

# CREDIT EDA CASE STUDY

By Rajat Sachan

# Introduction

This assignment focuses on applying EDA (Exploratory Data Analysis) in a practical business context. It provides insight into risk analytics within the banking and financial services sector, showcasing how data analysis is used to mitigate financial lending risks. Through EDA techniques, you'll gain valuable skills while understanding how data plays a crucial role in minimizing potential losses when lending to customers.

# Business Understanding

- In the context of a consumer finance company, we aim to analyze loan application data to identify patterns in loan defaults. The company faces two key risks: approving loans to likely repayers (avoiding loss of business) and to likely defaulters (resulting in financial loss).
- The data distinguishes clients with payment difficulties and those who pay on time. There are four possible outcomes: approved, canceled, refused, and unused offers. Our goal is to use EDA to understand how client and loan attributes influence loan default tendencies, aiding the company in better lending decisions.

# Business Objectives

- This case study's objective is to detect patterns indicating clients' payment difficulties, enabling actions like loan denial, reduced loan amounts, or higher interest rates for risky applicants. The goal is to ensure deserving loan applicants are not rejected.
- In simpler terms, the company aims to identify key factors (driver variables) leading to loan defaults—strong indicators of default. This knowledge aids portfolio management and risk assessment. To better grasp the domain, it's advisable to conduct some independent research on risk analytics, focusing on variable types and their importance.

# Workflow Walkthrough

- Data Understanding
- Data Cleaning and Manipulation
- Handling Outlier
- Addressing Data Imbalance
- Analyzing Data: Univariate, Segmented, Bivariate, and Correlation Analysis
- Analyzing Previous Applications with Similar Processing
- Merged Both Application and Previous Application
- Integrated Data Analysis: Univariate, Bivariate, and Correlation Analysis

# Initial Data Exploration

- Imported essential libraries: Pandas, Matplotlib, Numpy and Seaborn.
- Created two data frames from CSV files: "app\_df" and "pre\_app\_df".
- The "app\_data" dataframe has a dimension of 307,511 rows and 122 columns, while the "pre\_app\_df" consists of 1,670,214 rows and 37 columns.
- Both dataframes consist of varying data types, including int64, float64, and object.
- Examine summary statistics for numerical columns in the dataframes.

# Exploring Missing Data And Performing Data Cleaning And Manipulation

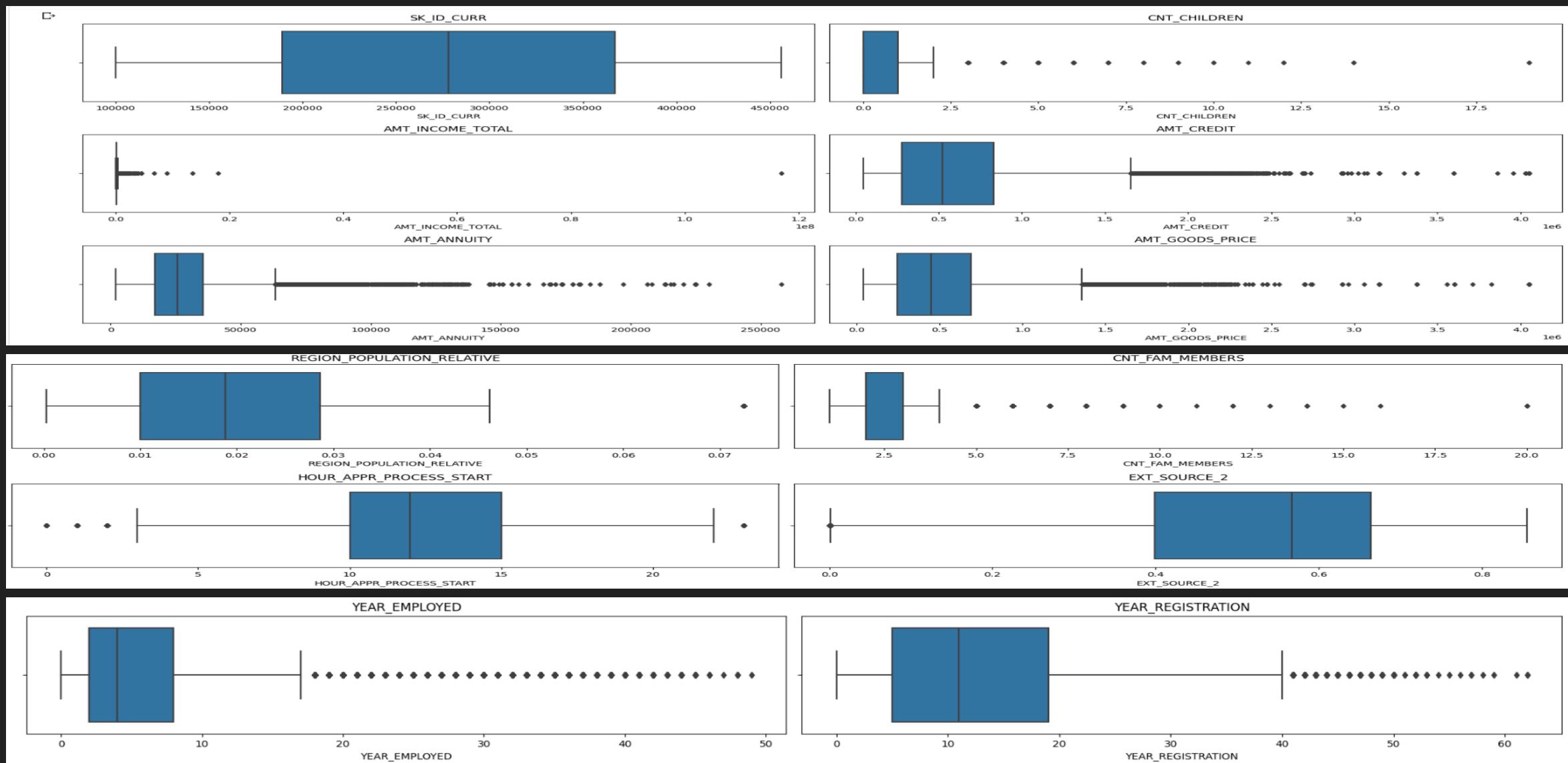
- Examining null values in each column and displaying percentage of null values.
- I assume that "OCCUPATION\_TYPE" significantly impacts loan approval by reflecting the applicant's financial stability and profession, justifying a 40% threshold for column removal.
- Removing columns with null values exceeding 40%.
- Some columns contain varying percentages of null values, ranging from 30% to 19%, or even as low as 0.33%. I decided to impute these null values using either the mean or mode, as appropriate.
- I identified certain columns containing NAN values or XNA entries. To tackle this, we removed the corresponding rows or replaced NAN or XNA values with the mode.

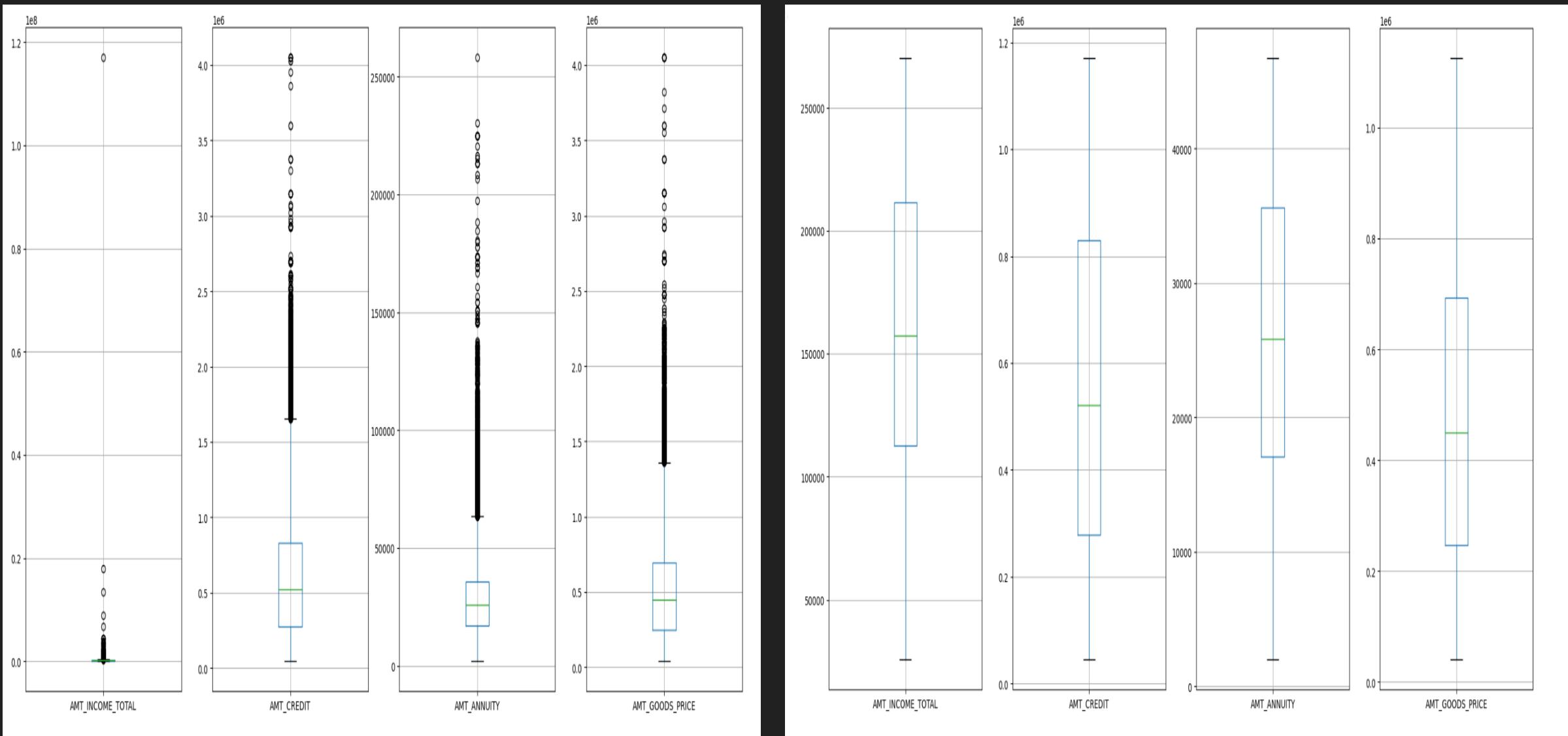
missing_percentage[missing_percentage > 30]
OWN_CAR_AGE 65.99
OCCUPATION_TYPE 31.35
EXT_SOURCE_1 56.38
APARTMENTS_AVG 50.75
BASEMENTAREA_AVG 58.52
YEARS_BEGINEXPLUATATION_AVG 48.78
YEARS_BUILD_AVG 66.50
COMMONAREA_AVG 69.87
ELEVATORS_AVG 53.30
ENTRANCES_AVG 50.35
FLOORSMAX_AVG 49.76
FLOORSMIN_AVG 67.85
LANDAREA_AVG 59.38
LIVINGAPARTMENTS_AVG 68.35
LIVINGAREA_AVG 50.19
NONLIVINGAPARTMENTS_AVG 69.43
NONLIVINGAREA_AVG 55.18
APARTMENTS_MODE 50.75
BASEMENTAREA_MODE 58.52
YEARS_BEGINEXPLUATATION_MODE 48.78
YEARS_BUILD_MODE 66.50
COMMONAREA_MODE 69.87
ELEVATORS_MODE 53.30
ENTRANCES_MODE 50.35
FLOORSMAX_MODE 49.76
FLOORSMIN_MODE 67.85
LANDAREA_MODE 59.38
LIVINGAPARTMENTS_MODE 68.35
LIVINGAREA_MODE 50.19
NONLIVINGAPARTMENTS_MODE 69.43
NONLIVINGAREA_MODE 55.18
APARTMENTS_MEDI 50.75
BASEMENTAREA_MEDI 58.52
YEARS_BEGINEXPLUATATION_MEDI 48.78
YEARS_BUILD_MEDI 66.50
COMMONAREA_MEDI 69.87
ELEVATORS_MEDI 53.30
ENTRANCES_MEDI 50.35
FLOORSMAX_MEDI 49.76
FLOORSMIN_MEDI 67.85

# Outlier Detection and Mitigation: Capping and Winsorization Techniques

- Initial observation: We detected the presence of outliers in several numerical columns, which could adversely affect our data analysis.
- Exploratory visualization: To identify these outliers, we created box plots for the numerical columns, enabling us to visually pinpoint the extreme values.
- Mitigation strategy: Following the identification of outliers, we implemented two techniques to handle them effectively: capping and winsorization.
- Capping: For extreme values exceeding a predefined threshold, we applied capping, replacing them with the threshold value to limit their impact on subsequent analyses.
- Winsorization: As an alternative approach, we employed winsorization, which involved replacing extreme outliers with the nearest non-outlier values, preserving the overall data distribution while reducing the influence of outliers.
- We included box plots below, which clearly visualize the presence of outliers in the data, reaffirming the importance of addressing these data points during our analysis.
- Negative Data Values: Some columns, including "DAYS\_BIRTH," "DAYS\_EMPLOYED," "DAYS\_REGISTRATION," and "DAYS\_ID\_PUBLISH," contain negative data values.
- Conversion to Absolute Values: In order to rectify this issue, we transformed these columns by converting the negative values into their absolute counterparts.

- Box Plot Analysis: Detecting Outliers in Data





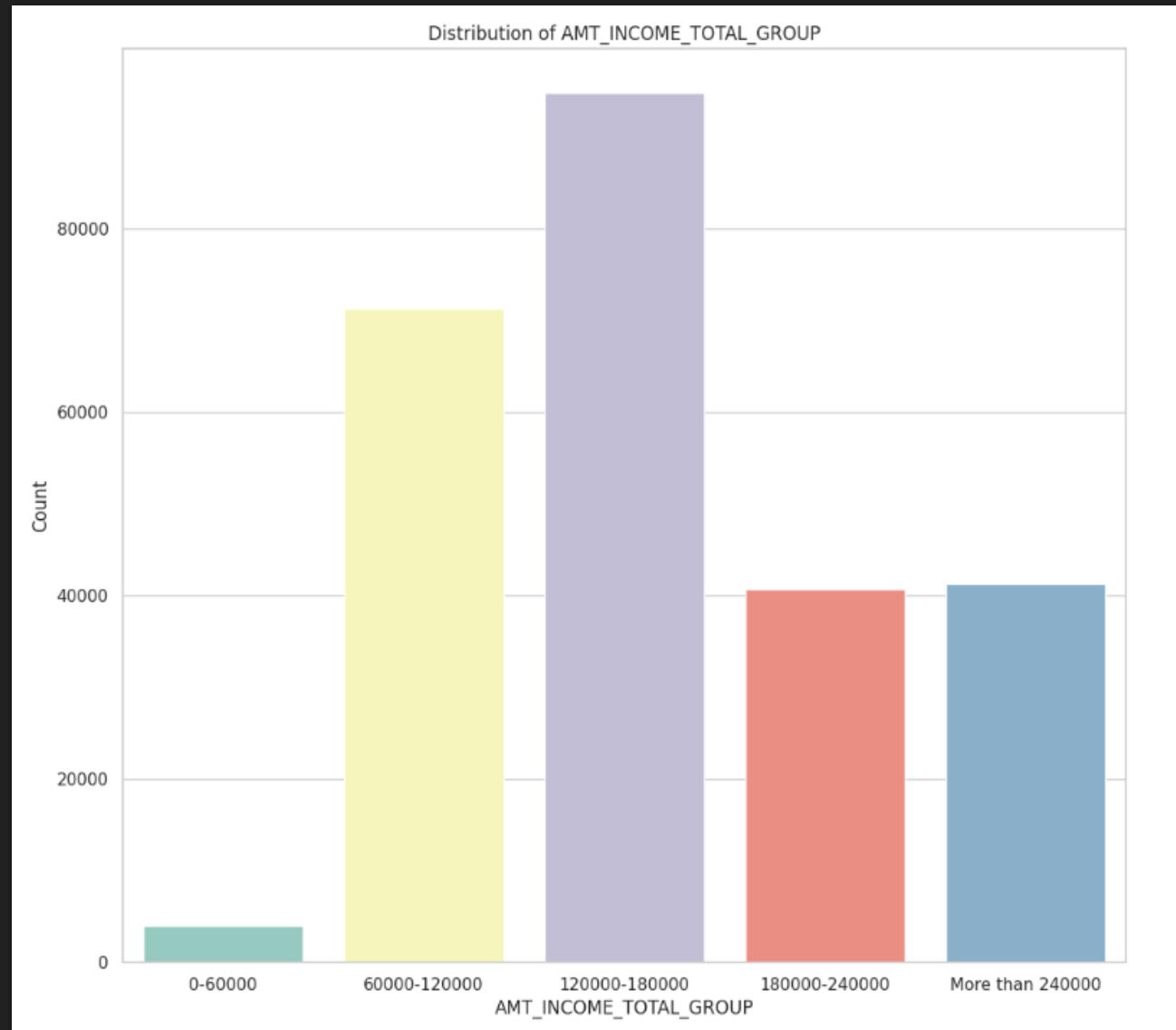
- These financial fields exhibit far-reaching outliers, signifying extreme values that warrant special attention and outlier management techniques.
- Outlier Display: The left side displays outliers.
- Post-Outlier Handling: On the right side, you can see the data after outlier handling.

# Univariate Analysis of Application Data

- Univariate Analysis of Numerical Columns: We conducted a univariate analysis focusing on our numerical columns.
- Binning Data: To facilitate this analysis, we initially divided the data into bins.
- Bar Graph Visualization: Subsequently, we created bar graphs to visually represent the data distribution for specific columns.
- Columns Included: Some of the columns analyzed in this manner included "AMT\_GOODS\_PRICE," "AMT\_INCOME\_TOTAL," "AMT\_CREDIT," and "AMT\_ANNUITY," among others. There are bar graphs displayed beneath this slide.
- Also created count plot and histogram visualizations for fields with object data types.
- In our data analysis, we've used count plots and histograms to visualize various dataset fields, gaining valuable insights for further analysis.

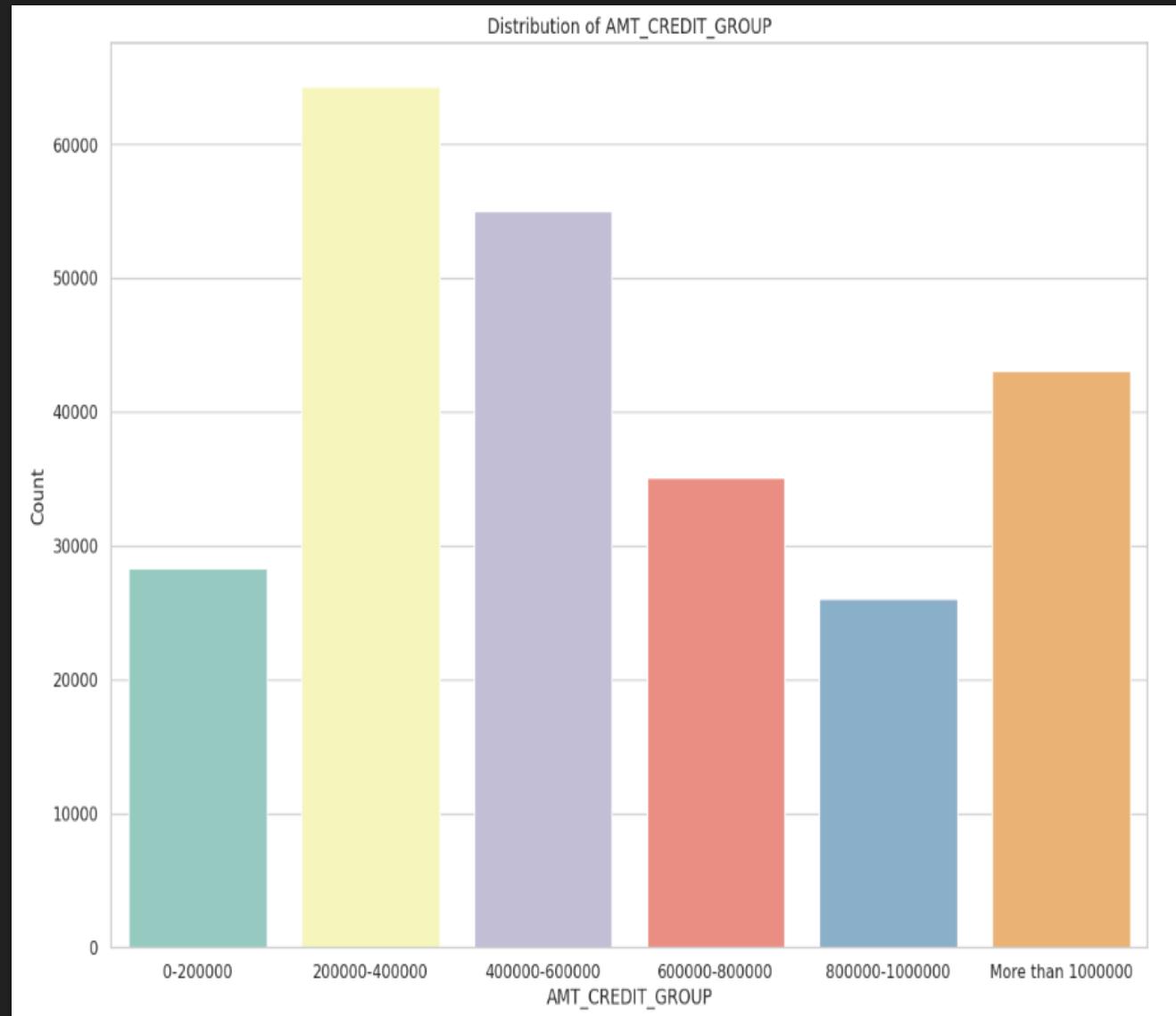
## Key Observations from AMT\_INCOME\_TOTAL Analysis

- 1. Income Group Analysis:** When we grouped the income data into AMT\_INCOME\_TOTAL\_GROUP bins, a prominent trend emerged.
- 2. Dominant Income Range:** The analysis clearly indicates that the largest number of applicants belong to the income range of 120,000 to 180,000 within the AMT\_INCOME\_TOTAL\_GROUP.
- 3. Lowest Applicant Count:** In contrast, the smallest number of applicants can be observed in the income bracket of 0 to 60,000, as revealed by the graph.



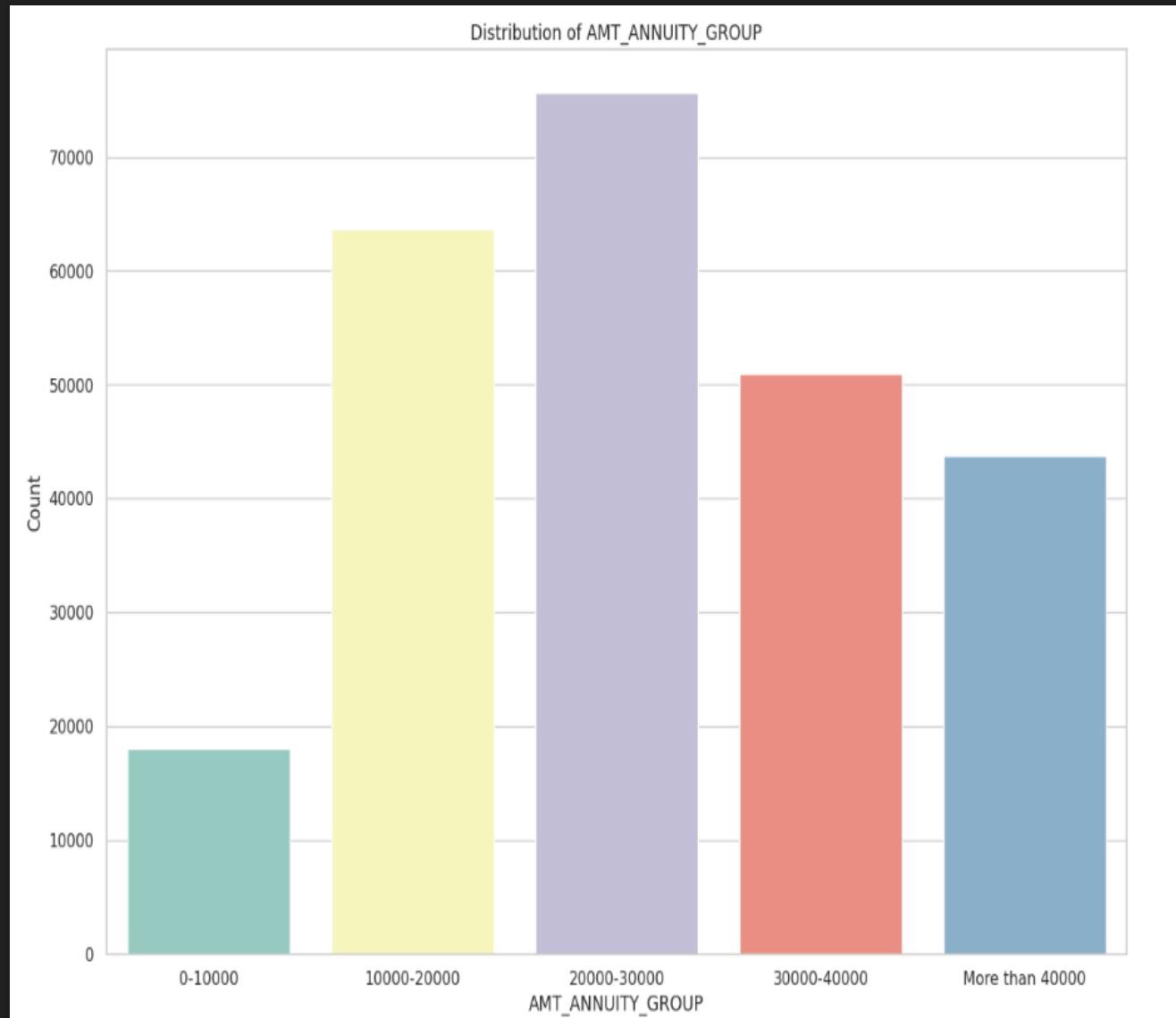
## Key Insights from AMT\_CREDIT Analysis

- 1. AMT\_CREDIT Grouping:** After categorizing AMT\_CREDIT into AMT\_CREDIT\_TOTAL\_GROUP bins, a distinct pattern emerges from the analysis.
- 2. Dominant Credit Range:** The analysis highlights that the largest number of applicants falls within the credit amount range of 200,000 to 400,000 within the AMT\_CREDIT\_TOTAL\_GROUP.
- 3. Lowest Applicant Count:** Conversely, the graph illustrates that the smallest number of applicants can be found in the credit amount bracket of 800,000 to 1,000,000.



## Key Observations from AMT\_ANNUITY Analysis

- 1. AMT\_ANNUITY Grouping:** After segmenting AMT\_ANNUITY into AMT\_ANNUITY\_GROUP categories, a distinct pattern becomes evident in our analysis.
- 2. Dominant Annuity Range:** The graph clearly indicates that the highest number of applicants falls within the loan annuity range of 20,000 to 300,000 within the AMT\_ANNUITY\_GROUP.
- 3. Lowest Applicant Count:** In contrast, the graph illustrates that the smallest count of applicants can be identified within the loan annuity range of 0 to 10,000.

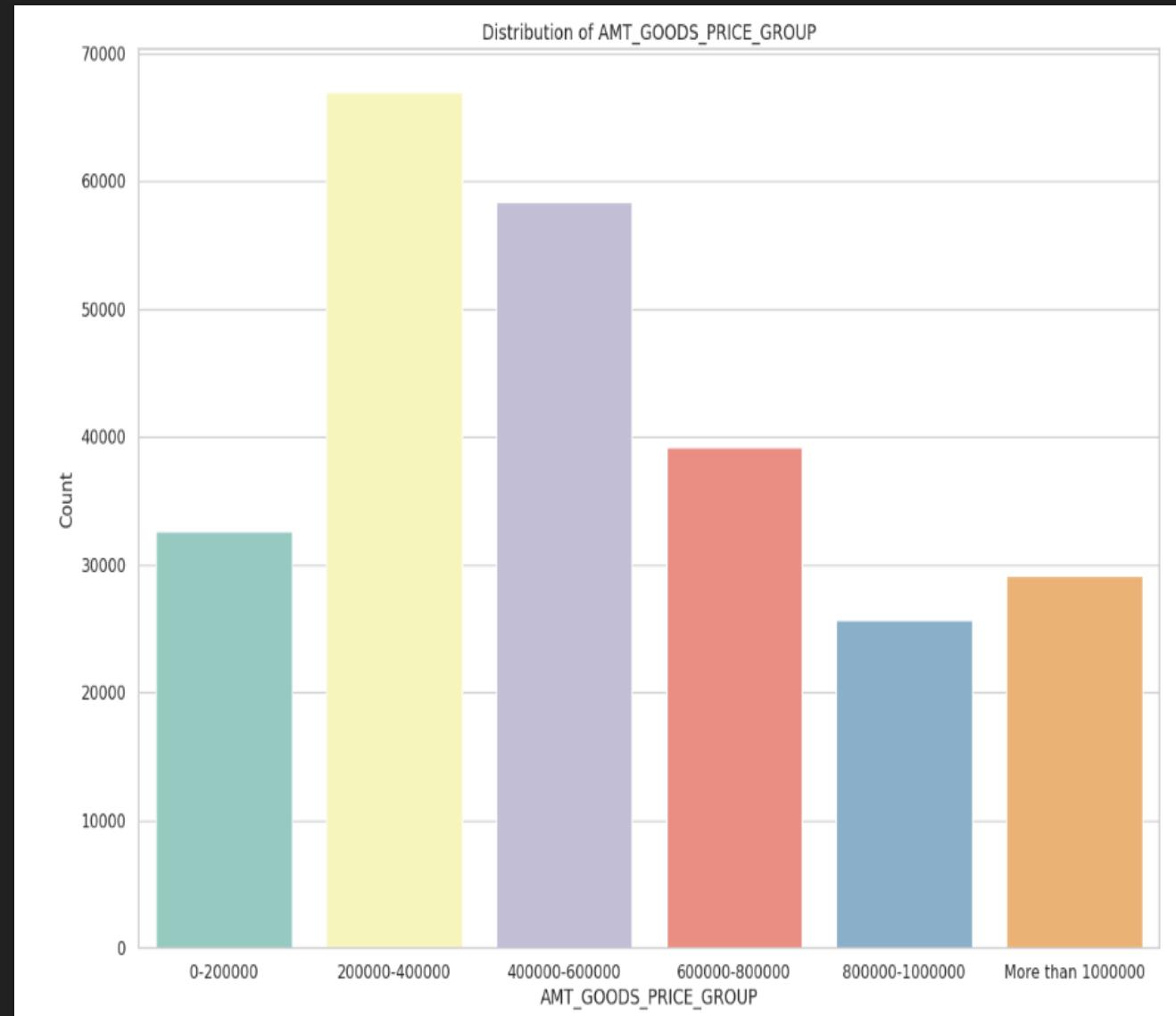


## Key Insights from AMT\_GOODS\_PRICE Analysis

**1. AMT\_GOODS\_PRICE Grouping:** After categorizing AMT\_GOODS\_PRICE into AMT\_GOODS\_PRICE\_GROUP categories, our analysis brings to light distinct trends.

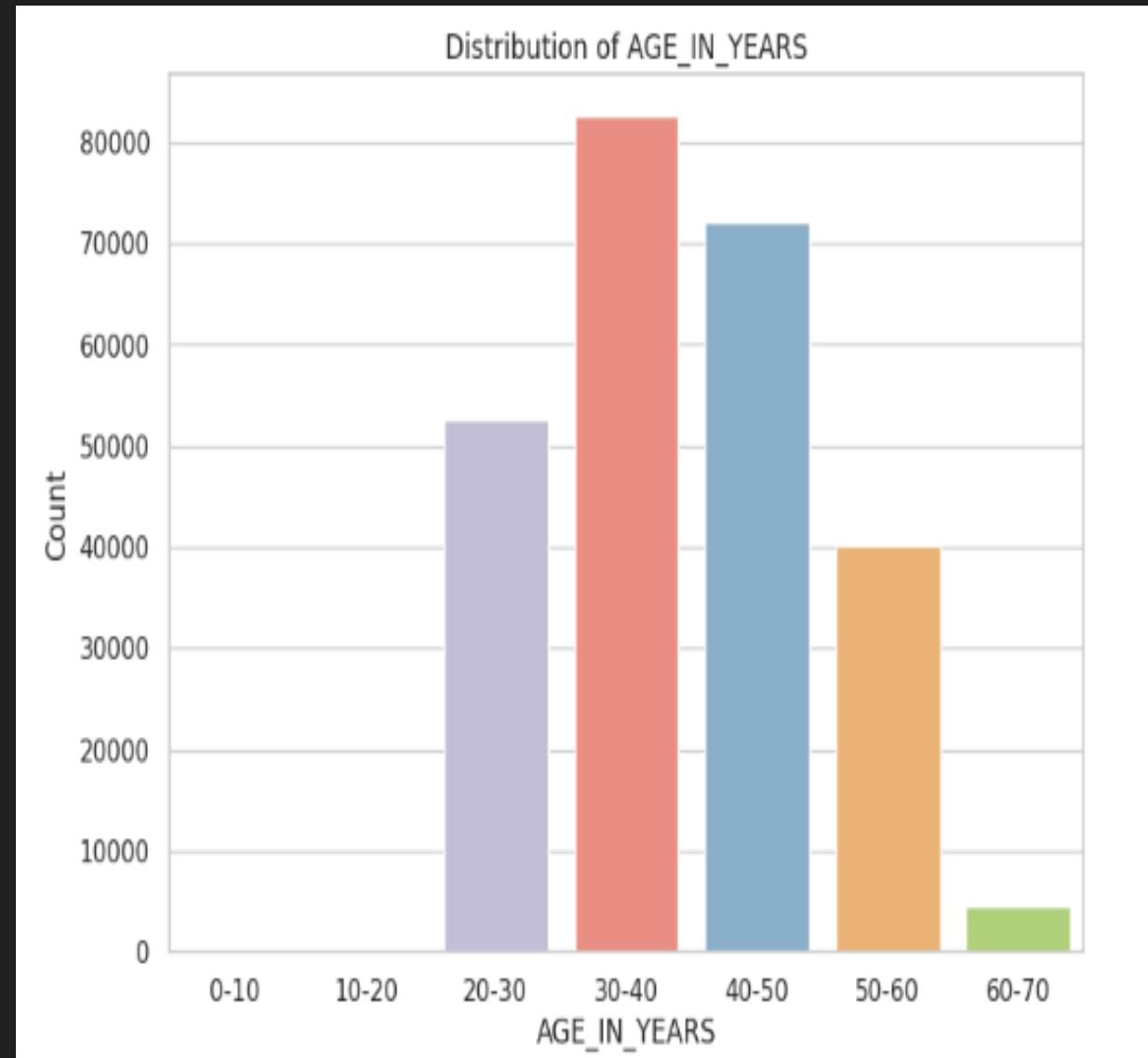
**2. Dominant Goods Price Range:** The graph prominently showcases that the highest number of applicants can be found within the price range of 200,000 to 400,000 for goods in the AMT\_GOODS\_PRICE\_GROUP.

**3. Lowest Applicant Count:** Conversely, the graph demonstrates that the smallest applicant count falls within the price range of 800,000 to 1,000,000 for goods within this grouping.



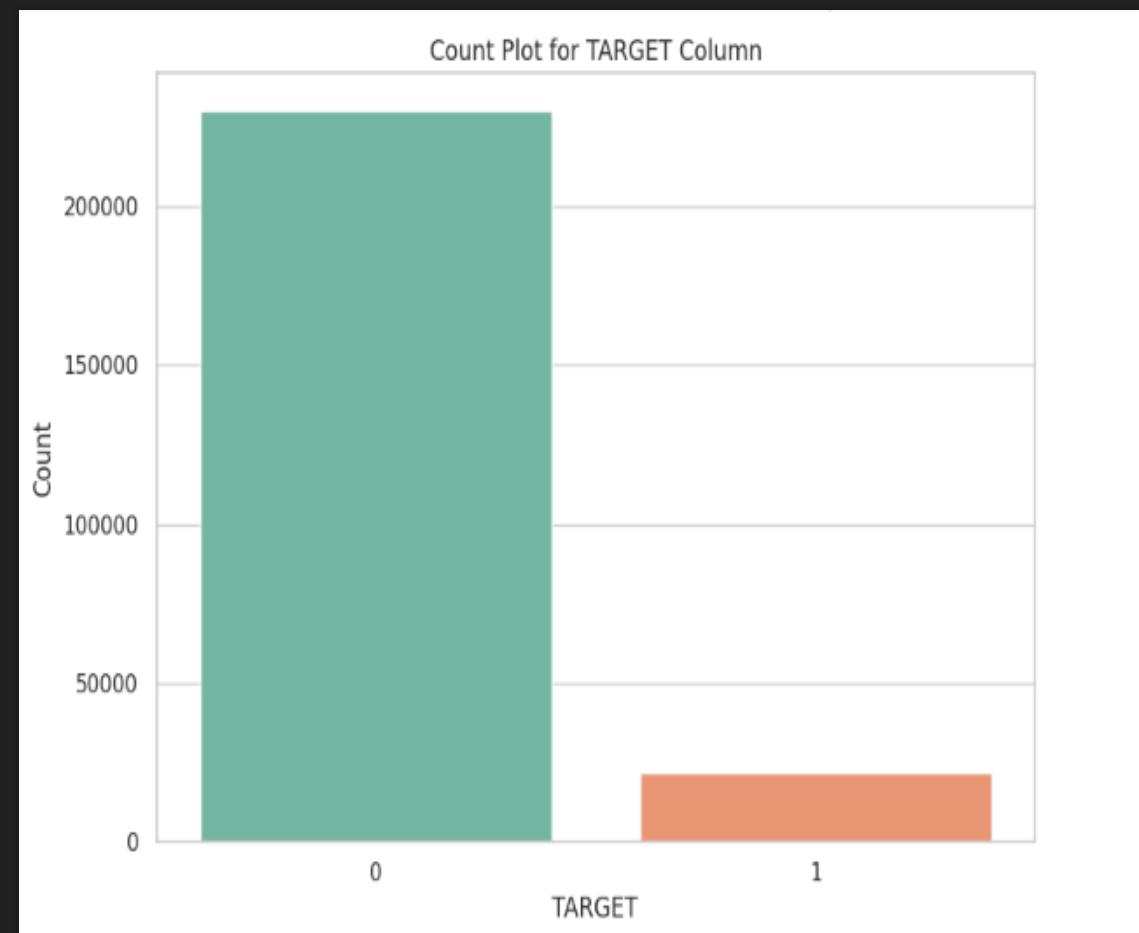
## Insights from Age Data Distribution Analysis

- 1. Age Distribution Analysis:** When examining the count graph for the "Age\_in\_years" category, a notable pattern emerges.
- 2. No Applicants Below 20:** Interestingly, the graph reveals that there are no applicants aged between 0 and 20 years.
- 3. Dominant Age Group:** The most significant observation is that the majority of applicants belong to the age category of 30 to 40 years, reflecting the highest number of applicants within this age range.

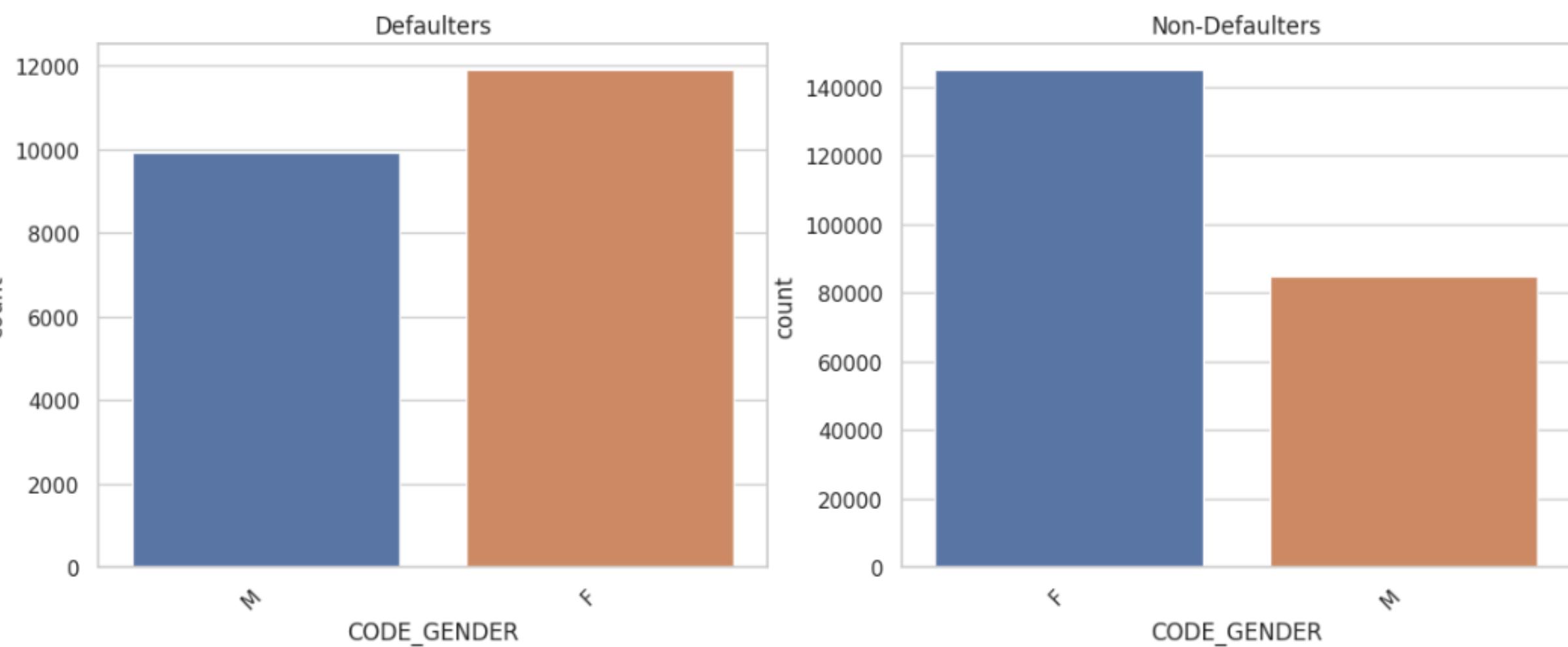


# Data Imbalance

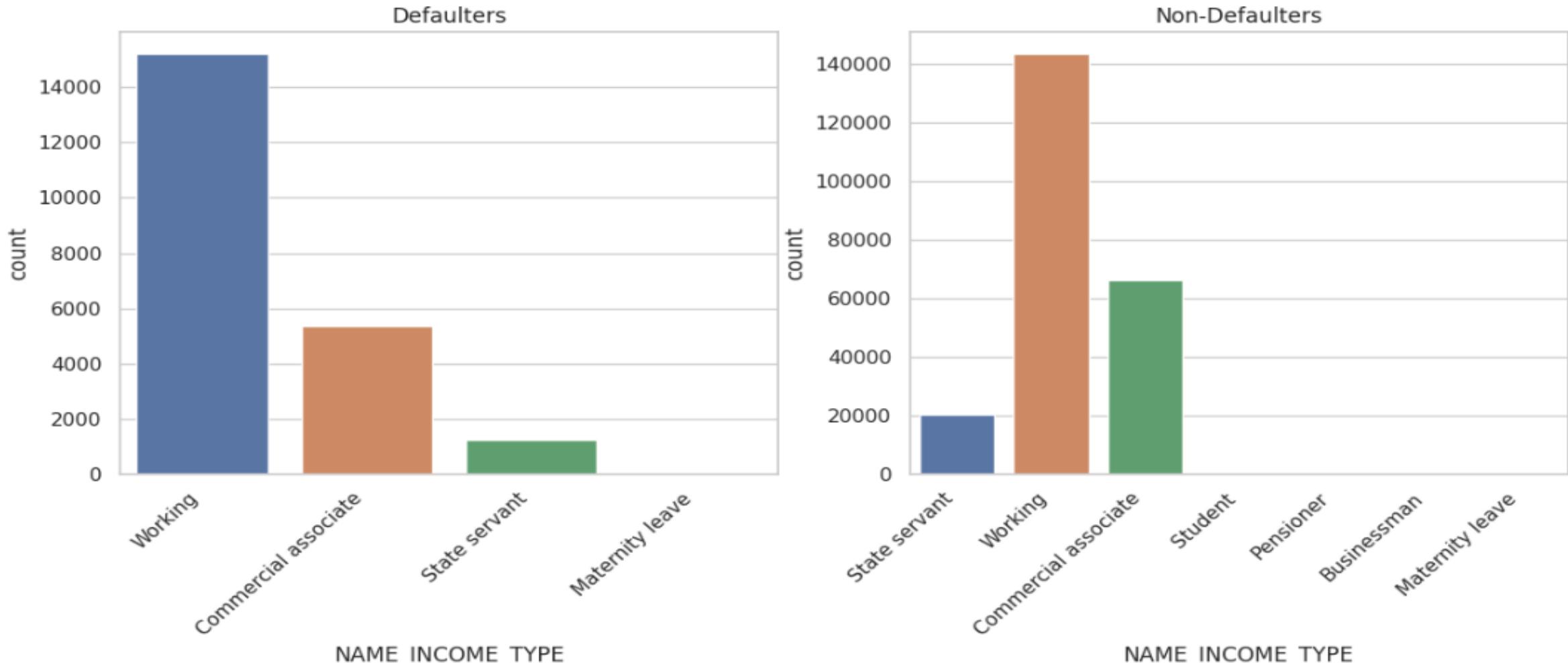
- **Understanding Data Imbalance:** Data imbalance occurs when there's a significant disparity in the number of instances representing different categories within a dataset.
- **Identification of Data Imbalance:** The analysis of the target column indicates the presence of data imbalance, where the target variable takes on values of 0 or 1.
- **Data Split:** To address this, we divided the data frame into two separate subsets: one for Target 1 and another for Target 0.
- **Target Variable 1:** Target variable 1 represents clients with payment difficulties. These individuals had late payments for at least one of the first Y installments of the loan, with a delay of more than X days, as defined in our sample.
- **Target Variable 0:** Target variable 0 encompasses all other cases, signifying clients without payment difficulties.
- **Data Imbalance in the Target Variable:** The target variable highlights a data imbalance, where a value of 1 corresponds to clients who encountered payment difficulties, characterized by late payments on certain early installments. Meanwhile, a value of 0 encompasses all other cases.
- **Imbalance Ratio Calculation:** After performing calculations, it was determined that the data imbalance ratio stands at 10.55.
- **Separate Analysis:** We conducted separate data analyses for these two datasets, involving univariate and bivariate analyses as well as correlation graph plotting, to gain a comprehensive understanding of each dataset independently.



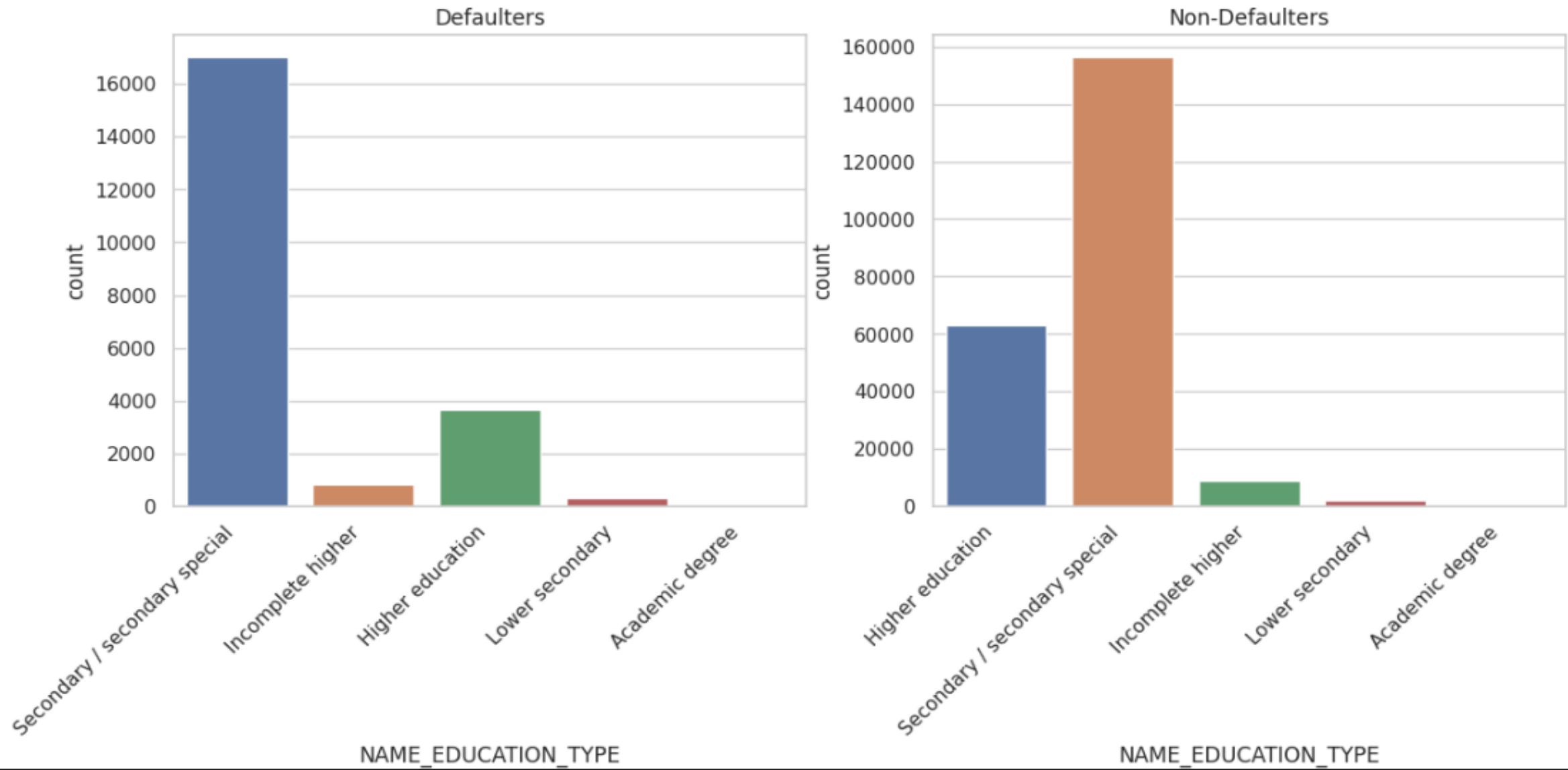
# Segmented Univariate analysis and Bivariate analysis



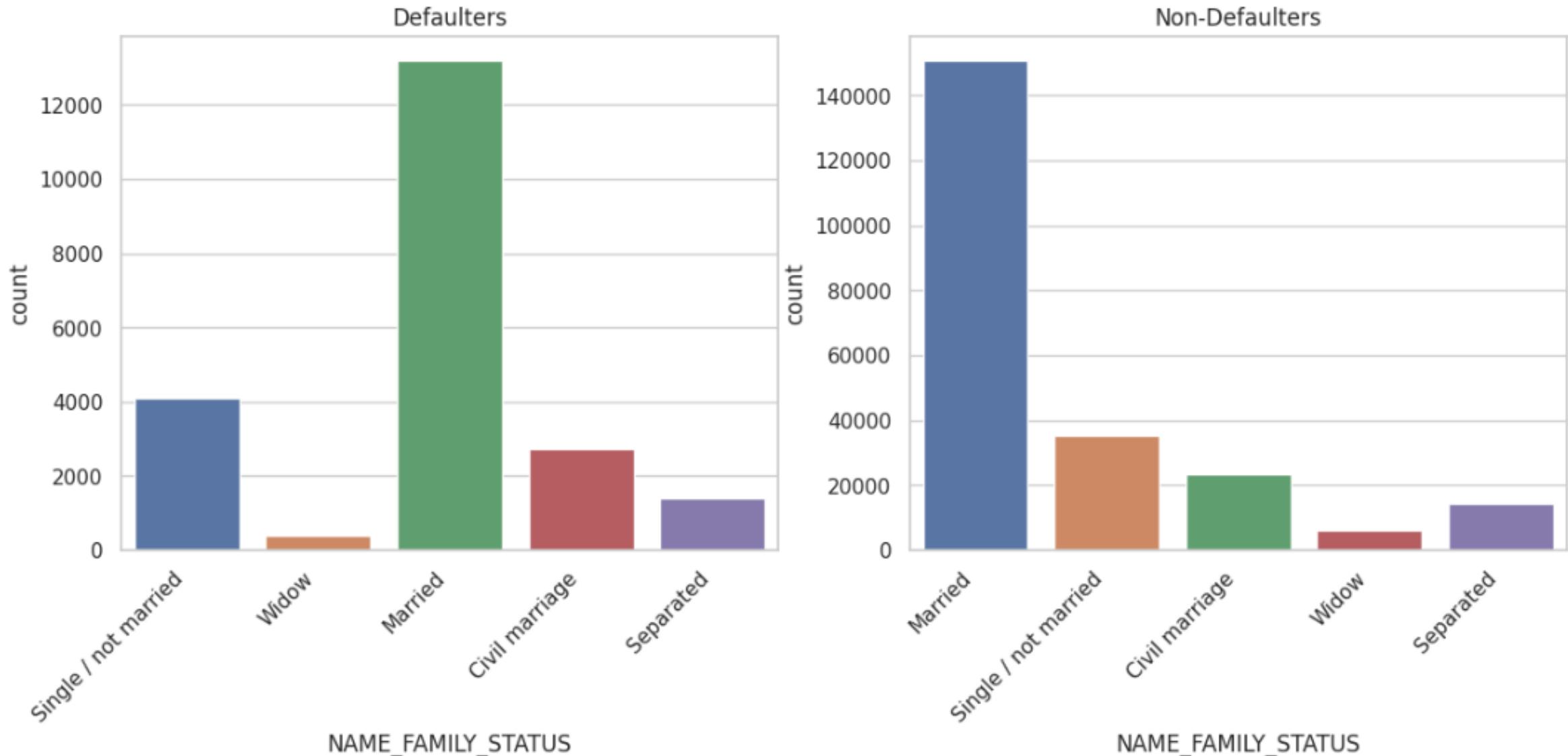
- The plot clearly illustrates that females have a higher ratio in both Target0 and Target1 groups.
- However, within Target0 (Non-Defaulters), females have a higher ratio than in Target1 (Defaulters).
- Consequently, it is advisable for the bank to prioritize granting loans to females, as they exhibit a higher likelihood of successful payments.
- In contrast, males have a higher ratio in Target1 compared to Target0, but their representation among Non-Defaulters is lower than among Defaulters.



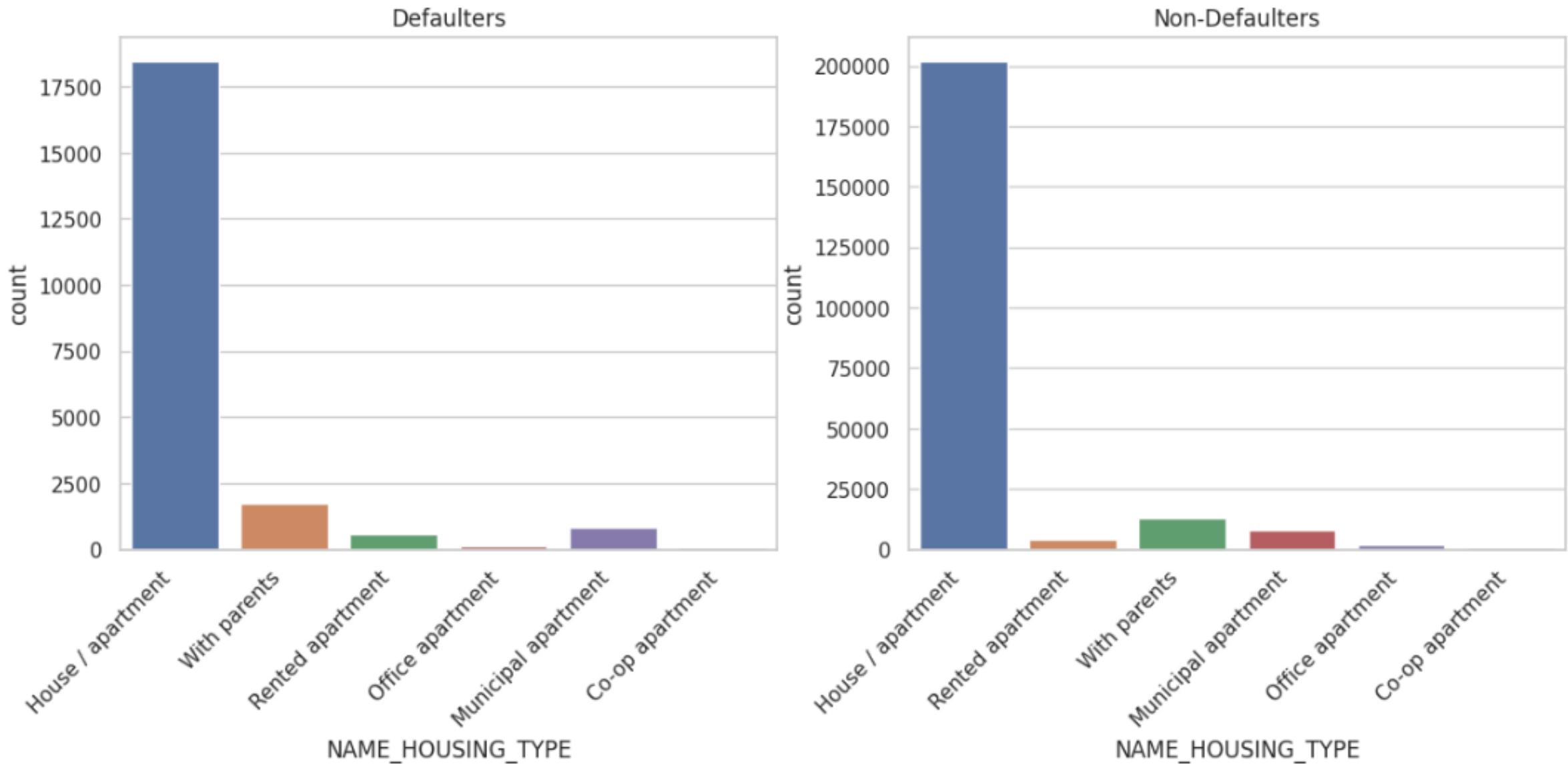
- **Income Types in Both Target Groups:** In both Target 0 (Non-Defaulters) and Target 1 (Defaulters), the prevalence of "working" income type is notably higher.
- **Low Frequency of State Servants:** It's worth noting that the frequency of state servants is comparatively low. Therefore, the bank should carefully evaluate loan applications from this category.
- **Loan Consideration for Working Class:** Considering the working class with "income type" as "working," the bank may need to offer loans with higher interest rates as this group seems to have a higher representation among both defaulters and non-defaulters.



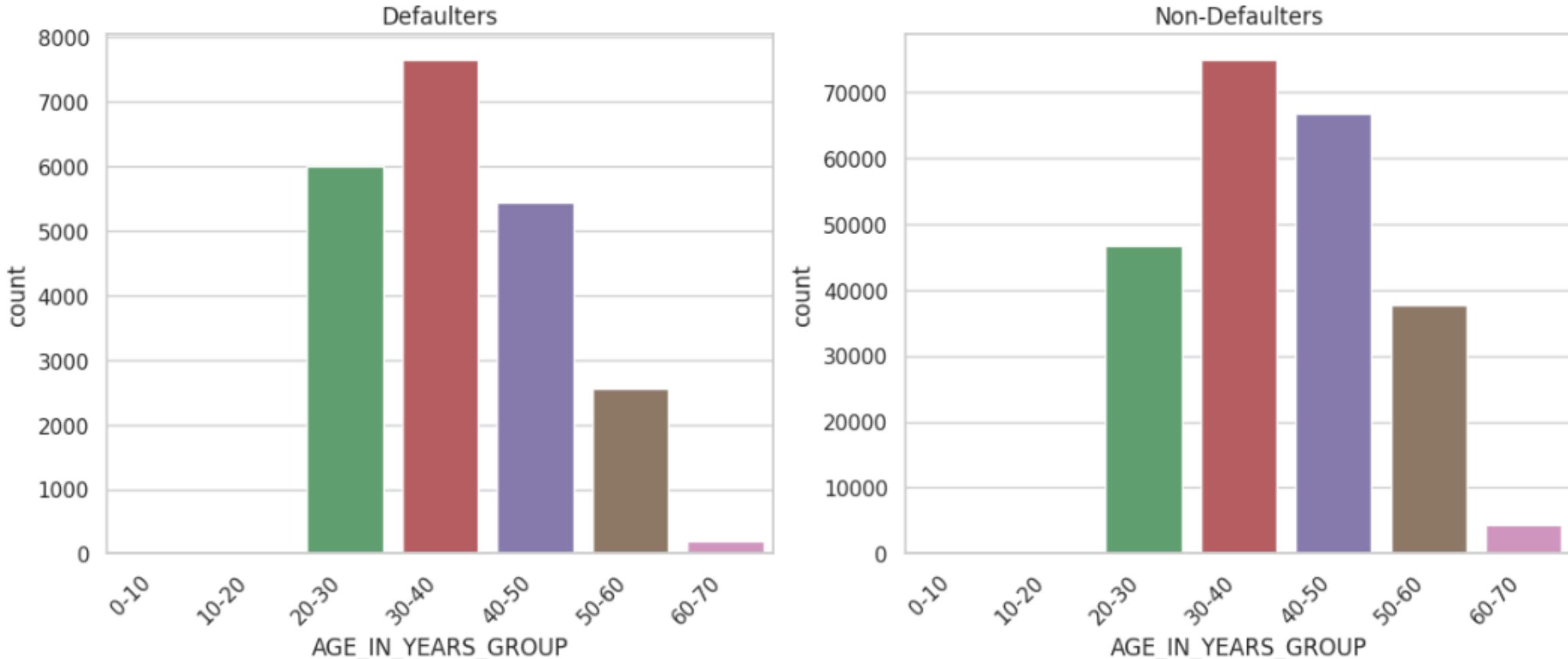
- **Education Type Insights:** Secondary/secondary special education types show higher representation in both Target 0 (Non-Defaulters) and Target 1 (Defaulters). However, it's essential to note that Target 0 has a significantly larger count than Target 1 in this category.
- **Consideration for Loans:** While secondary/secondary special education types may be considered, other factors should also be taken into account when evaluating loan applications.
- **Lowest Academic Degree Representation:** Academic degree education types have the lowest representation in both Target 0 (Non-Defaulters) and Target 1 (Defaulters).



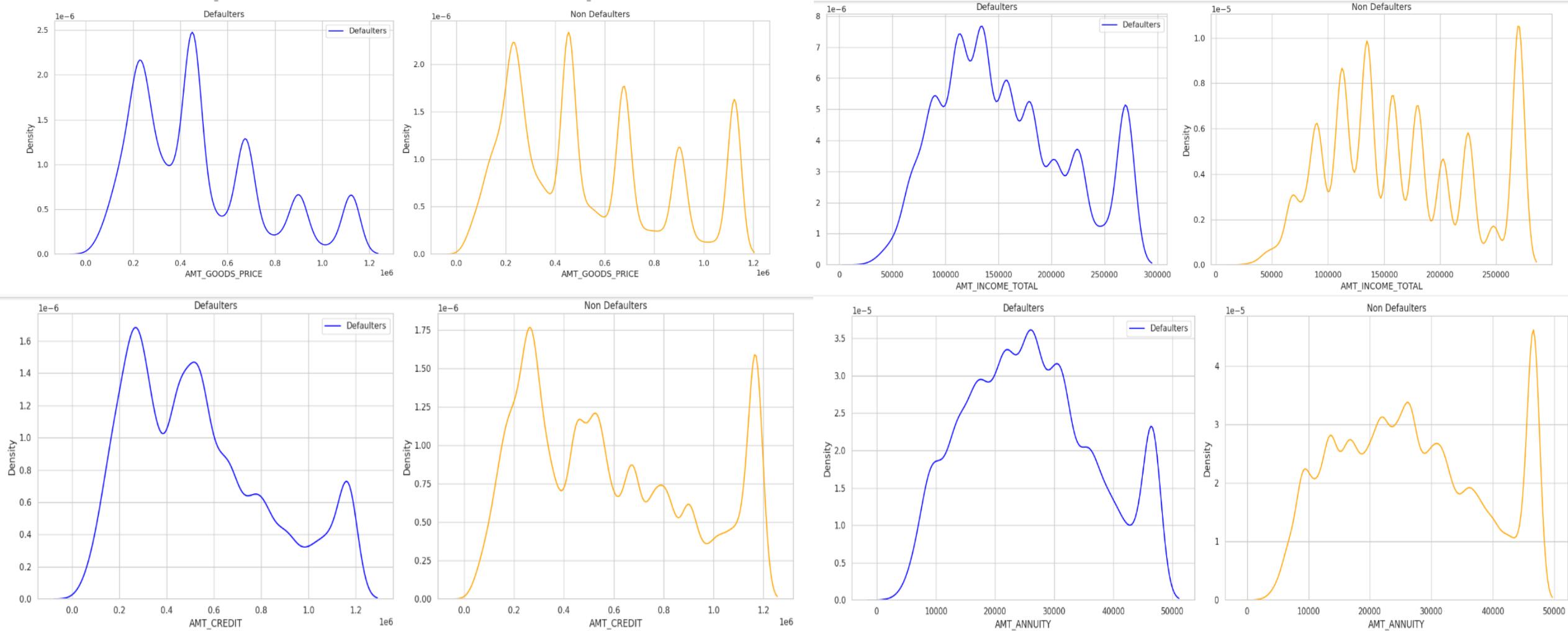
- **Family Status Insights:** Marriage family status is more prevalent in both Target 0 (Non-Defaulters) and Target 1 (Defaulters). However, it's important to note that Target 0 has a significantly larger count than Target 1 in this category.
- **Loan Evaluation:** While considering applicants with marriage family status may be relevant, it's advisable for the bank to factor in other criteria as well.
- **Lowest Representation:** Widow family status have the lowest representation in both Target 0 (Non-Defaulters) and Target 1 (Defaulters).



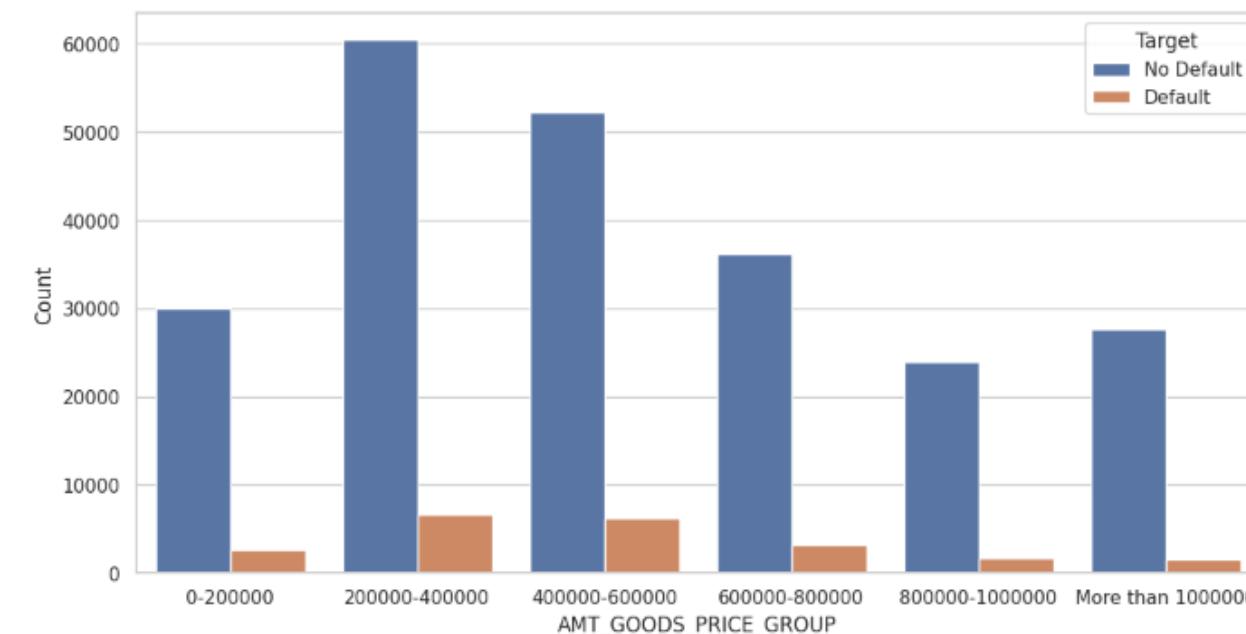
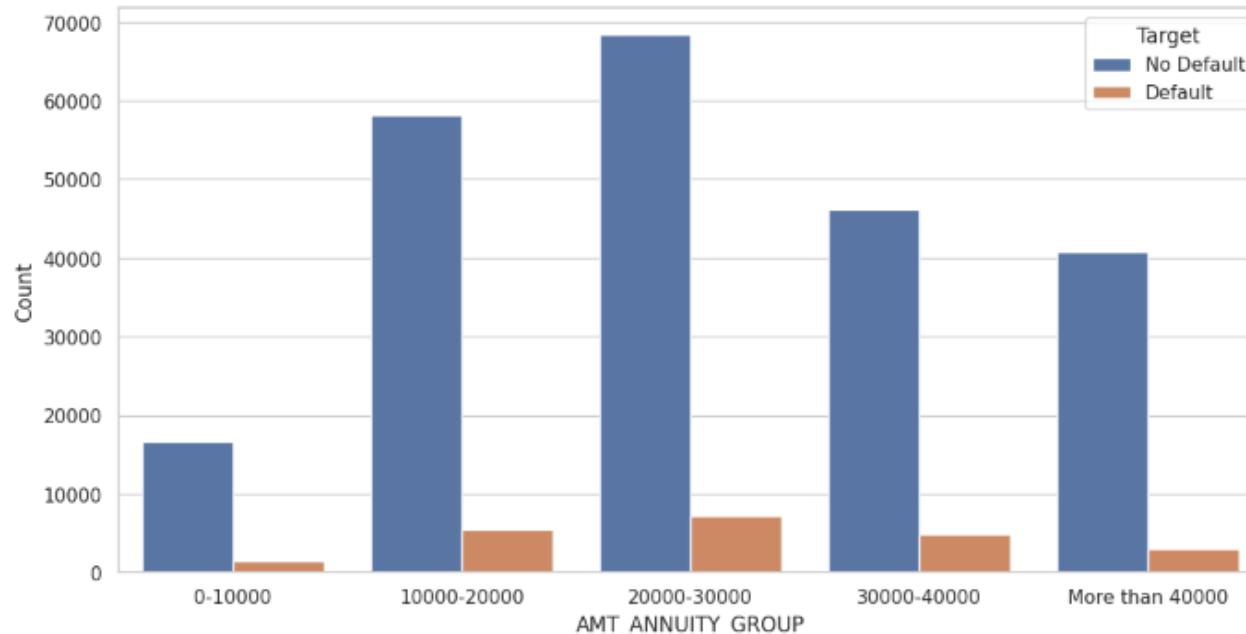
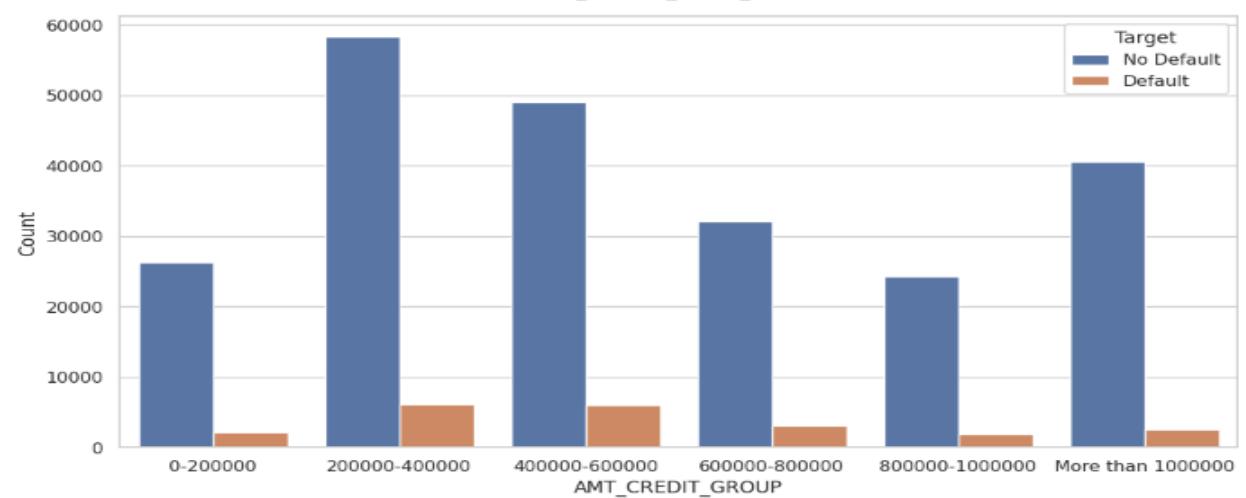
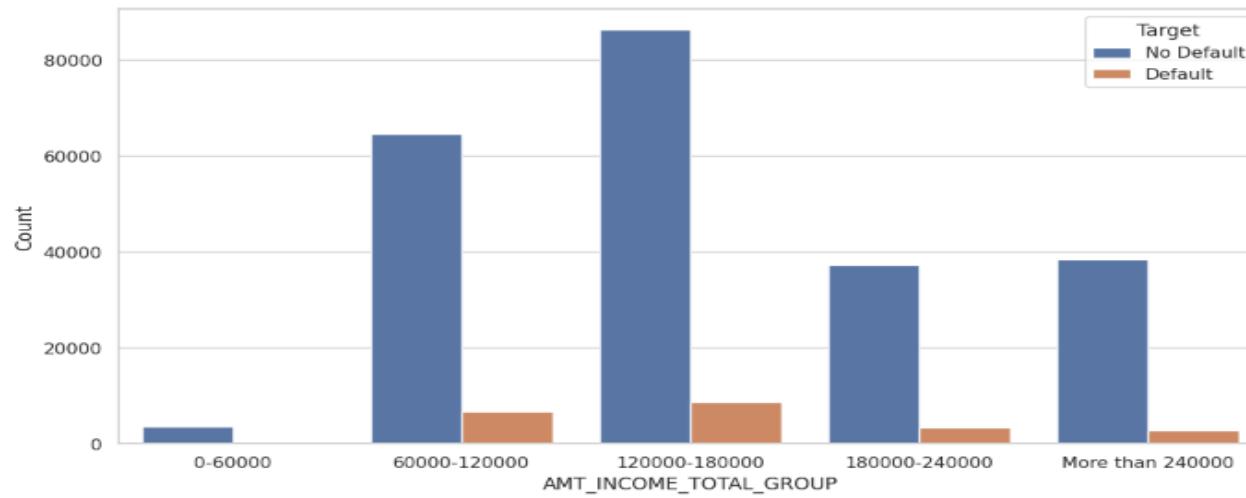
- **House Type Insights:** House/apartment house types show a higher presence in both Target 0 (Non-Defaulters) and Target 1 (Defaulters). However, it's crucial to note that Target 0 significantly outnumbers Target 1 in this category.
- **Loan Evaluation:** While applicants with house/apartment house types may be considered, it's advisable for the bank to take into account other relevant factors as well.
- **Lowest Representation:** Co-op apartment house types have the lowest representation in both Target 0 (Non-Defaulters) and Target 1 (Defaulters).



- **Age Group Insights:** Notably, there are no applications within the 0-20 age range. The 30-40 age group exhibits higher representation in both Target 0 (Non-Defaulters) and Target 1 (Defaulters). However, it's important to acknowledge that the count of Target 0 significantly surpasses that of Target 1 within this age bracket.
- **Loan Evaluation:** While considering applicants from the 30-40 age group may be pertinent, the bank should also take into consideration other relevant factors.
- **Lowest Representation:** The 60-70 age group has the lowest representation in both Target 0 (Non-Defaulters) and Target 1 (Defaulters).



- **Financial Group Representation:** Certain financial groups are represented in both Target 0 (Non-Defaulters) and Target 1 (Defaulters).
- **Density Plot Analysis:** Utilizing density plots for both Target 0 (Non-Defaulters) and Target 1 (Defaulters), we can perform a comprehensive analysis of these groups.



- Common Financial Group Representation:** Certain financial groups are shared between Target 0 (Non-Defaulters) and Target 1 (Defaulters).
- Count Plot Analysis:** We've utilized count plots for both target groups to streamline data examination, enabling us to assess differences in the number of defaulters.

# Plotted several bivariate graphs to conduct our analyses

## **Default Rates by Education and Gender:**

- The highest default rates are observed among males with lower secondary education, followed by those with secondary/secondary special education.
- Default rates are notably higher among males compared to females, particularly within the age range of 20 to 30 years.

## **Total Income vs. Gender:**

- Looking at the 'AMT\_INCOME\_TOTAL\_GROUP' vs. Gender graph, males exhibit a higher default rate compared to females, especially in the income range of 0 to 60,000.

## **Credit Amount vs. Gender:**

- Analyzing the 'AMT\_CREDIT\_GROUP' vs. Gender graph, a similar trend emerges with higher default rates among males, particularly in the credit amount range of 400,000 to 600,000.

## **Loan Annuity vs. Gender:**

- In the 'AMT\_ANNUITY\_GROUP' vs. Gender graph, once again, males show a higher propensity for default, particularly within the loan annuity range of 10,000 to 40,000.

### **Goods Price vs. Gender:**

- Shifting focus to the 'AMT\_GOODS\_PRICE\_GROUP' vs. Gender graph, males consistently exhibit higher default rates, especially in the goods price range of 200,000 to 600,000.

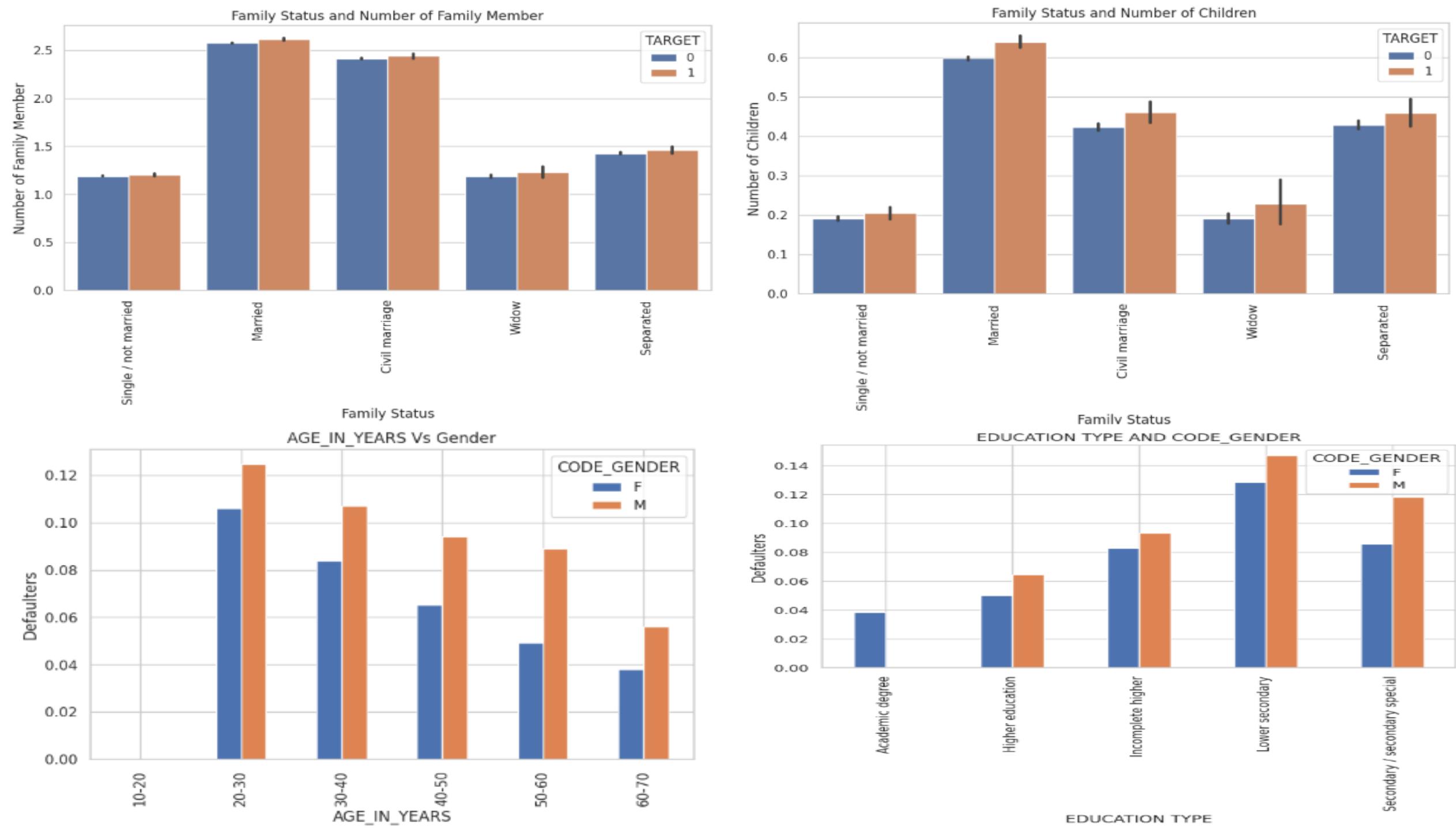
### **Family Size and Default Risk:**

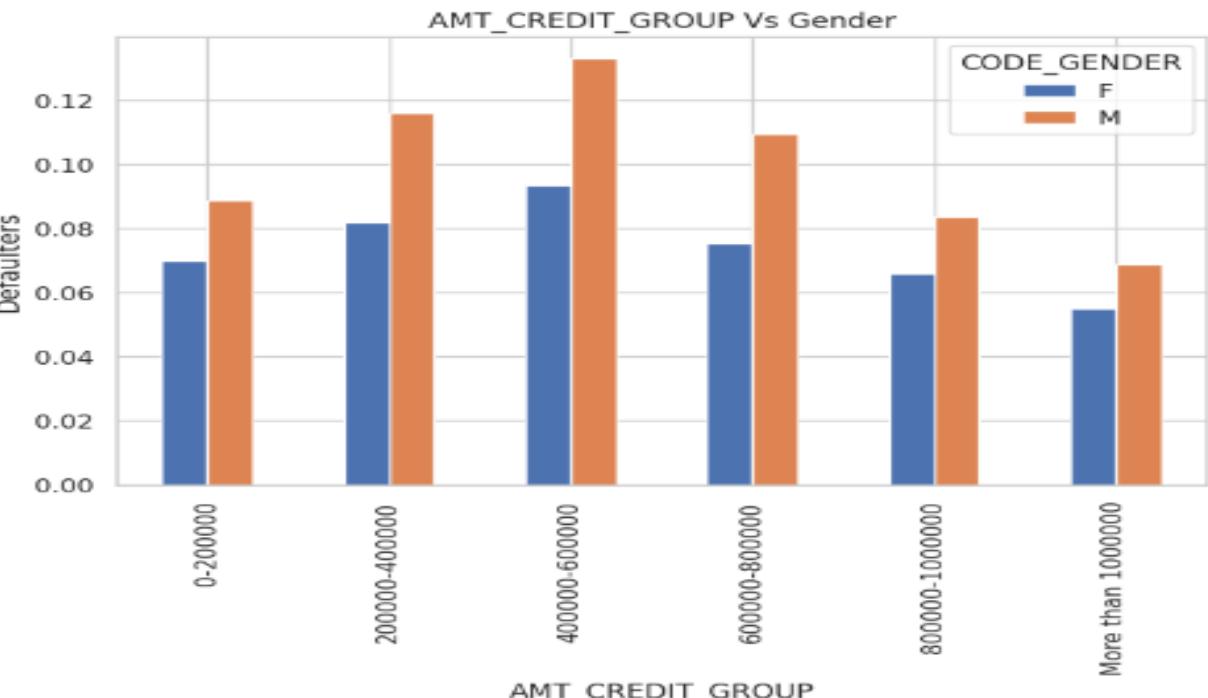
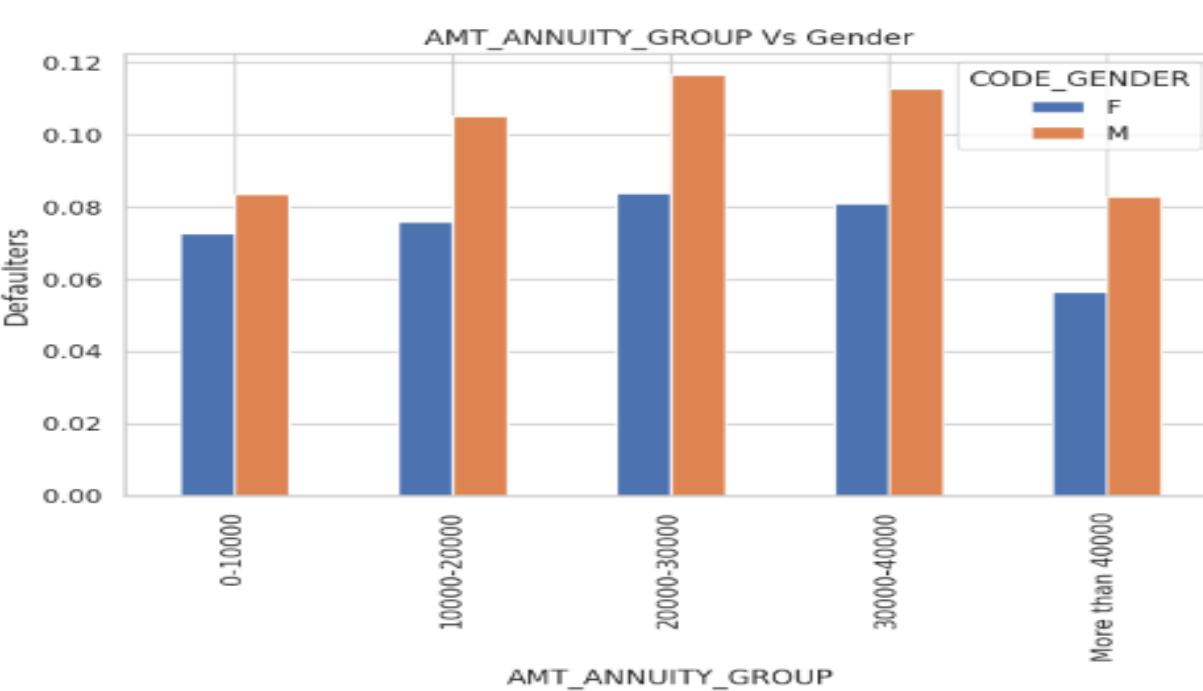
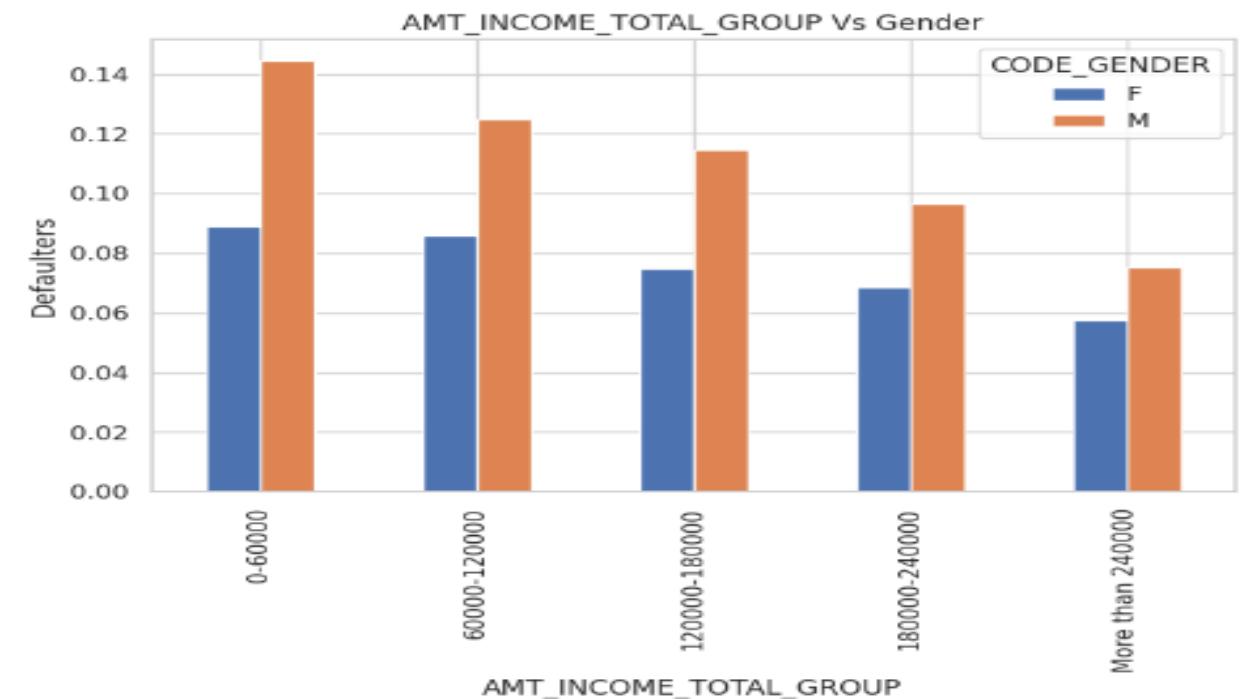
- Married clients with five or more children face an elevated risk of default, likely due to economic strains associated with supporting multiple dependents.
- Married clients with larger family sizes exhibit increased risks in both defaulting and non-defaulting categories, possibly due to the impact of additional family members on their economic situation.

### **Income Source and Family Size:**

- Individuals receiving income through Maternity Leave tend to have a higher likelihood of default when they have a greater number of family members.
- Similarly, individuals receiving income via Maternity Leave have a higher tendency to default when they have a larger number of children.

**Conclusion:** Across various financial aspects such as income, credit, loan annuity, and goods price, males consistently exhibit higher default rates compared to females. This trend holds true across different income and expense ranges. Additionally, family size, especially in cases of marriage and larger families, is associated with an increased risk of default.





# Correlation of Non Defaulters

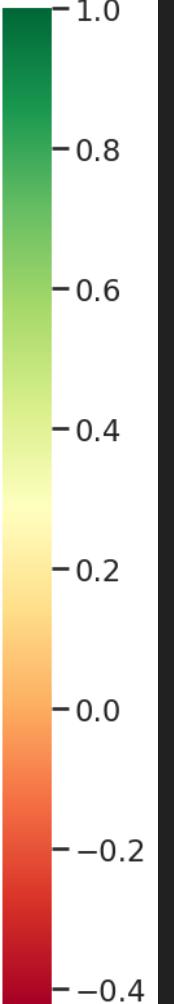
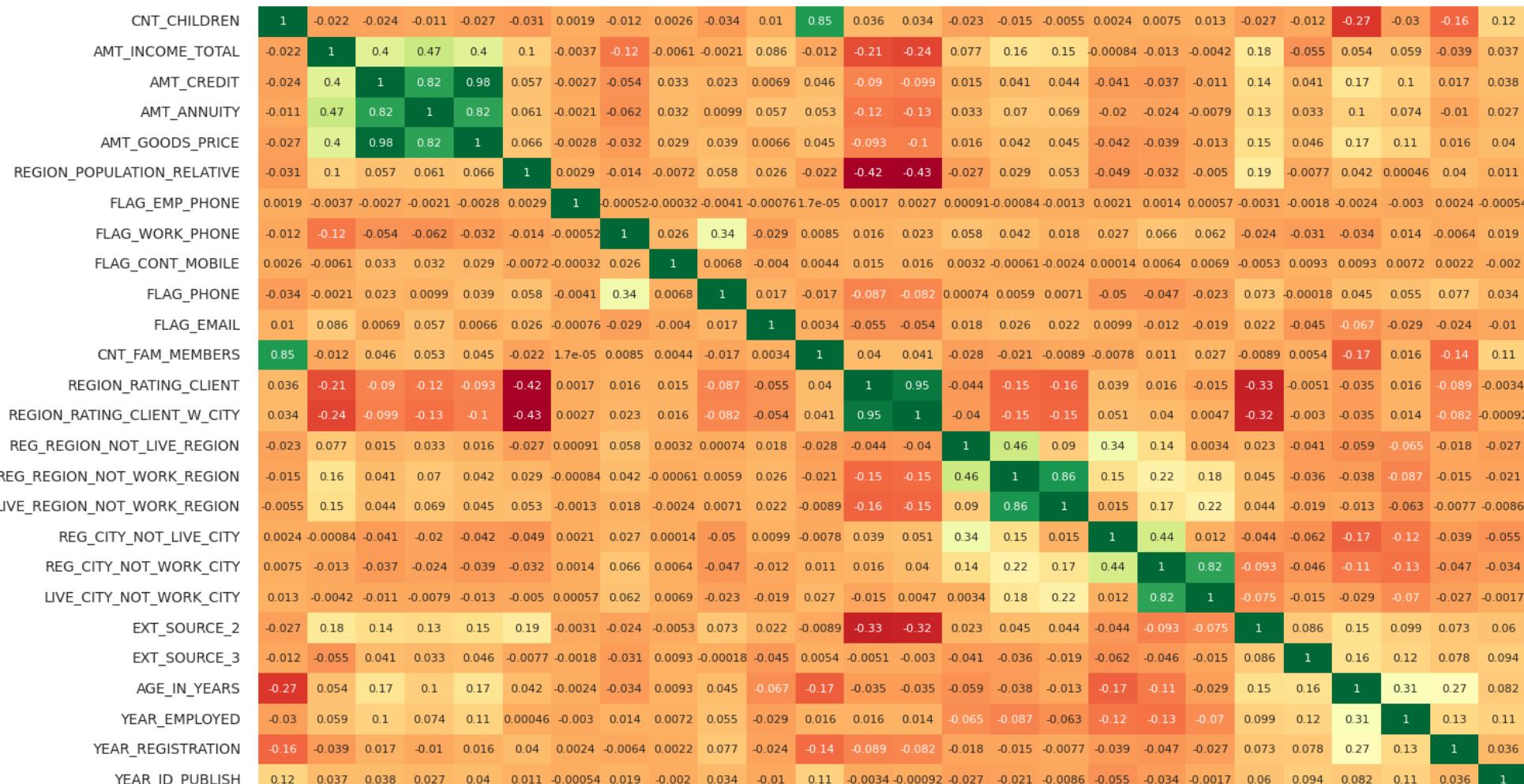
## Highly Correlated Variables:

- There's a strong correlation (0.98) between AMT\_CREDIT and AMT\_GOODS\_PRICE.
- REGION\_RATING\_CLIENT\_W\_CITY and REGION\_RATING\_CLIENT are highly correlated (0.95).
- CNT\_FAM\_MEMBERS and CNT\_CHILDREN show a notable correlation (0.85).
- AMT\_ANNUITY and AMT\_CREDIT exhibit a strong correlation (0.82).

## Relationships:

- Higher AMT\_CREDIT is associated with clients who have fewer children, while lower AMT\_CREDIT is associated with those having more children.
- Increased income is correlated with clients having fewer children, while lower income is linked to clients with more children.
- Higher income tends to be associated with clients having fewer family members, whereas lower income is related to clients with larger families.
- In densely populated areas, clients tend to have fewer children, as CNT\_CHILDREN is inversely proportional to REGION\_POPULATION\_RELATIVE.
- In densely populated areas, AMT\_CREDIT tends to be higher, indicating a direct relationship between AMT\_CREDIT and REGION\_POPULATION\_RELATIVE.
- AMT\_INCOME\_TOTAL is typically higher in densely populated areas, suggesting an inverse relationship between income and REGION\_POPULATION\_RELATIVE.

Correlation for Target 0



# Correlation of Defaulters

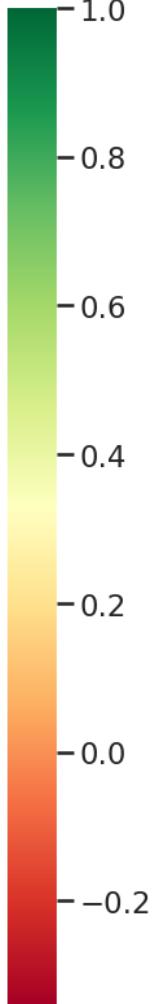
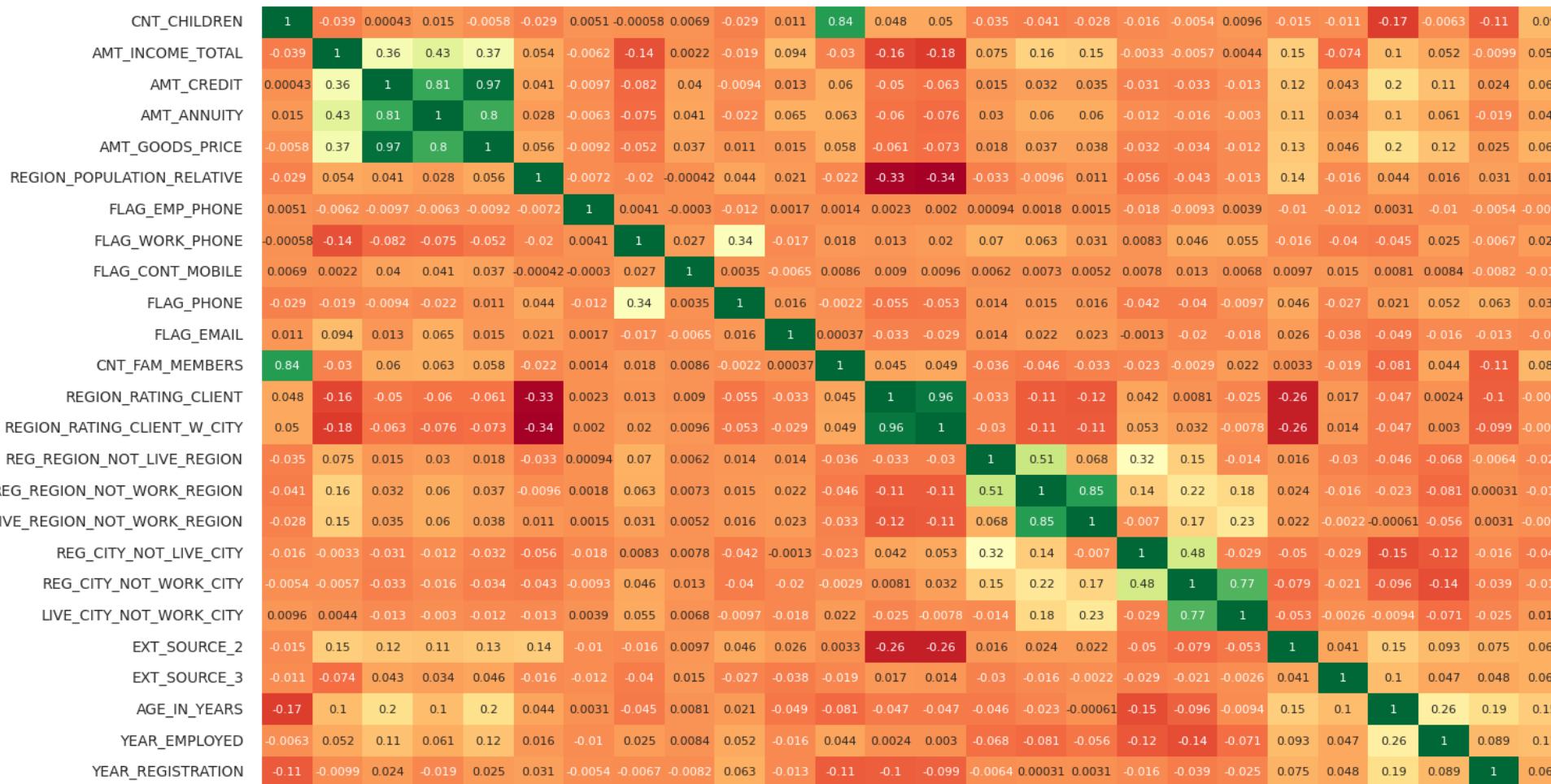
## Highly Correlated Variables:

- AMT\_CREDIT and AMT\_GOODS\_PRICE display a strong correlation of 0.97.
- REGION\_RATING\_CLIENT\_W\_CITY and REGION\_RATING\_CLIENT are highly correlated, with a coefficient of 0.96.
- CNT\_FAM\_MEMBERS and CNT\_CHILDREN show a notable correlation, registering at 0.84.
- AMT\_ANNUITY and AMT\_CREDIT exhibit a strong correlation of 0.81.

## Relationship Insights:

- Clients with higher AMT\_CREDIT tend to have fewer children, whereas those with lower AMT\_CREDIT tend to have more children.
- Increased income correlates with clients having fewer children, while lower income correlates with clients having more children.
- Higher income is linked to clients having fewer family members, while lower income is associated with clients having more family members.
- In densely populated areas, clients tend to have fewer children, as CNT\_CHILDREN inversely correlates with REGION\_POPULATION\_RELATIVE.
- Densely populated areas show a tendency for higher AMT\_CREDIT, suggesting a direct relationship between AMT\_CREDIT and REGION\_POPULATION\_RELATIVE.
- AMT\_INCOME\_TOTAL tends to be higher in densely populated areas, indicating an inverse relationship between income and REGION\_POPULATION\_RELATIVE.

Correlation for Target 1



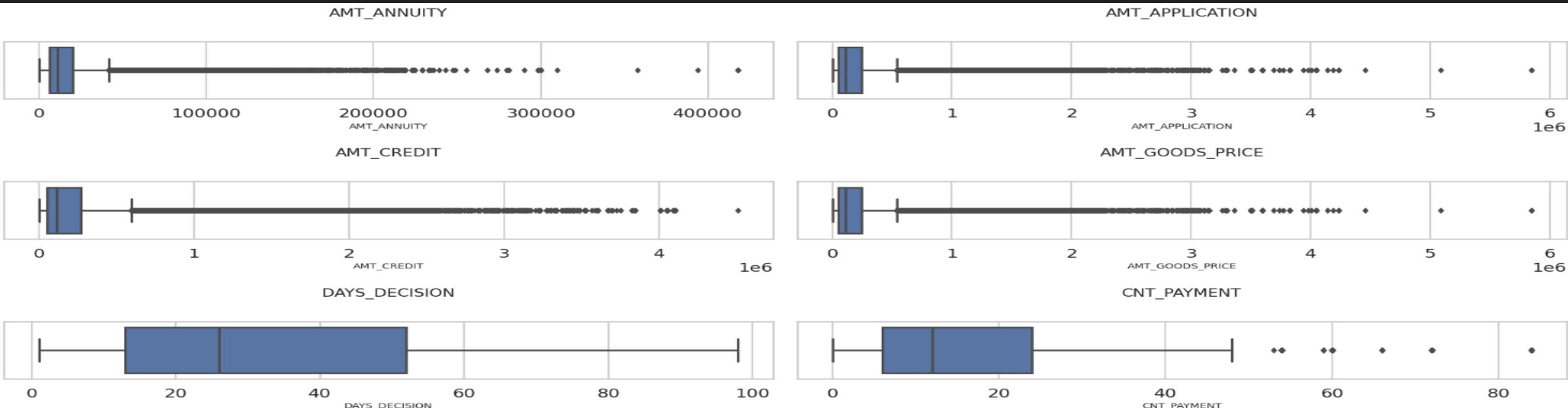
# Previous application

## Exploring Missing Data and Data Cleaning:

- We began by examining null values in each column and calculated the percentage of missing data.
- Columns with null values exceeding 30% were removed, while others with varying percentages of null values were assessed, ranging from 23% to 0.33%.
- Columns lacking significance were removed, and for the remaining null values, we opted for imputation using either the mean or mode as appropriate.
- Inspecting and addressing individual columns with missing data in a sequential manner.
- Additionally, we identified columns containing NAN or XNA entries and addressed them by either removing corresponding rows or replacing NAN, XNA, or XAP values with the mode.

## Outlier Detection and Mitigation Using Capping and Winsorization:

- Initial observation revealed the presence of outliers in numerous numerical columns, prompting concern for their impact on data analysis.
- To identify these outliers, we created box plots for numerical columns, providing a visual representation of extreme values.
- To mitigate the influence of outliers, we implemented two techniques: capping and winsorization.
- Capping involved replacing extreme values exceeding a predefined threshold with the threshold value, limiting their effect on subsequent analyses.
- Winsorization, as an alternative approach, replaced extreme outliers with the nearest non-outlier values, preserving the overall data distribution.
- Box plots were included to visualize outlier presence, emphasizing the necessity of addressing these data points.
- Furthermore, we resolved negative values in specific columns by converting them into their absolute counterparts to rectify this issue.

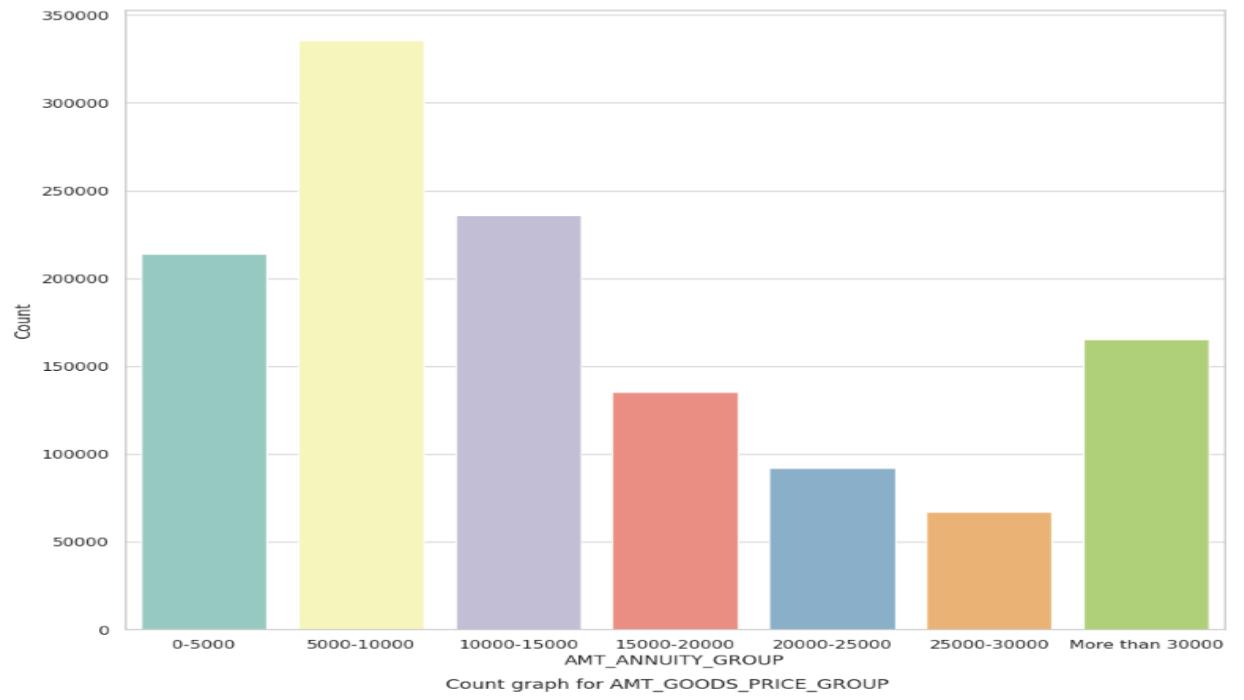


# Exploratory Data Analysis Insights for Previous Application Data

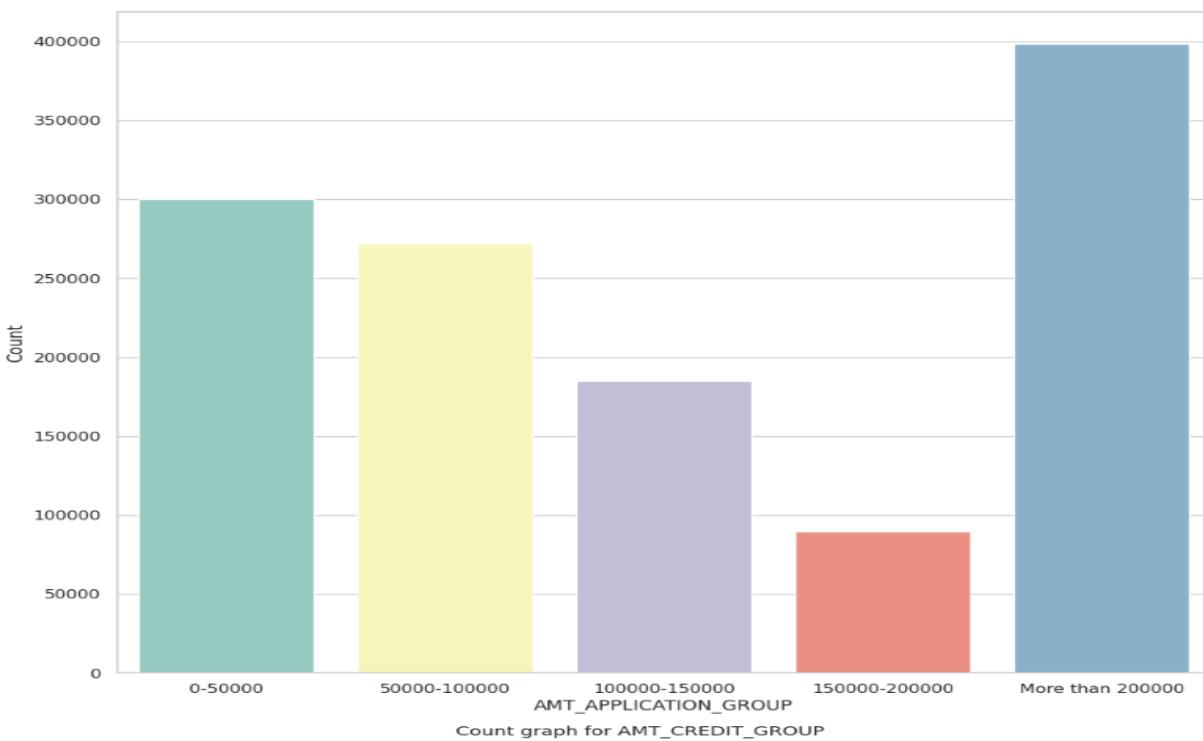
These observations were based on exploratory data analysis using various types of graphs for the previous application data.

- The graph revealed that 'DAYS\_DECISION' exhibited a negative correlation ratio with other variables.
- Assumption: Renaming the column from 'DAYS\_DECISION' to 'MONTHS\_DECISION' was suggested.
- Utilized Various Graph Types for Analyzing Previous Application Data, Presented Below:

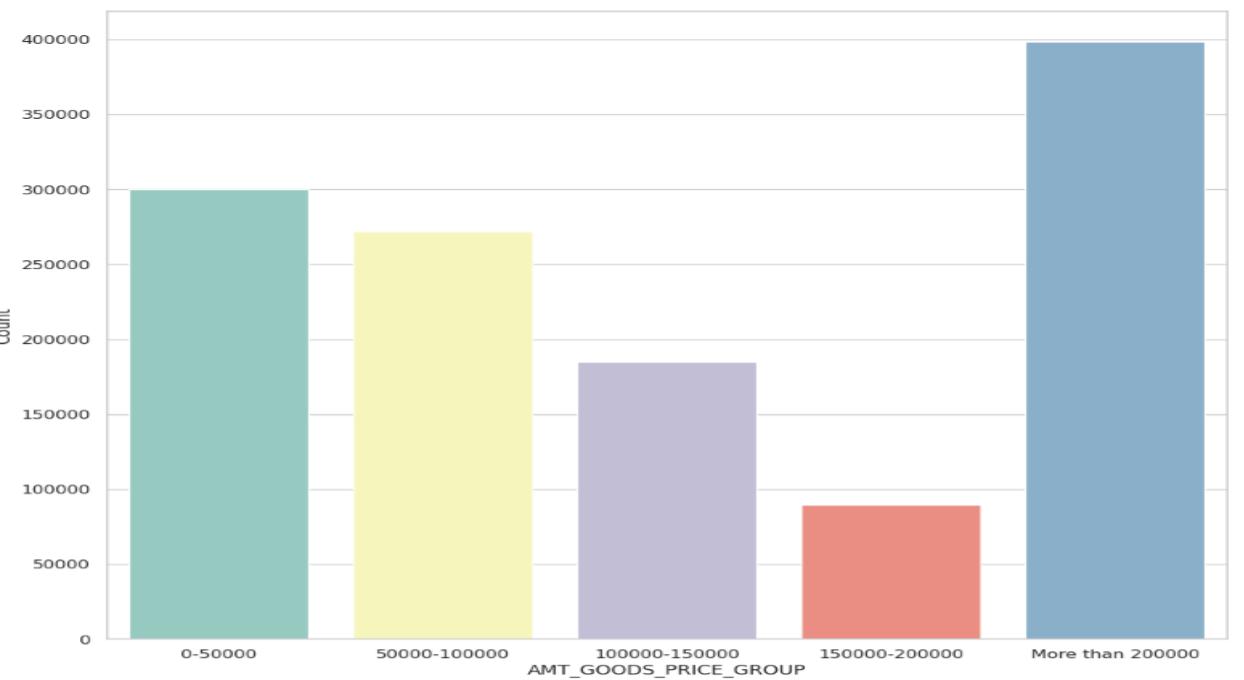
Count graph for AMT\_ANNUITY\_GROUP



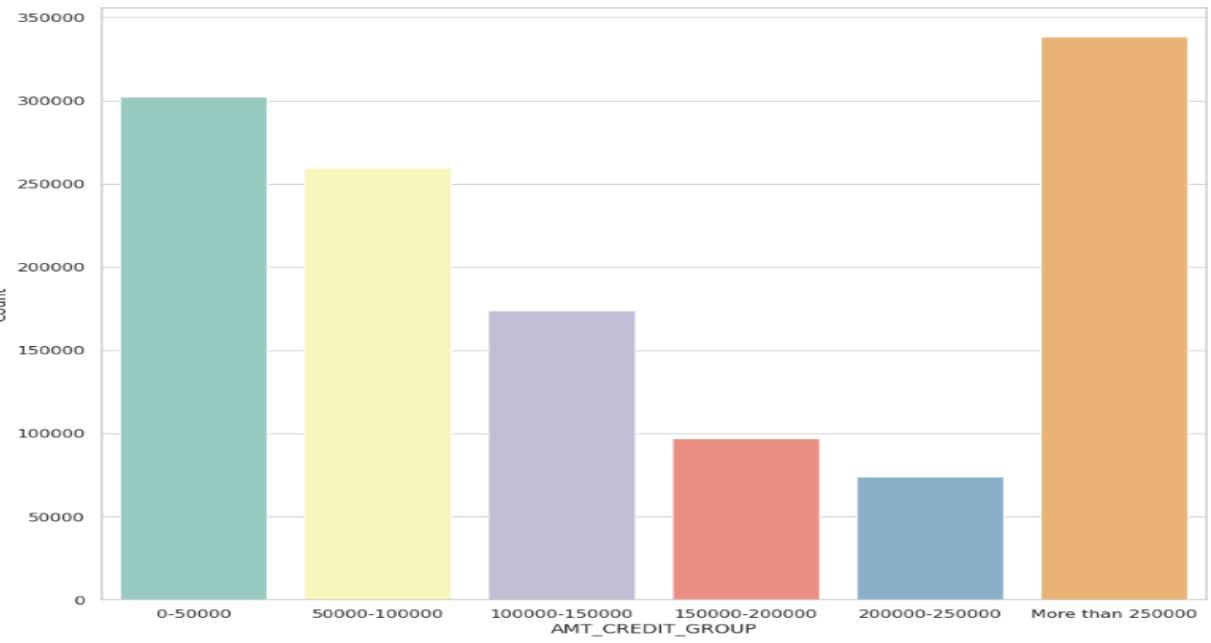
Count graph for AMT\_APPLICATION\_GROUP



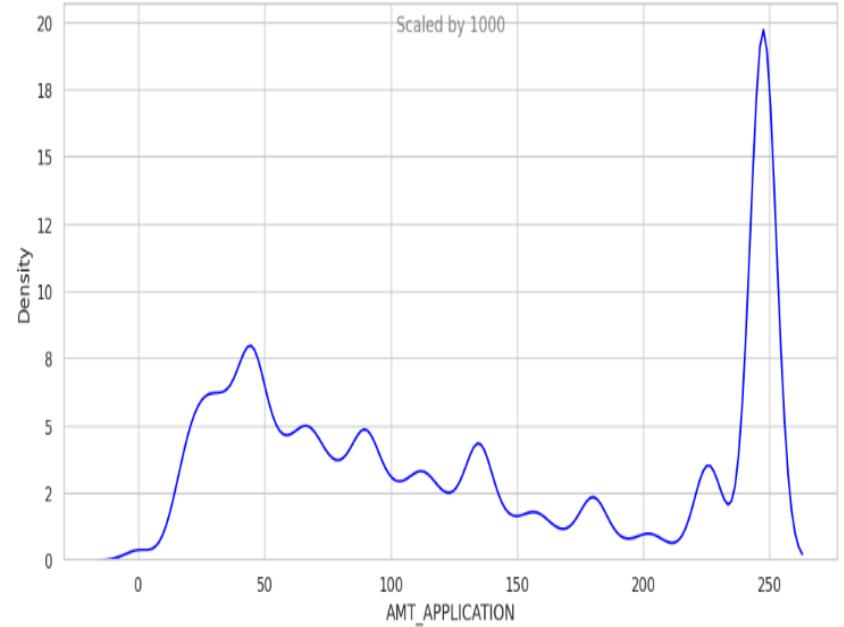
Count graph for AMT\_GOODS\_PRICE\_GROUP



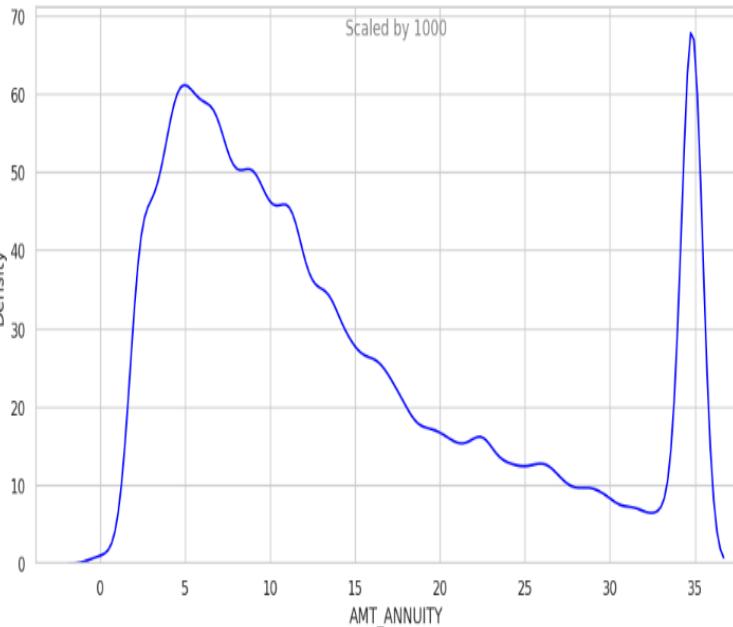
Count graph for AMT\_CREDIT\_GROUP



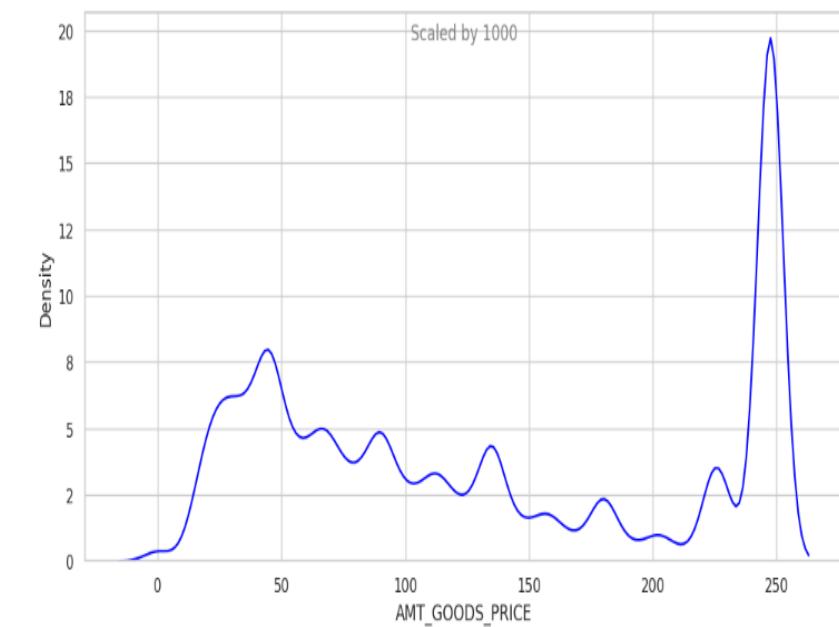
Density Plot for AMT\_APPLICATION



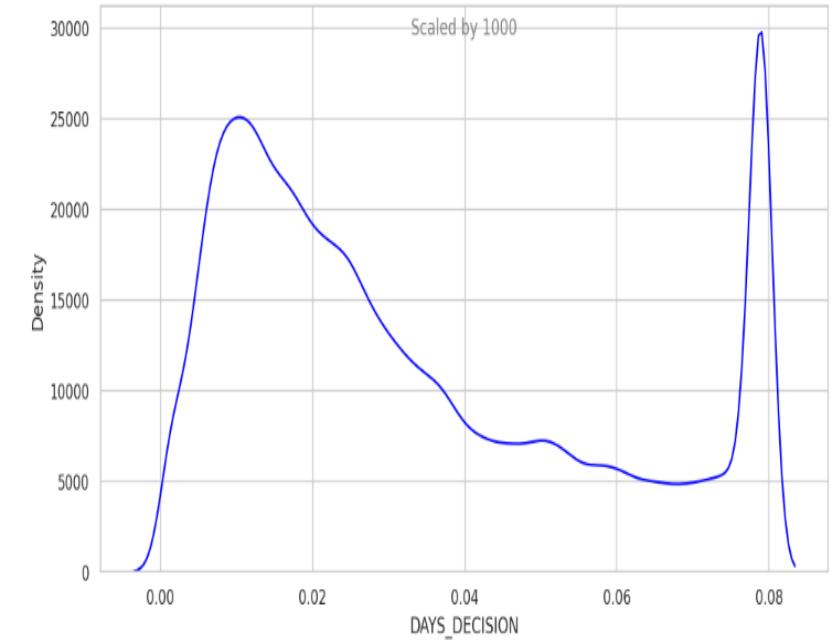
Density Plot for AMT\_ANNUITY



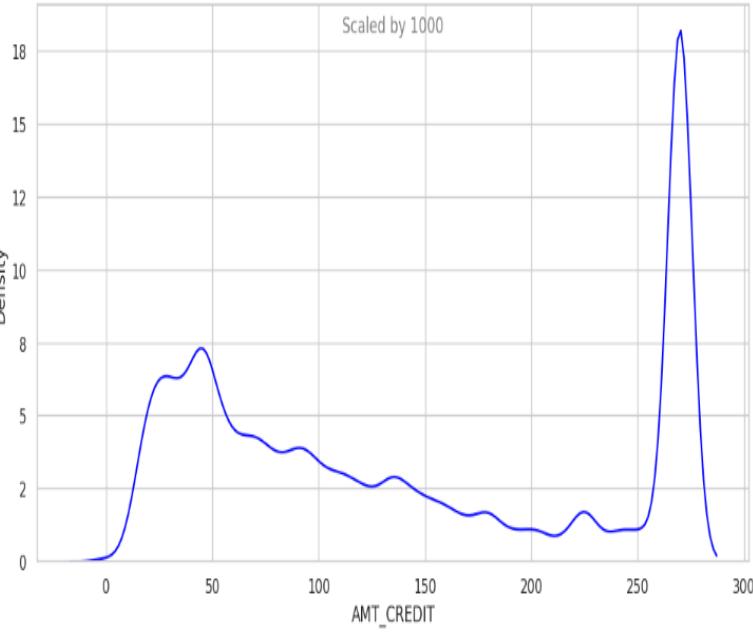
Density Plot for AMT\_GOODS\_PRICE



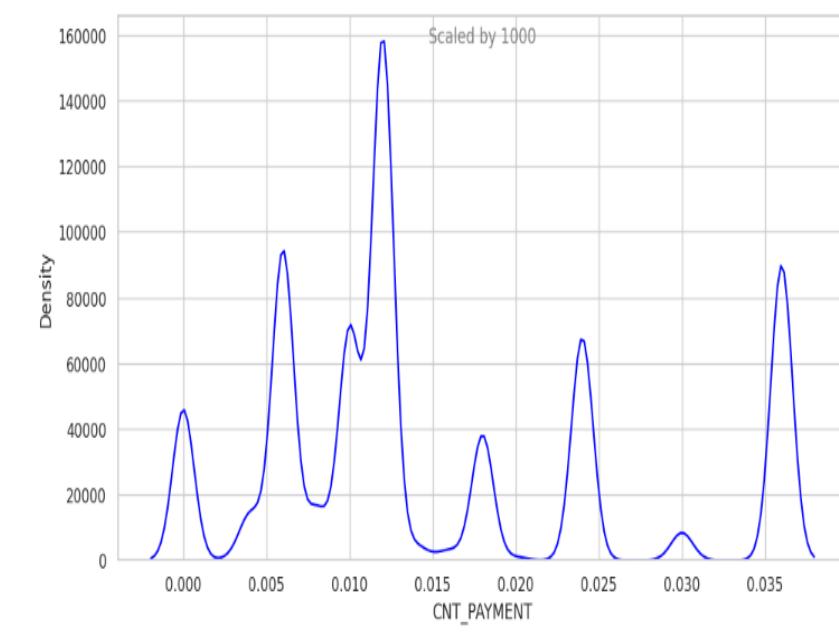
Density Plot for DAYS\_DECISION

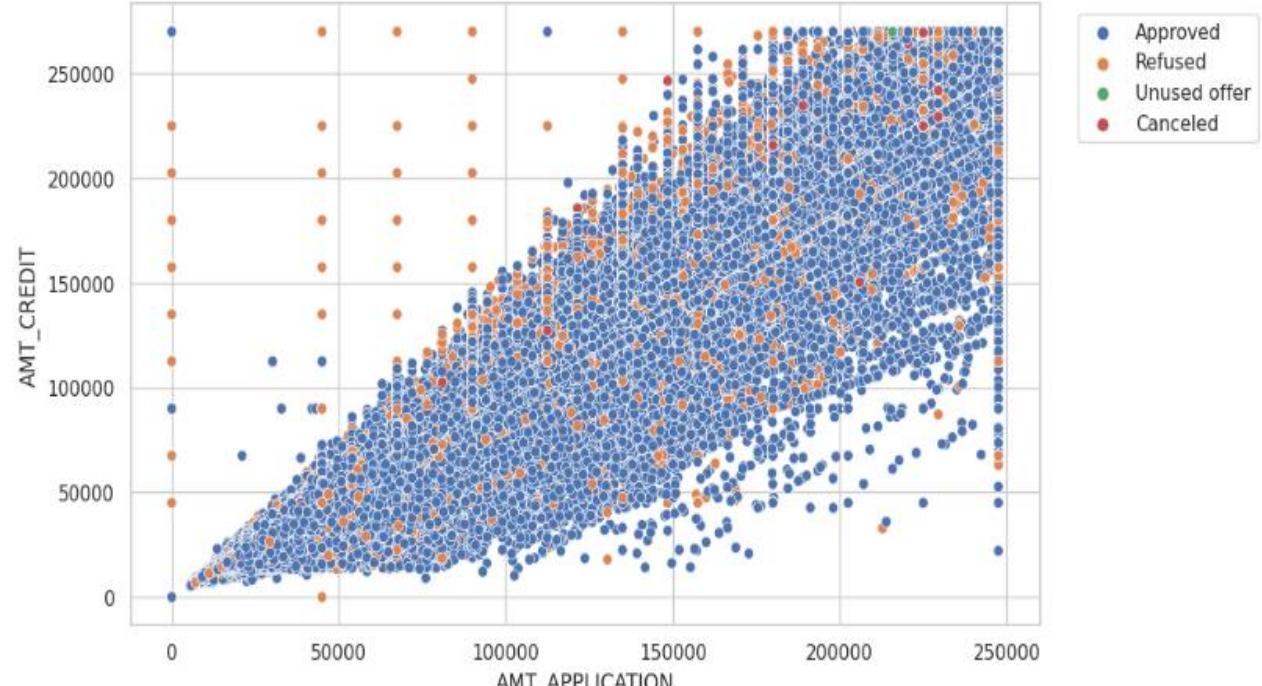
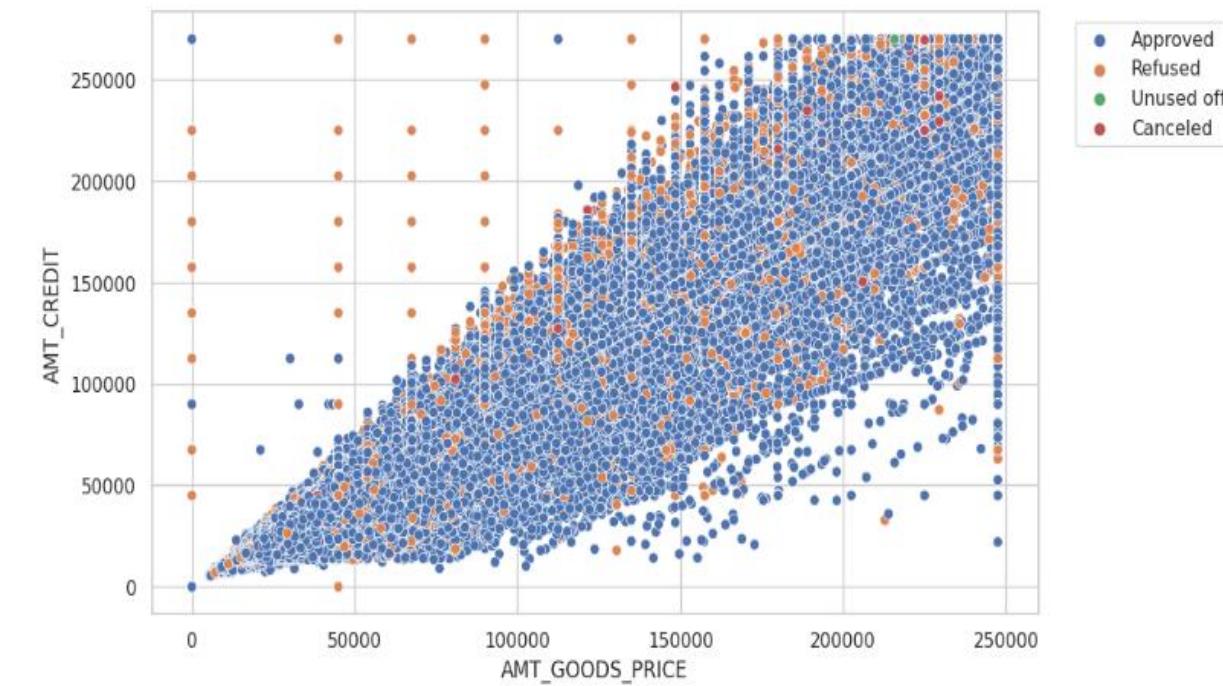
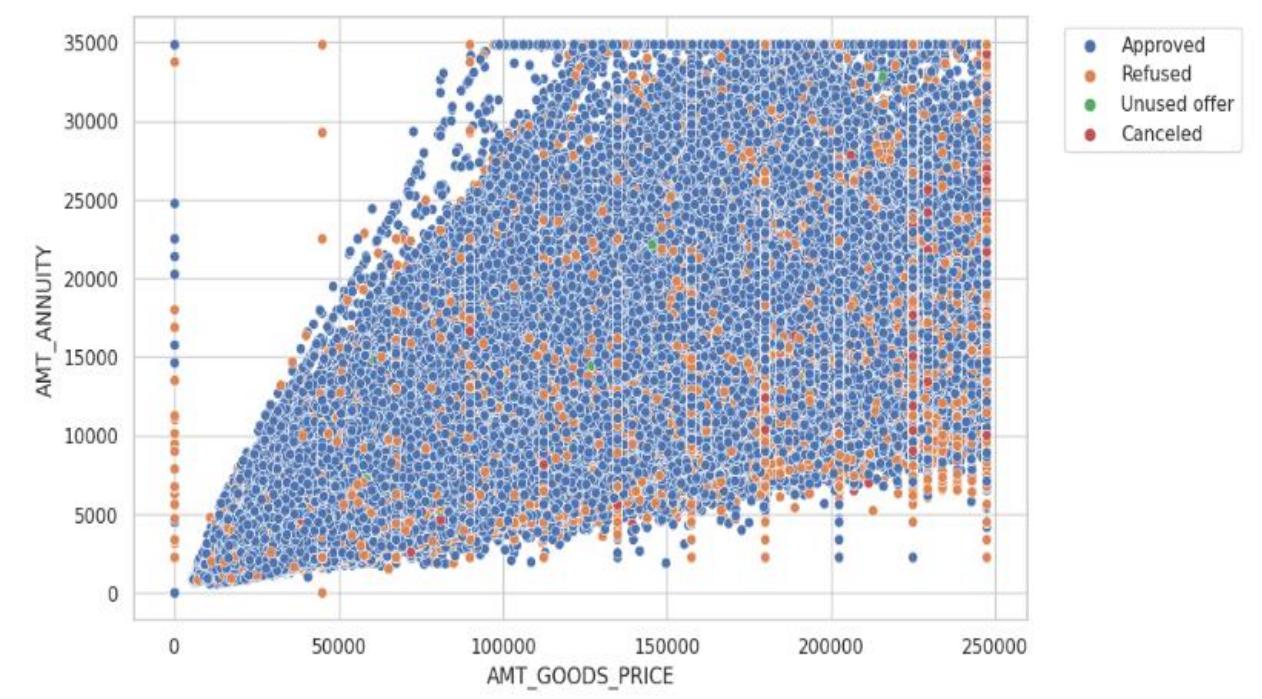
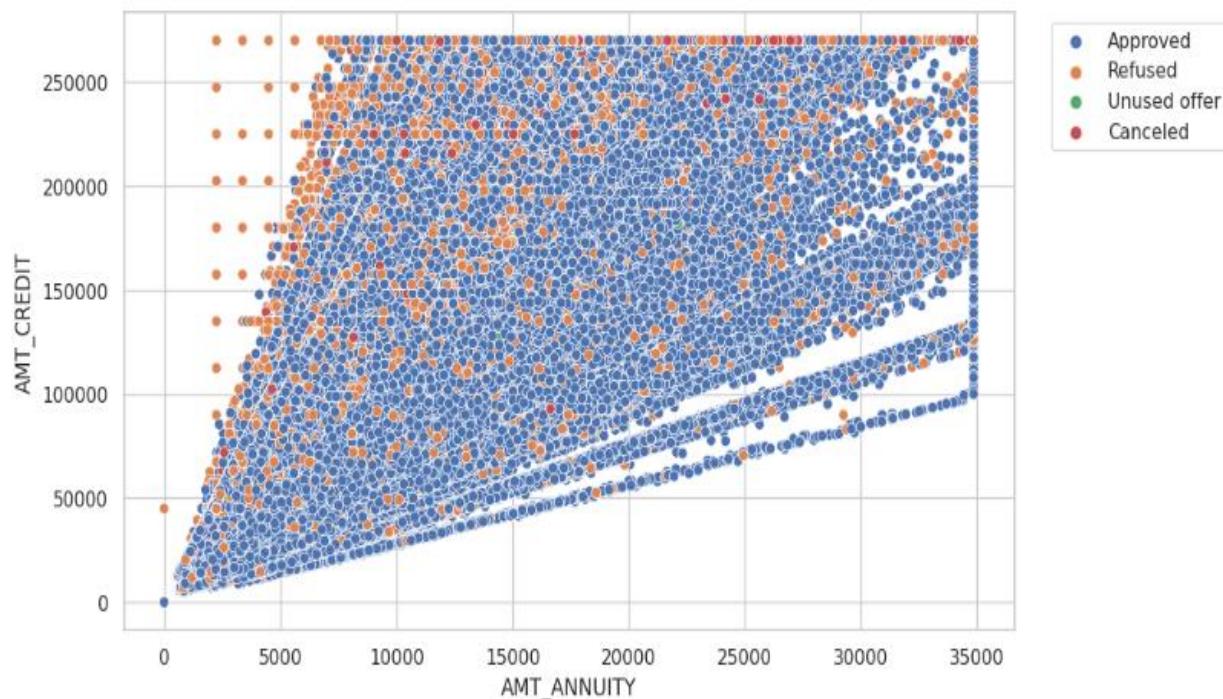


Density Plot for AMT\_CREDIT

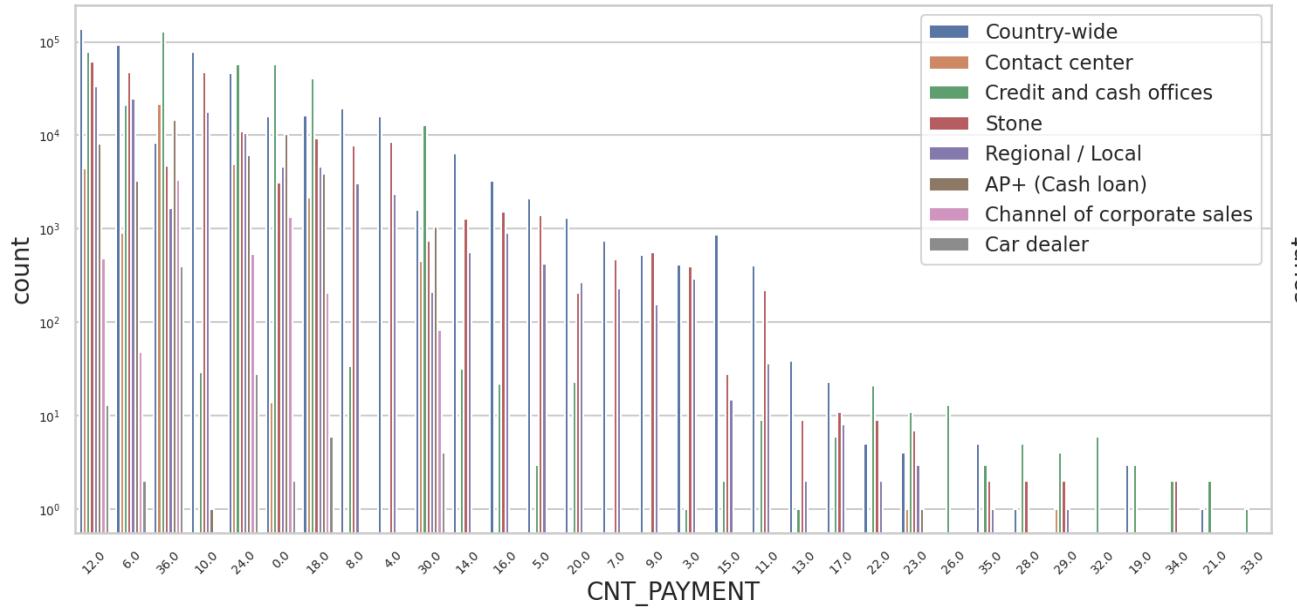


Density Plot for CNT\_PAYMENT

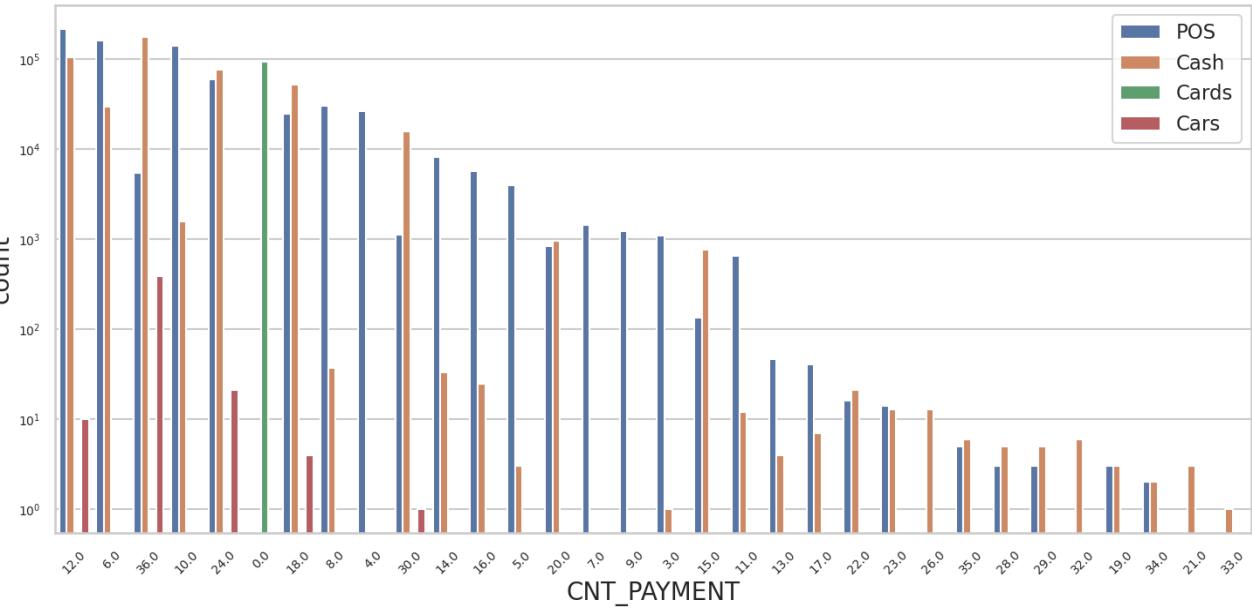




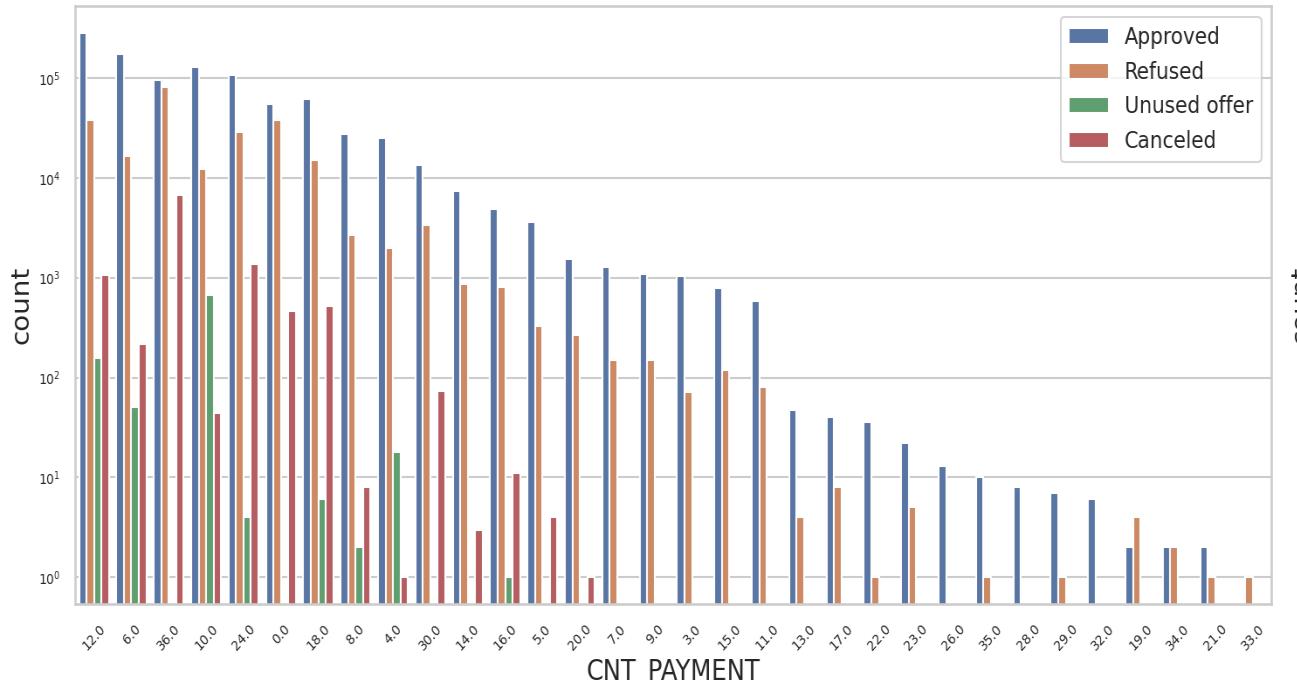
Distribution of CHANNEL\_TYPE with purposes



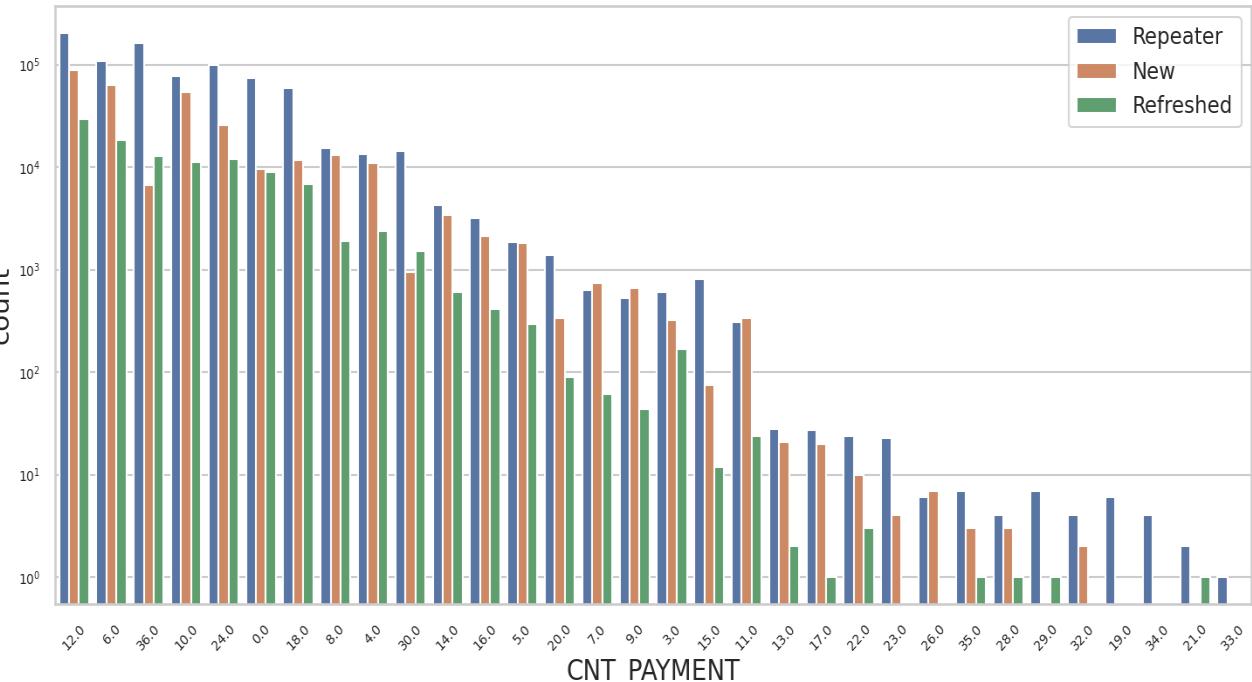
Distribution of NAME\_PORTFOLIO with purposes



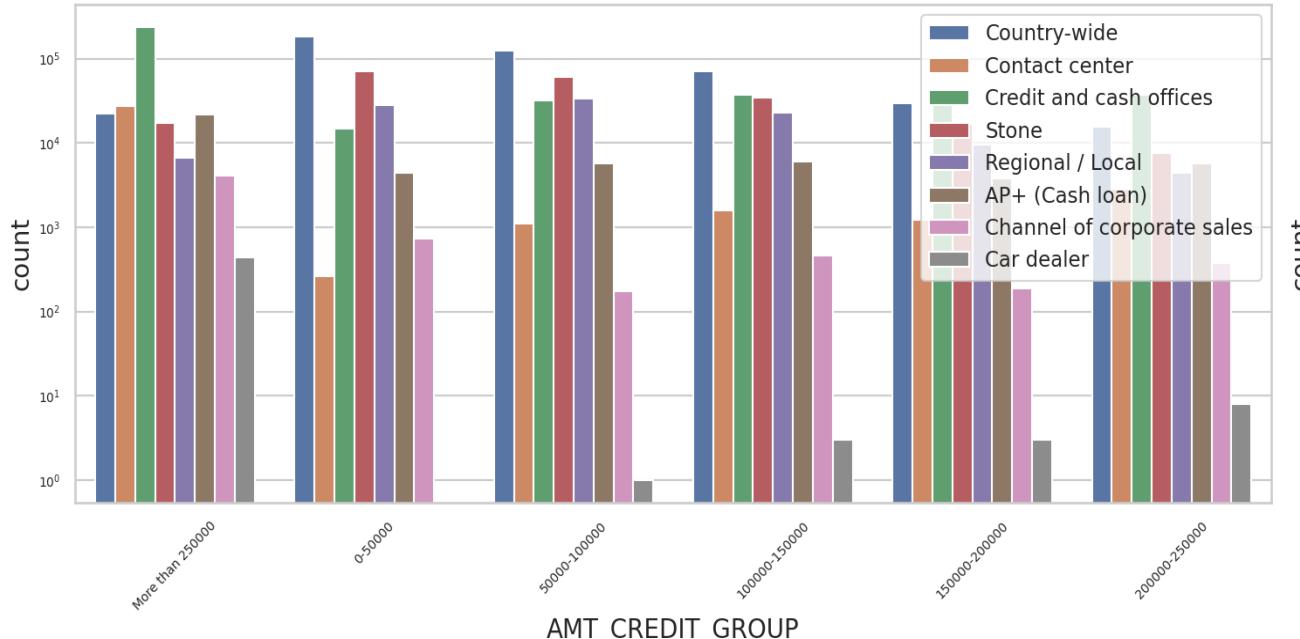
Distribution of NAME\_CONTRACT\_STATUS with purposes



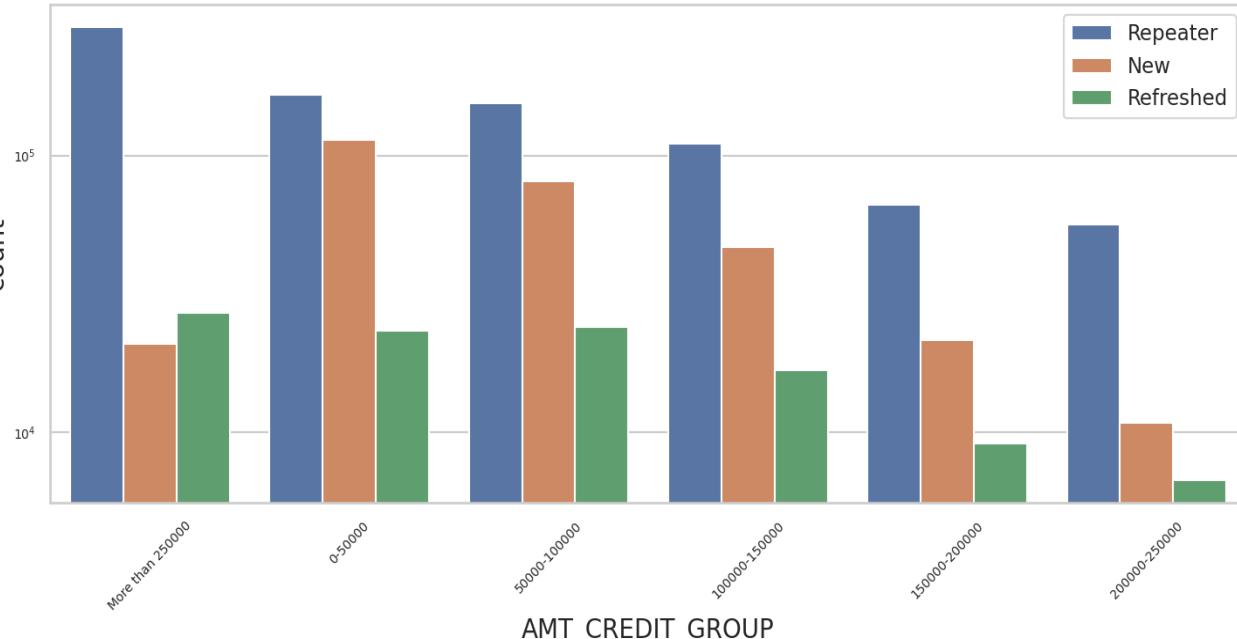
Distribution of NAME\_CLIENT\_TYPE with purposes



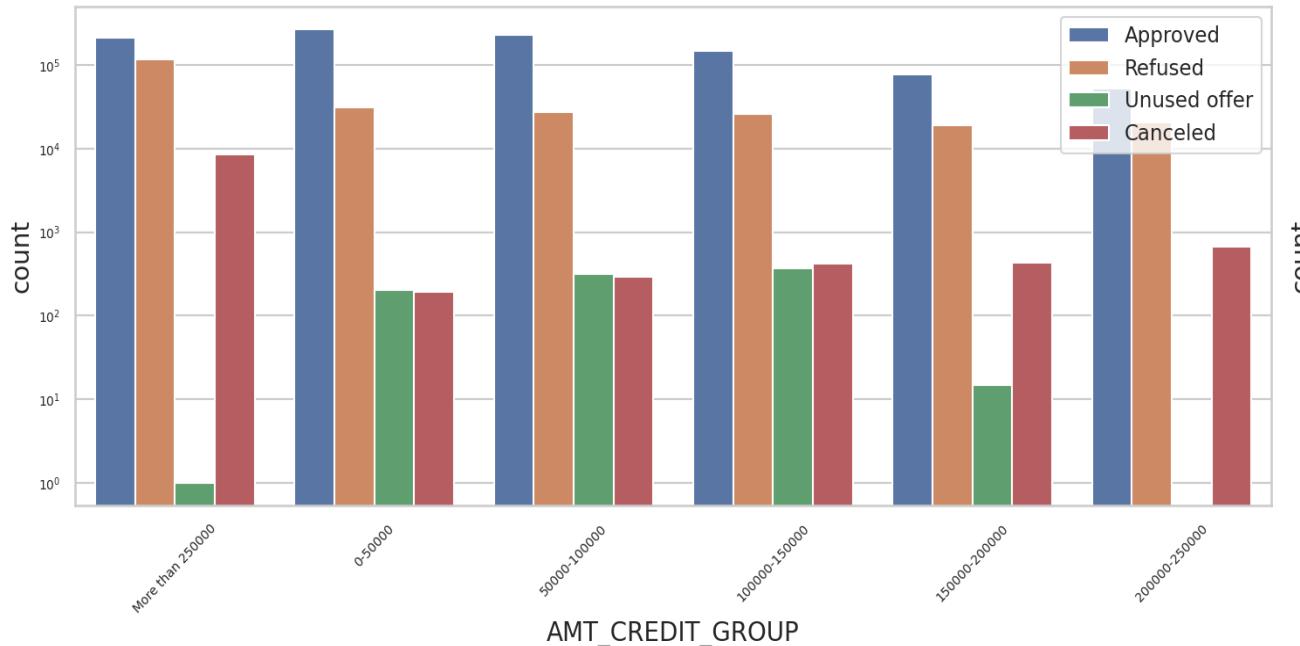
Distribution of CHANNEL\_TYPE with purposes



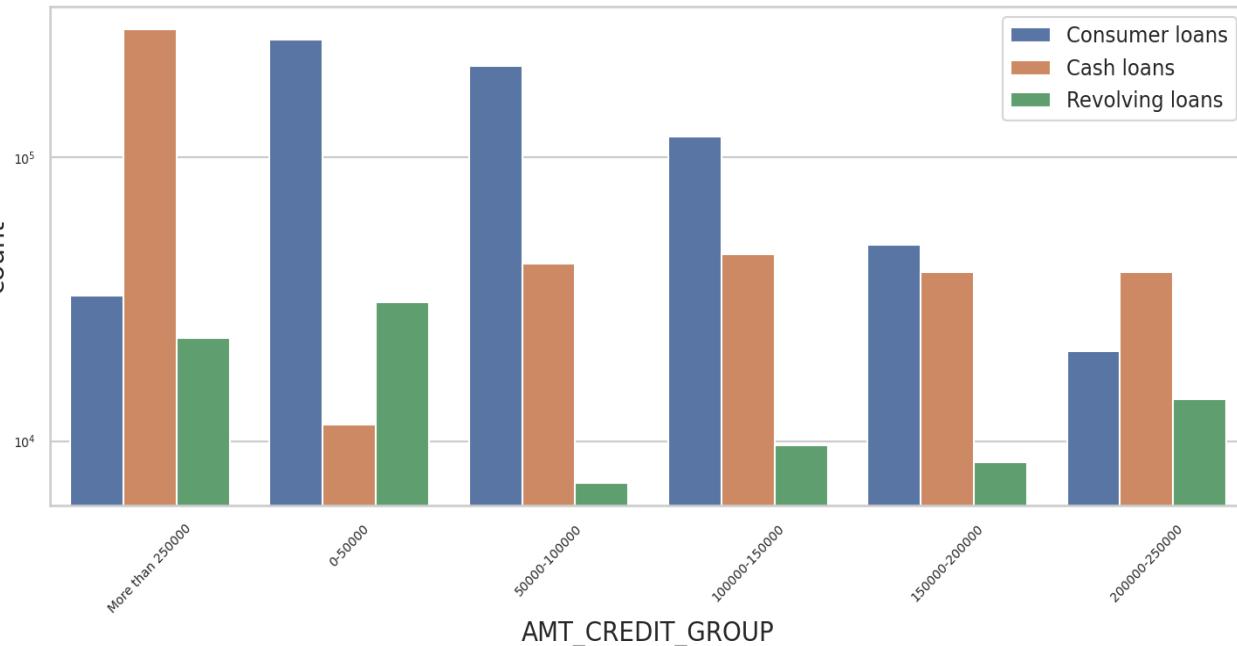
Distribution of NAME\_CLIENT\_TYPE with purposes

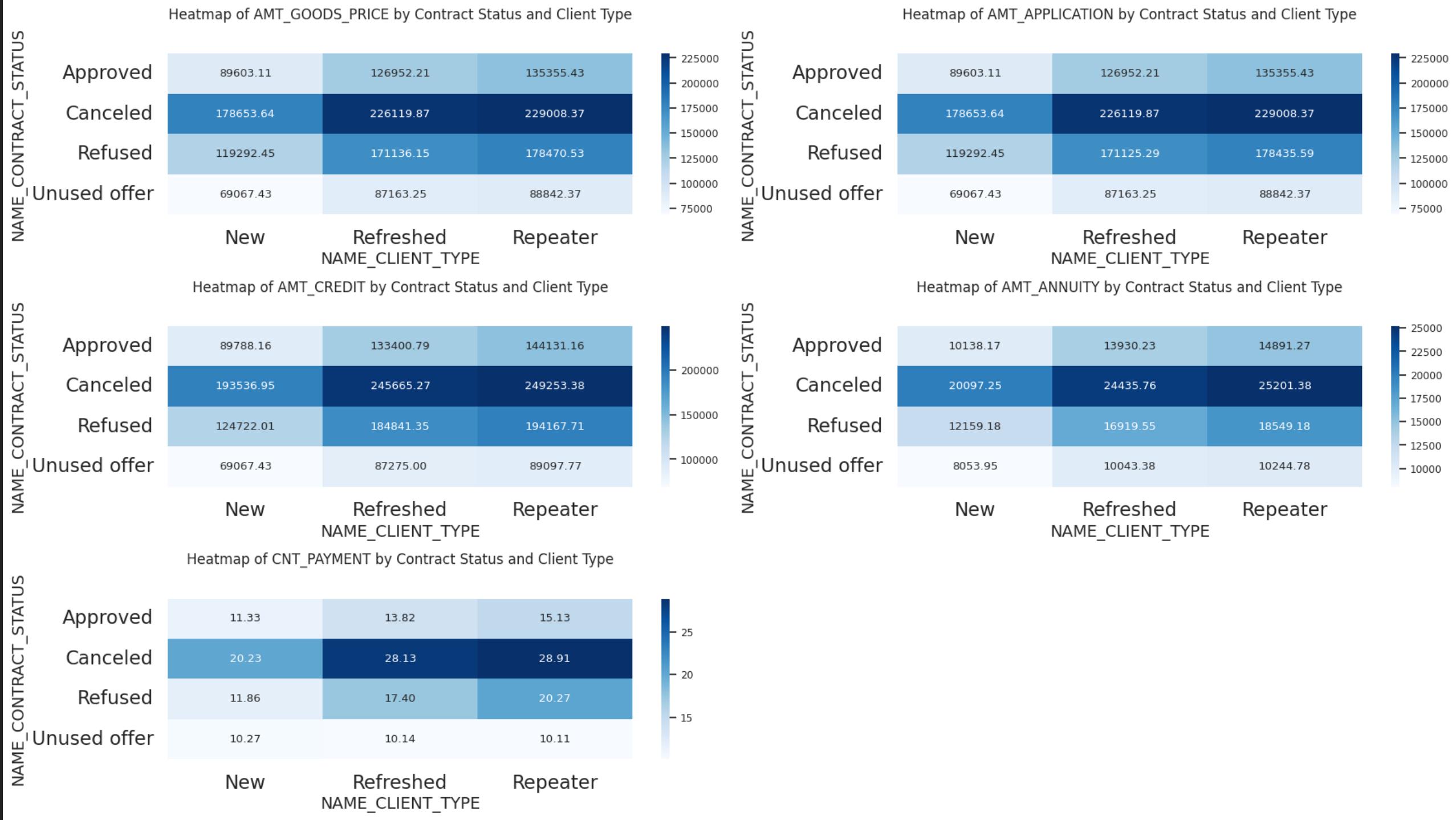


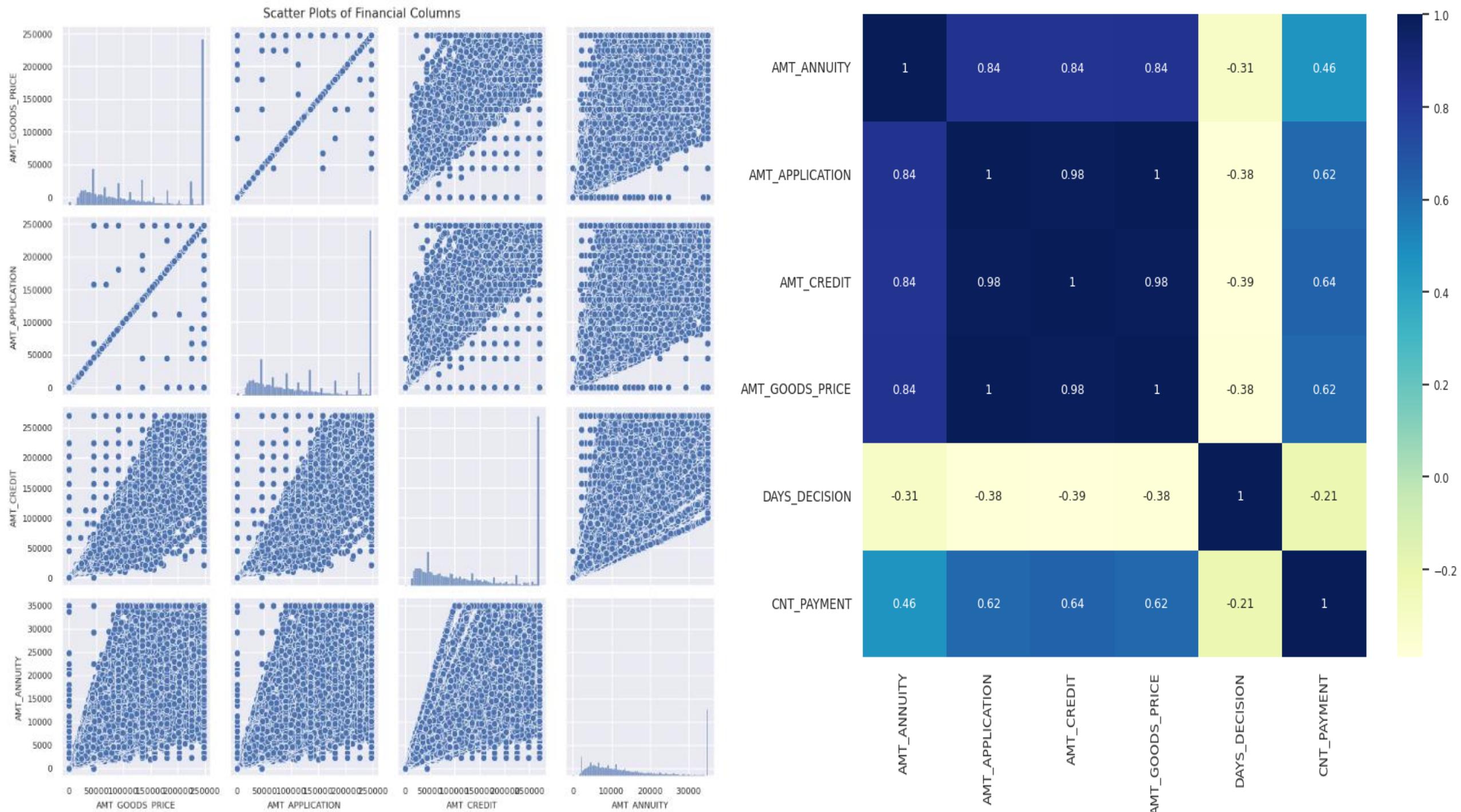
Distribution of NAME\_CONTRACT\_STATUS with purposes



Distribution of NAME\_CONTRACT\_TYPE with purposes







Clarifications on the types of graphs used for the previous application data:

### **Histogram Analysis:**

- Histograms were created for numerical columns in the dataset, revealing distinct density distributions.

### **Bar Plot Analysis (Bivariate):**

- In bar graphs comparing "AMT\_CREDIT" and "AMT\_GOODS\_PRICE\_GROUP" across various categorical columns, a consistent trend was observed. Moving right on the graphs indicated an increase in represented values.

### **Scatter Plot Analysis (Bivariate):**

- Scatter plots between numerical columns and the "NAME\_CONTRACT\_STATUS" field showed a linear relationship, suggesting a consistent correlation pattern. The majority of data points were associated with approved loans.

### **Count Graph Analysis (Bivariate):**

- Count graphs illustrated relationships between numerical values and different categorical values.

### **Heat Map Analysis:**

- Heat maps revealed notable points in various financial aspects. For instance, in aspects like 'AMT\_GOODS\_PRICE', 'AMT\_APPLICATION', 'AMT\_CREDIT', 'AMT\_ANNUITY', and 'CNT\_PAYMENT,' there were exceptionally low unused offers for new client types. Conversely, cancelled application amounts were high for Repeater client types, indicating more favorable policies or interest rates for repeat applicants.

### **Pair Plot Analysis:**

- The Pair Plot analysis highlighted strong correlations between 'AMT\_GOODS\_PRICE,' 'AMT\_ANNUITY,' 'AMT\_APPLICATION,' 'AMT\_CREDIT,' and 'CNT\_PAYMENT.' These correlations were expected because higher item values were associated with larger loan amounts and annual installments. Additionally, a notable correlation existed between 'AMT\_CREDIT' and 'AMT\_GOODS\_PRICE' due to the inherent link between the credit granted and the price of financed goods.

# Merged Both DataSets

## Default Rate Among Approved Loans:

Among approved loans, there is an 8.2% default rate.

## Concerning Trends in Previous Applications:

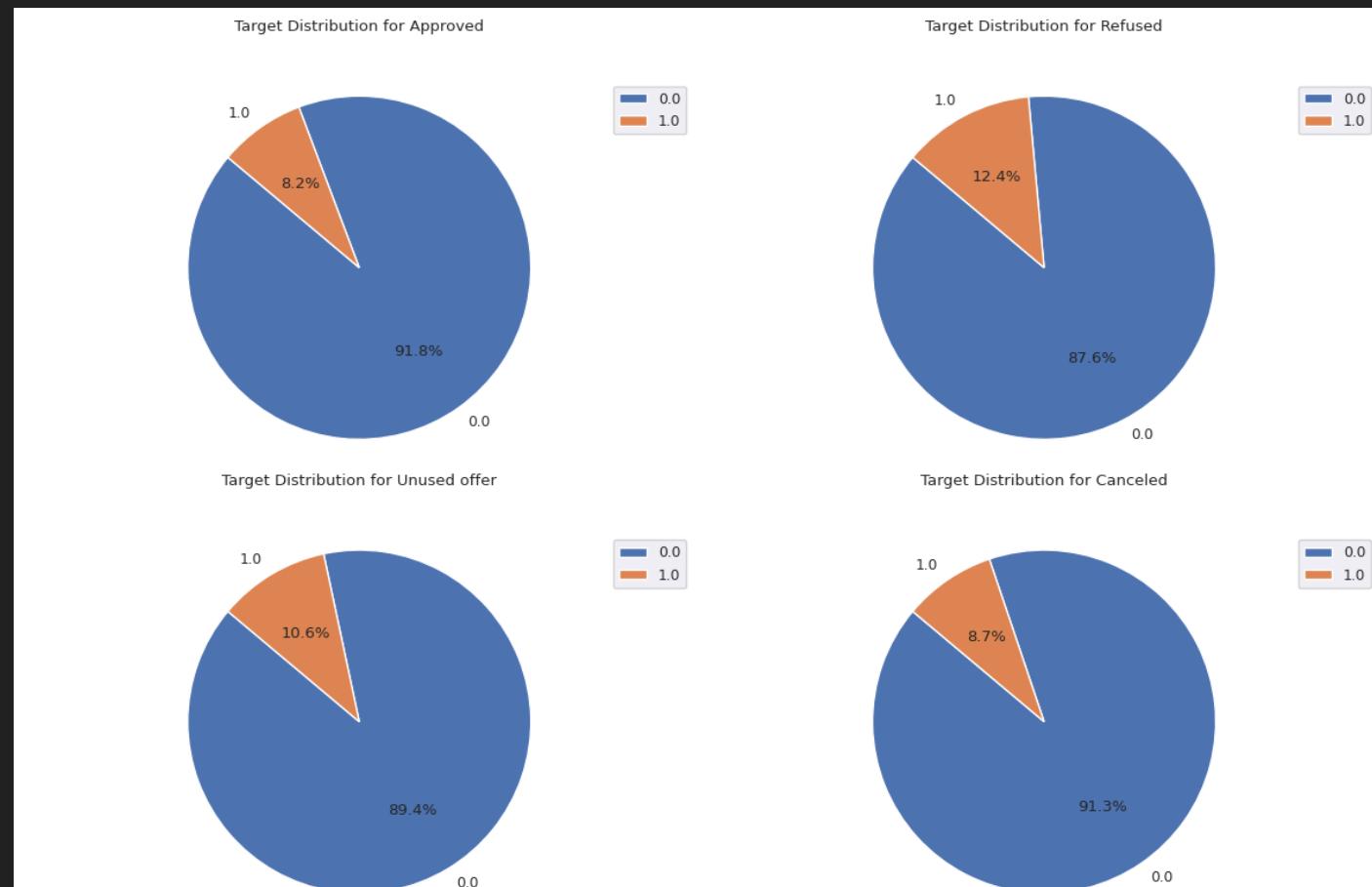
Alarming trends emerge when defaults are observed in previous applications that were either refused, cancelled, or unused.

## Financial Company Actions:

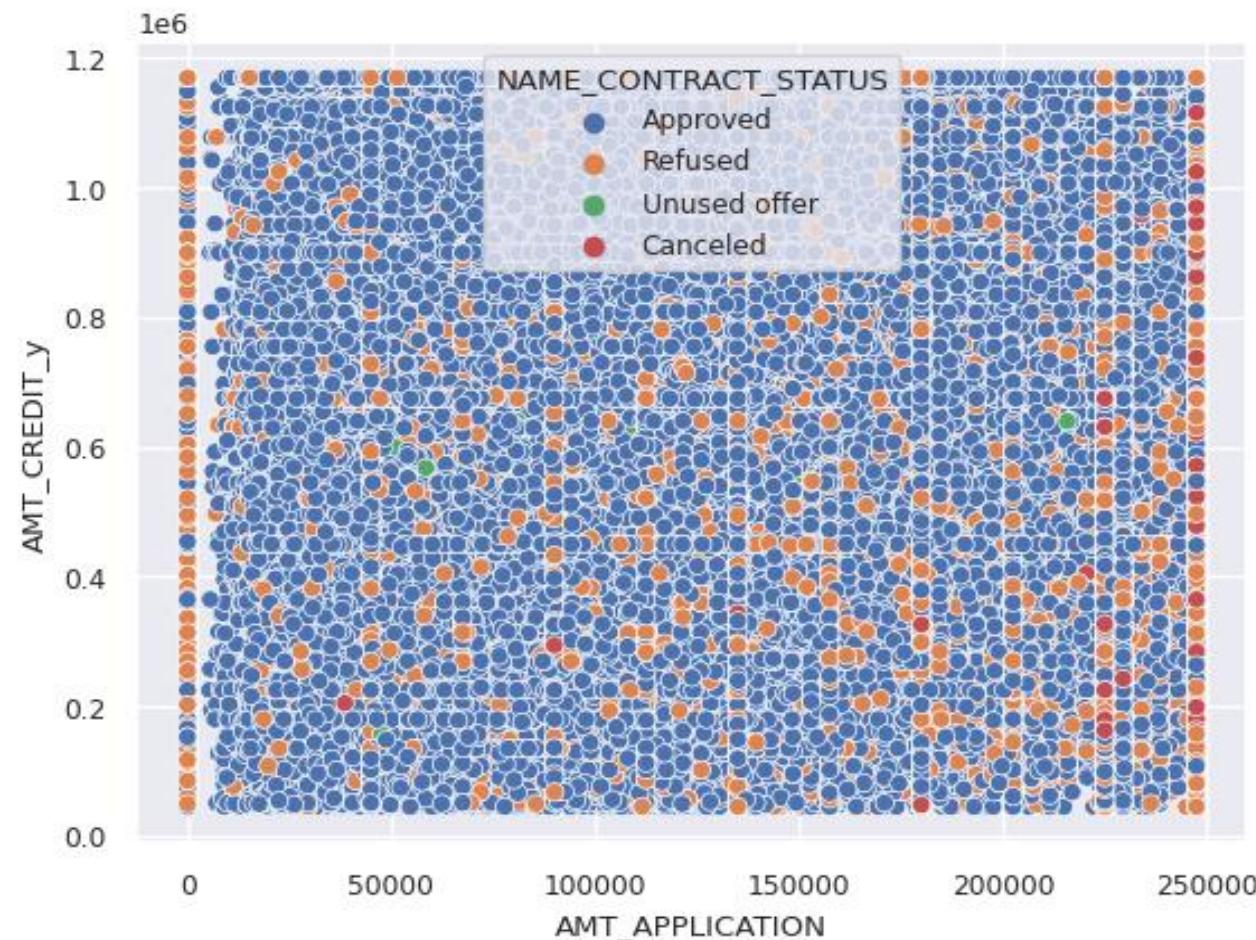
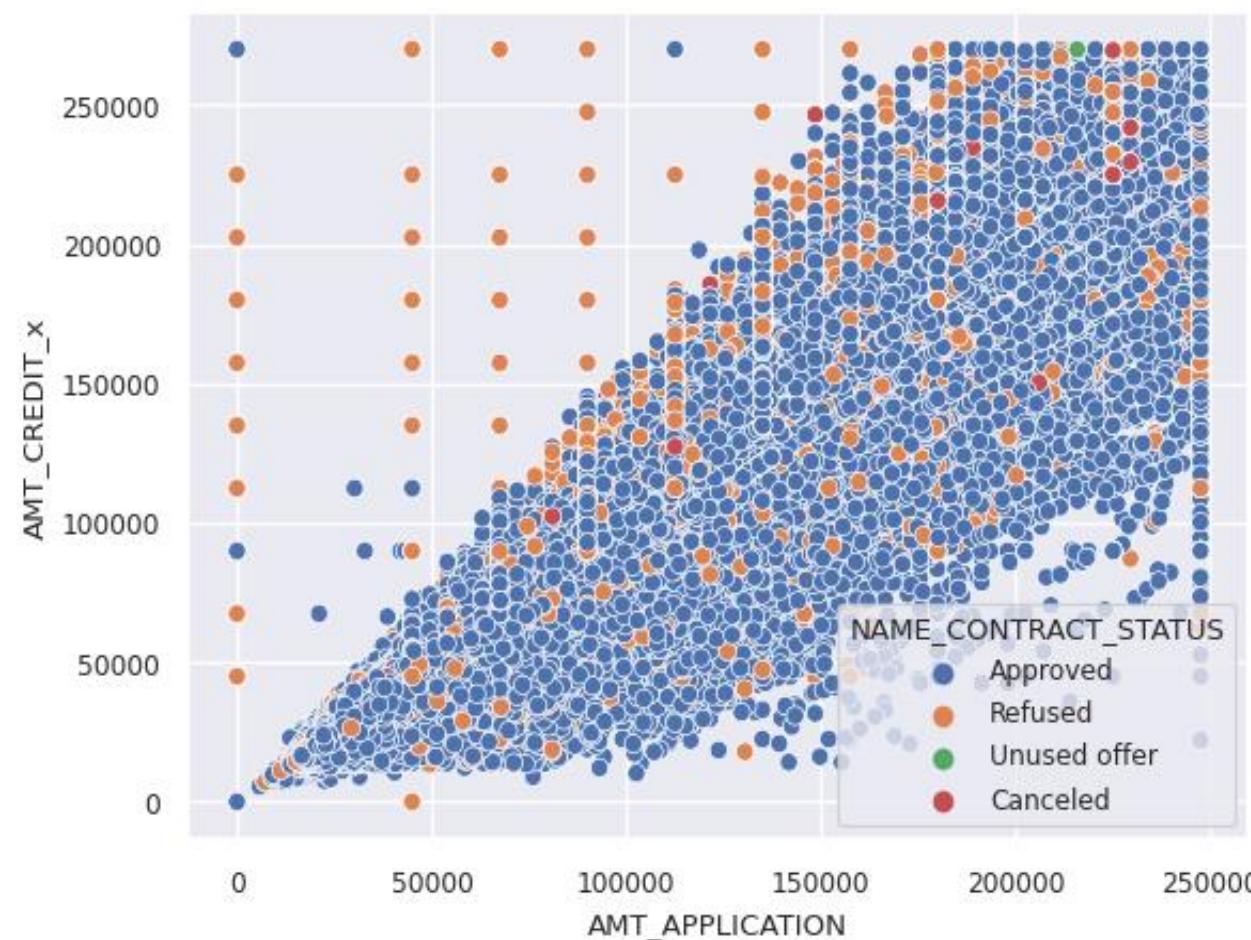
These observations suggest that despite the financial company's refusal or cancellation of previous applications, the current approved loans are still experiencing defaults.

## Persistent Default Risk:

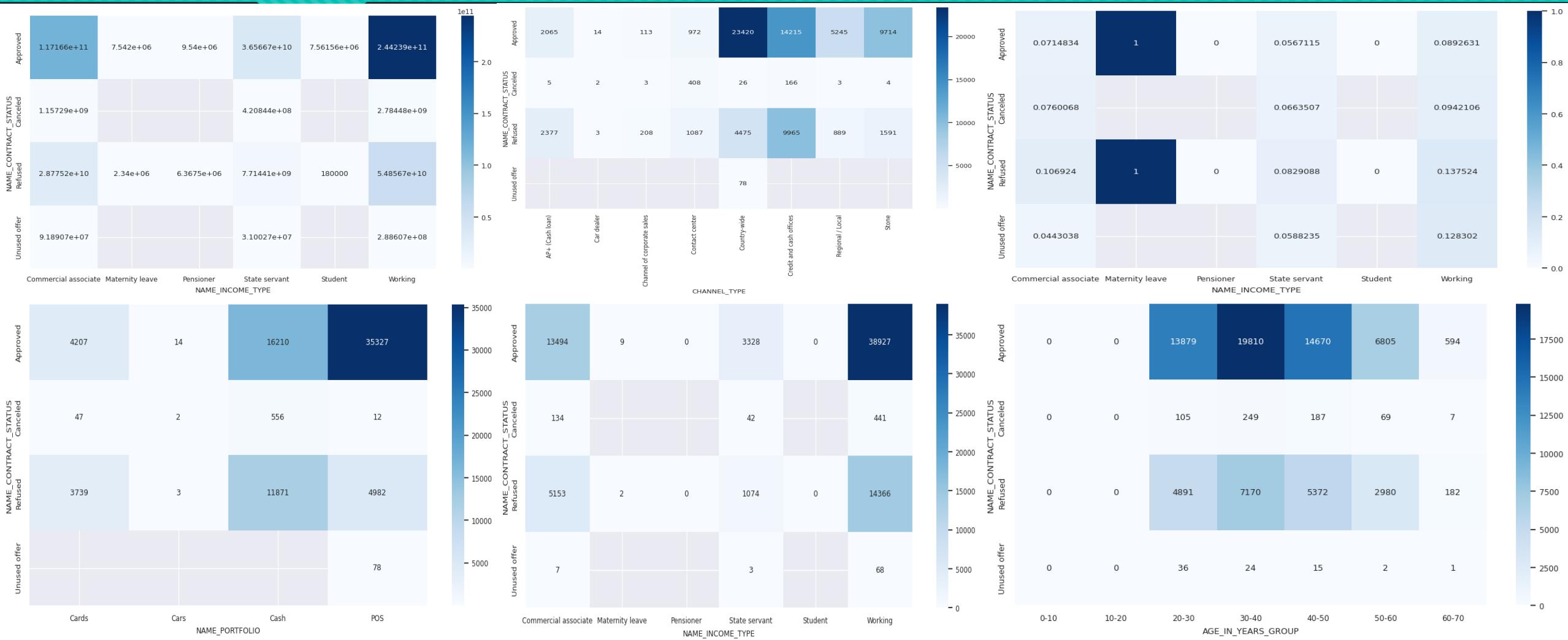
The presence of defaults in such cases indicates a persistent risk of default, warranting further investigation and risk assessment.



# Exploring Relationships with Scatterplots



# Key Insights from Analysis and Correlation Graphs



Here are several key points derived from the graph analysis above:

### **Heatmap Analysis for Default Correlation:**

- A higher value in the upper matrix of the heatmap indicates a stronger correlation with default (Target 1).

### **Income Types and Default Trends:**

- Working applicants with an approved status have the highest number of defaults among various income types.

### **Concerning Trends in Previous Applications:**

- Notably, previous applications with refused, cancelled, or unused loans also show defaults, indicating persistent default risks despite prior rejections or cancellations.
- Within the working class applicants, 14,366 individuals who were previously refused have defaulted on their current loans.

### **Age Groups and Default Rates:**

- The age group 30-40 exhibits the highest number of defaulters among approved loans, followed by age groups 20-30 and 40-50, which also display significant default rates.
- Previous applications with refused and cancelled loans are correlated with defaults in current applications.

### **Credit Offerings and Applicant Categories:**

- A significant observation is the higher credit offerings to individuals in the working and commercial associate categories.
- Refused applicants, especially students, tend to receive smaller credit amounts.

### **Default Analysis by Applicant Categories:**

- Among approved loans, the POS associate category has the highest count of defaulters.
- Approved loans in the car portfolio show only 14 instances of default.

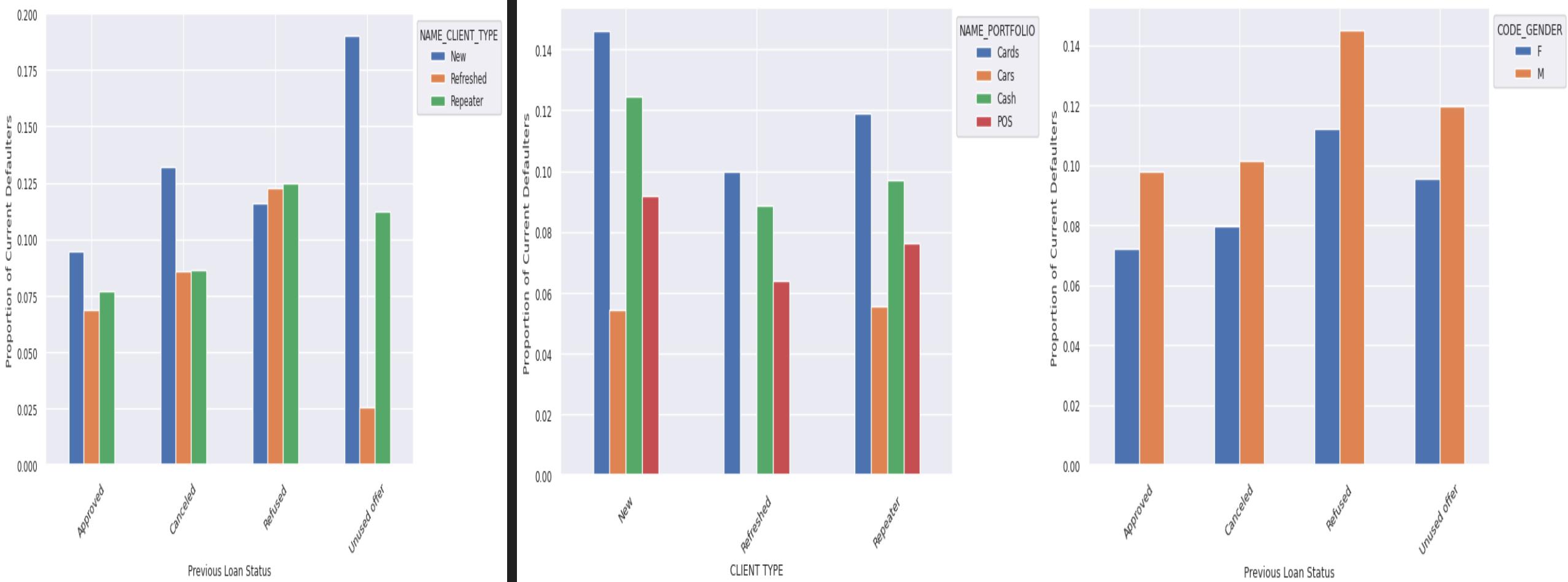
### **Default Cases by Country-wide and Car Dealer Categories:**

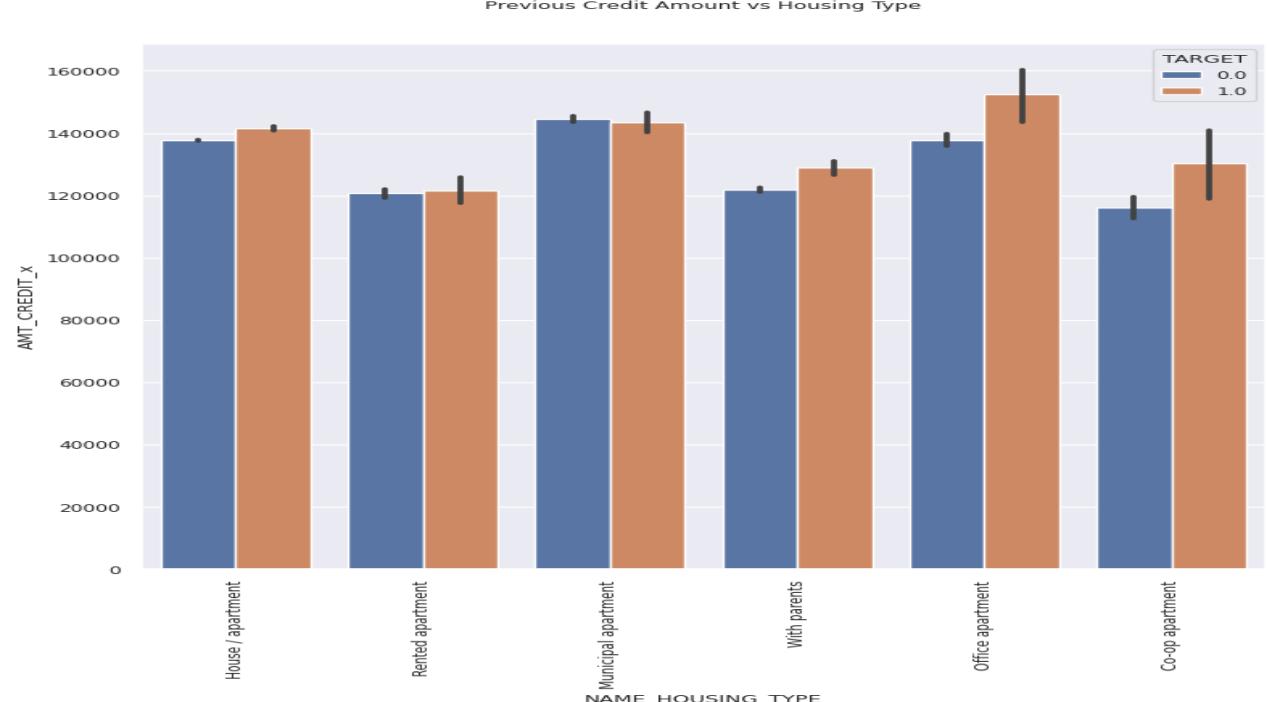
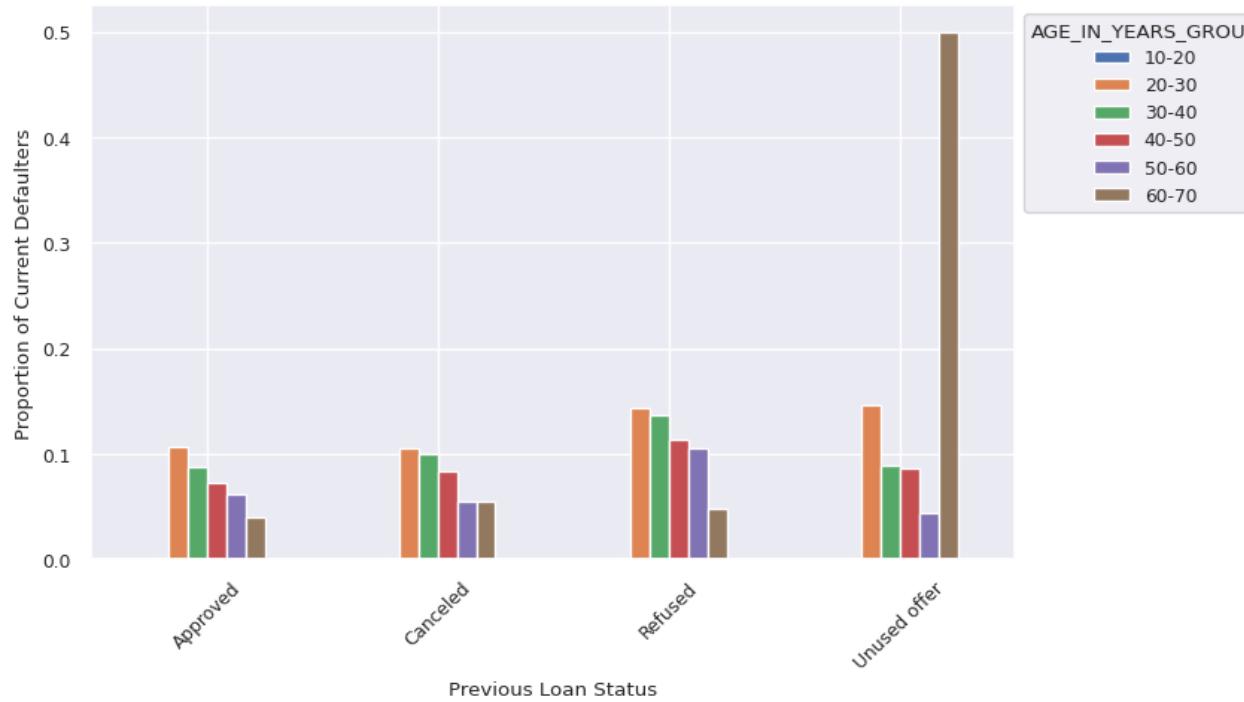
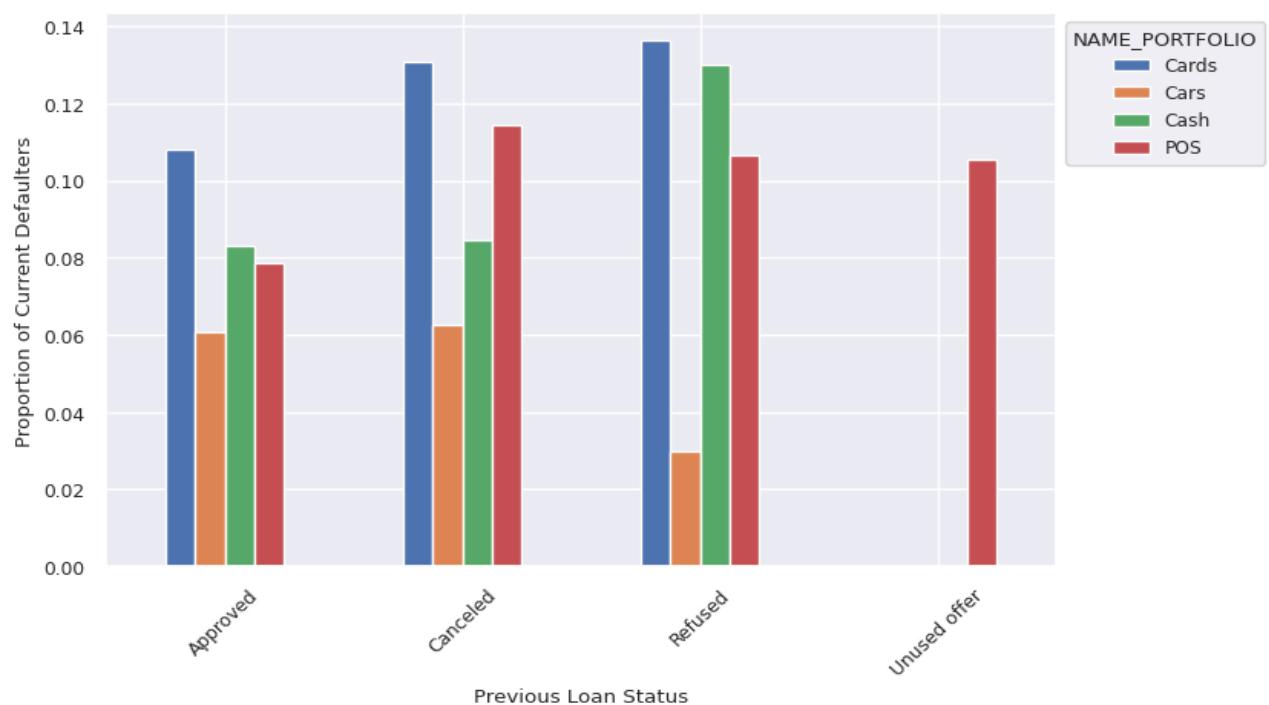
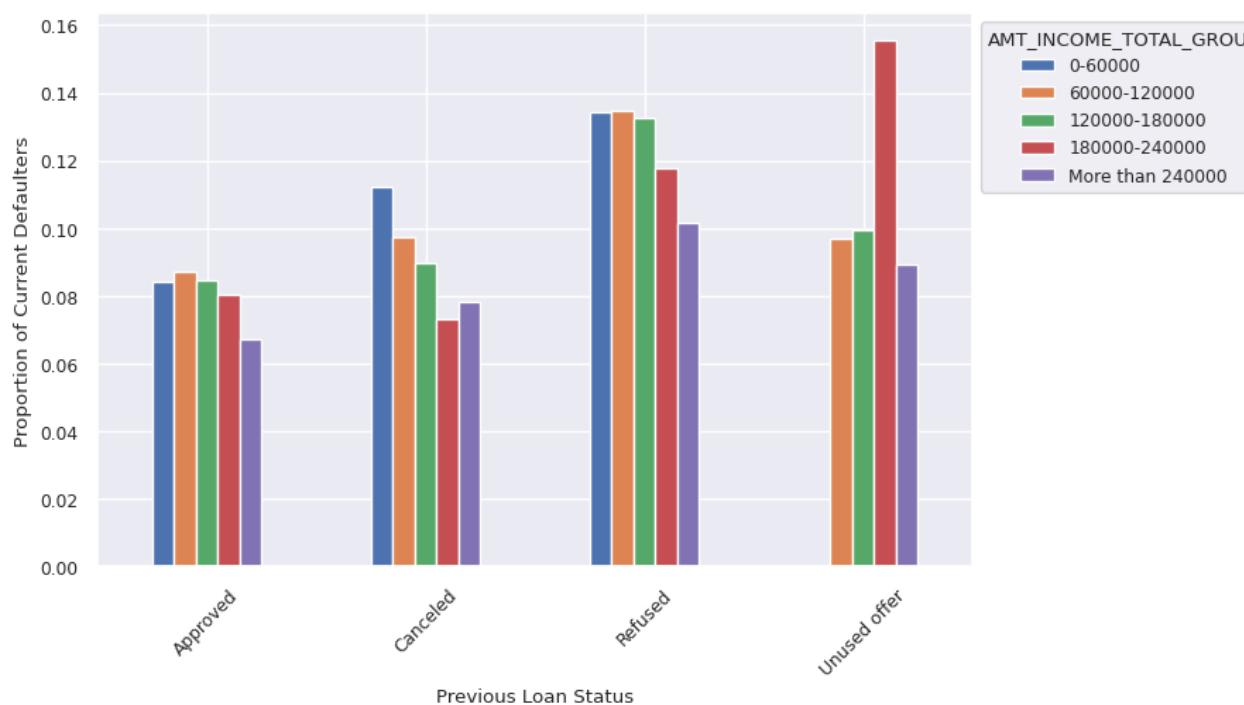
- Among approved loans, the country-wide associate category has the highest number of default cases.
- In the car dealer category of approved loans, there are only 14 instances of defaults.

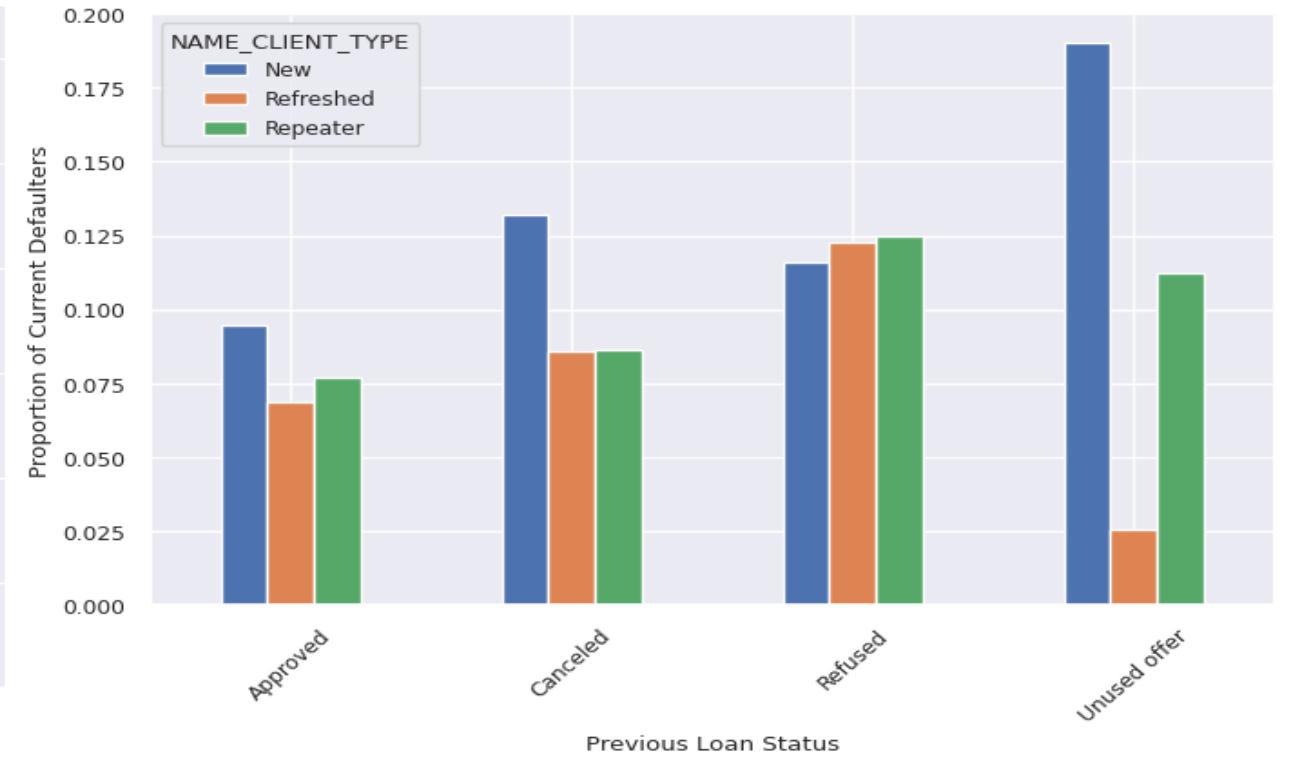
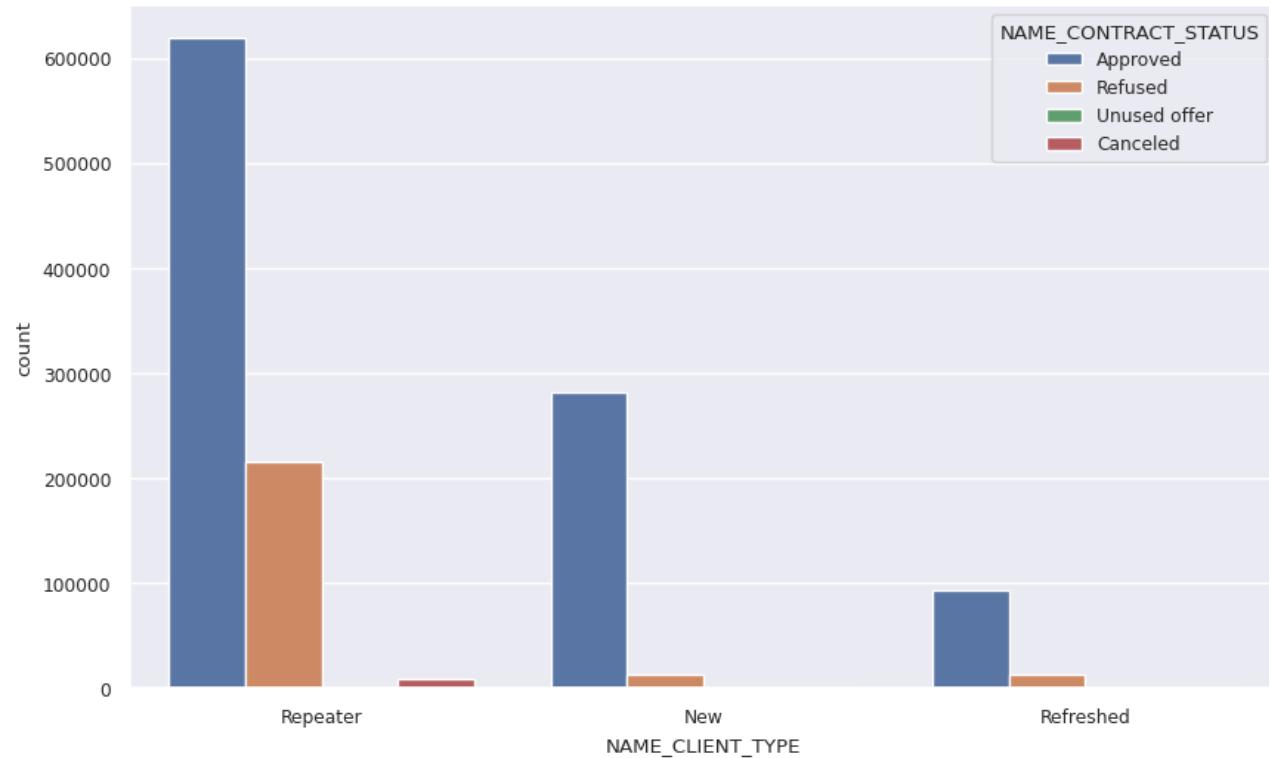
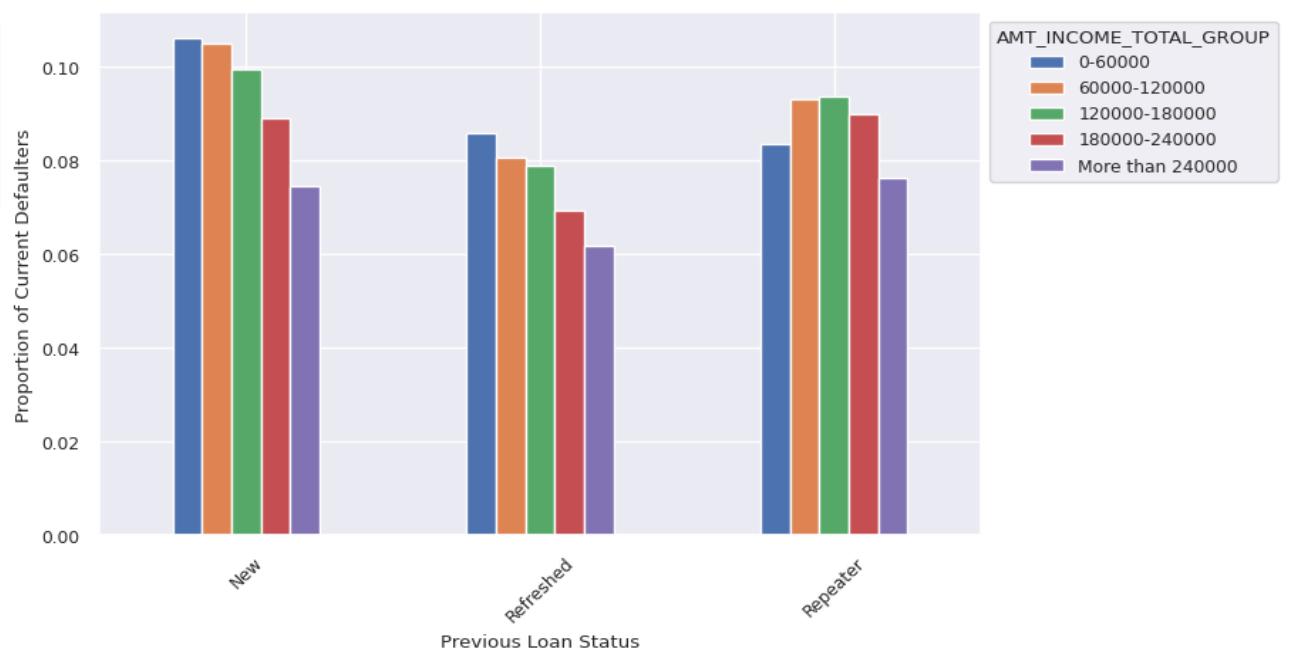
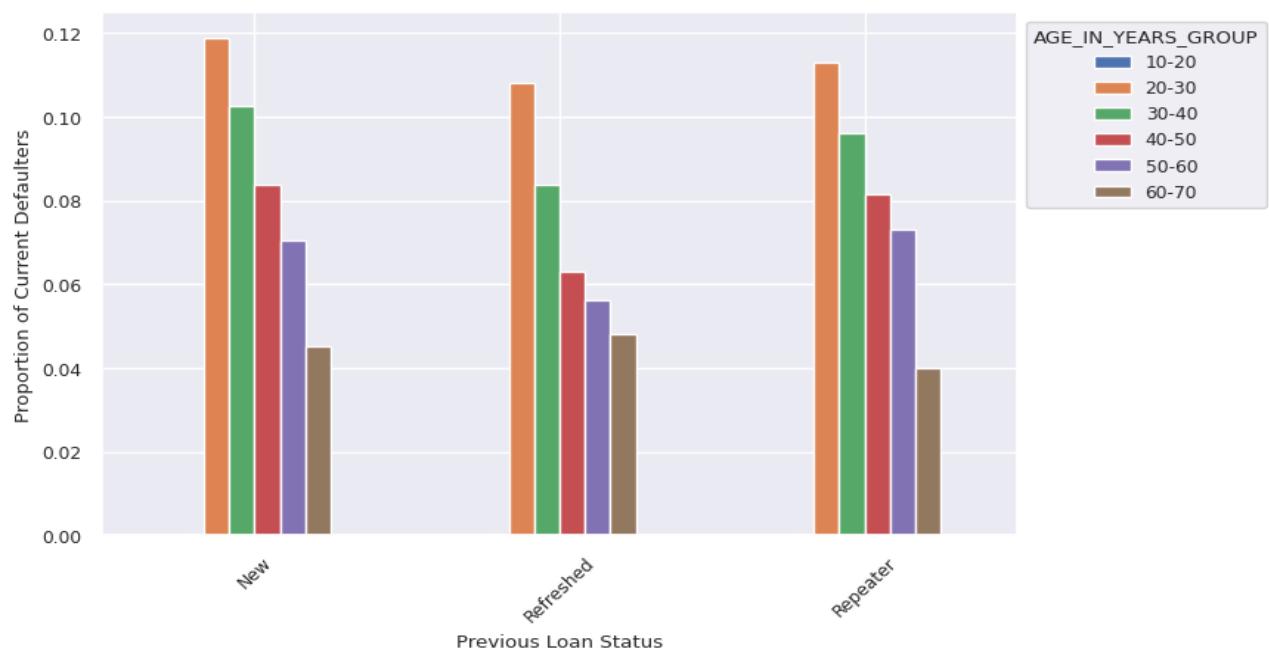
### **Maternity Leave and Other Categories:**

- When considering both approved and refused loans, the Maternity leave associate category stands out with the highest count of default cases.
- No instances of defaults are observed in the pensioner and student categories for both approved and refused loans.

# Exploring Column Relationships through Bivariate Bar Graphs







Here are several key points derived from the graph analysis above:

### **Default Patterns by Previous Loan Status and Gender:**

- Refused clients have a higher default rate compared to approved clients.
- Across all loan statuses, males have a higher default rate than females.

### **Default Patterns by Previous Loan Status and Client Type:**

- Defaulters are more prevalent among clients with a previous loan status of "Unused offers," especially among "New" clients.
- For previous "Approved" statuses, new clients show a higher default rate, followed by repeater clients.
- Among applicants with a previous status of "Refused," defaulters are more common among repeater clients.
- Among applicants with a previous status of "Canceled," new clients have a higher default rate.

### **Default Patterns by Age Group for Approved Clients:**

- For approved clients, younger applicants tend to have a higher default rate.
- Senior applicants among approved clients display a lower default rate.

### **Default Patterns by Income Range for Unused Offers:**

- Among unused offers, applicants with an income range of 180,000-240,000 have a higher default rate, while incomes exceeding 240,000 show the least defaults.
- Across all loan statuses, default rates are relatively consistent across various income groups.

### **Default Patterns by Loan Status and Portfolio Type:**

- Clients who refused applied loans for cards exhibit a higher default rate.
- For approved loan statuses, clients applying for cars have a lower default rate.
- Clients who refused applied loans for cars also show a lower default rate.

### **Housing Type and Default Rates:**

- Office apartments have a higher credit default rate (Target 1), while municipal apartments show a higher credit repayment rate (Target 0). This suggests caution when granting loans for office apartments.

### **Client Types and Approved Loans:**

- Repeater clients have a higher count of approved loans compared to New and Refreshed clients.

### **Default Patterns by Previous Loan Status and Client Type (Continued):**

- Among clients previously categorized as "Unused offers," there is a higher occurrence of defaulters, particularly among "New" clients.
- For clients with a previous loan status of "Approved," "New" clients have a higher count of defaults, followed by "Repeater" clients.
- Among applicants with a previous loan status of "Refused," instances of defaulters are higher among "Refreshed" clients.
- Among those with a previous loan status of "Canceled," the count of defaulters is higher among "New" clients.

### **Age Groups and Default Rates (Continued):**

- Across all previous loan statuses, younger applicants exhibit a higher frequency of defaults.
- Senior applicants show a lower incidence of defaults compared to other age groups, regardless of the previous loan status.

### **Income Ranges and Default Rates:**

- New applicants with an income range of 0-60,000 tend to have a higher rate of defaults.
- Among repeater applicants, those with an income range of 0-60,000 show a lower default rate compared to other client types.
- Repeater applicants with an income exceeding 240,000 exhibit a higher default rate compared to other income ranges.

### **Portfolio Types and Default Rates:**

- Among the cards portfolio, new applicants display the highest default rate.
- Within the cars portfolio, refreshed applicants have no instances of default.
- Among different client types in the cards portfolio, refreshed applicants exhibit the lowest default rate compared to others.

# Key Insights for Risk Assessment and Lending Strategies

## **Leading Factors in Defaults:**

- Medium income levels, age groups of 20-30 and 30-40, being male, and occupation types like business type 3 and self-employment are associated with higher default rates.

## **Family and Education Factors:**

- Married clients with five or more children are at an elevated risk of defaulting.
- Default rates are highest among males with lower secondary education.
- Maternity leave applicants tend to have a higher likelihood of default when they have a larger number of children.

## **Unused Applications and Age:**

- Unused applications are more prevalent in the higher age group, potentially explaining their underutilization.

## **Gender and Default Rates:**

- Female applicants have lower default rates, suggesting they could be given extra consideration.

## **Working Applicants and Caution:**

- While more defaulters are working applicants, this does not imply outright refusal. Further examination of other parameters is necessary.

### **Educational Background Consideration:**

- Banks should prioritize applicants with a background of 'Higher education' while exercising caution with those who have completed 'Secondary/secondary special,' as they tend to face difficulties in making payments.

### **Income Type and Risk:**

- It is advisable for banks to be cautious with clients with an income type of 'Working' due to their higher likelihood of payment difficulties. Instead, they should focus on clients classified as 'Commercial associate,' 'Pensioner,' and 'State servant.'

### **Housing Type and Payment Issues:**

- To reduce payment difficulties, banks should pay more attention to clients residing in 'House/apartment' housing types, as they have fewer payment issues.

### **Channel Types and Loan Outcomes:**

- Banks should note that the 'Country-wide' channel type results in more approved loans, whereas the 'Credit and cash offices' channel type is associated with a higher number of canceled and refused loans.