: Chi-Square Test Results for relationship between Marital status and acceptance of promotional campaigns

Test Name	Statistic	P-value
Chi-Square Test	53.66	2.75×10^{-9}

V. Conclusion

The main goal of this research was to analyze customer personality in order to assess the acceptance of campaigns by customers and to identify various factors influencing this acceptance. The findings revealed that only 14.96% of customers responded positively to these campaigns, which is a very low figure. Therefore, an analysis of different sections of the dataset was conducted to gather sufficient reasons for the business to improve its sales policies or consider replacing them with new strategies. The analysis showed that the best-selling products were wine, meat, and gold, with significant impacts from both in-store and online sales. This suggests that future campaigns should focus on these top-selling products. Additionally, no significant relationship was found between age and campaign acceptance among customers, indicating that age is not a significant factor. However, marital status and having children are important factors influencing campaign acceptance. In particular, individuals with children, especially those with teenagers, showed very low interest in the campaigns. Therefore, it is recommended to design future campaigns specifically targeting this group. Since online sales emerged as one of the top sales methods, and customer visits to the website have favorable statistics, it is advisable to display campaigns through the website to reach customers effectively. For future work, it is suggested to conduct a more in-depth analysis of the 24 complaints made by customers. Additionally, a closer examination of the second campaign, which had the lowest acceptance rate, should be carried out. More tests can be implemented for further analysis of the dataset, including a detailed review of five campaigns based on spending behaviors.