

# CREDIT EDA ASSIGNMENT

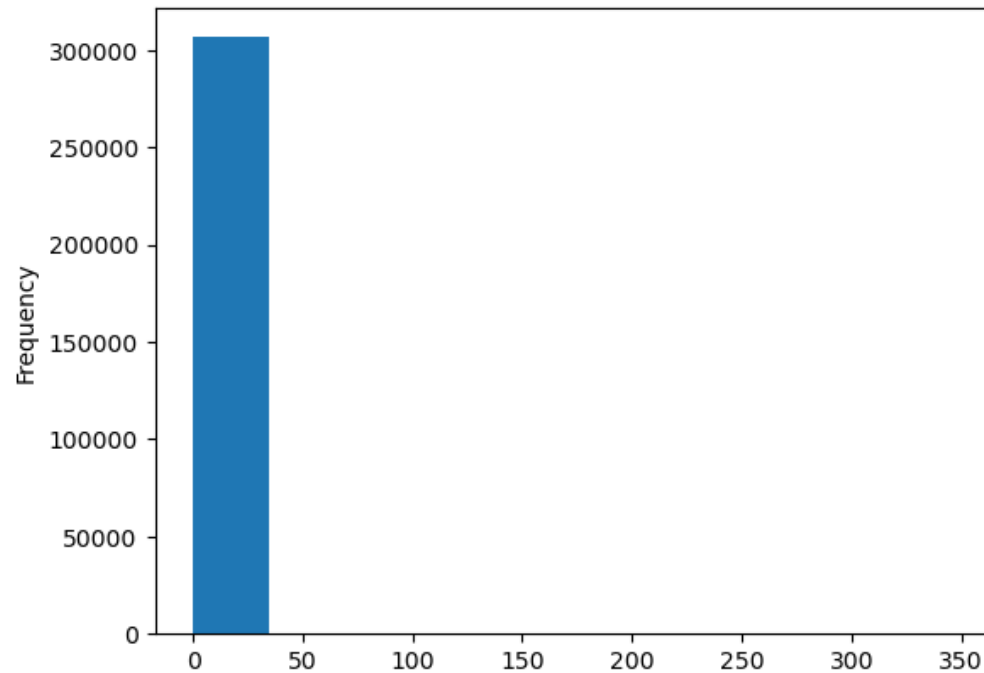
BY

HARSHAVARDHINI J S

# Application\_data.csv

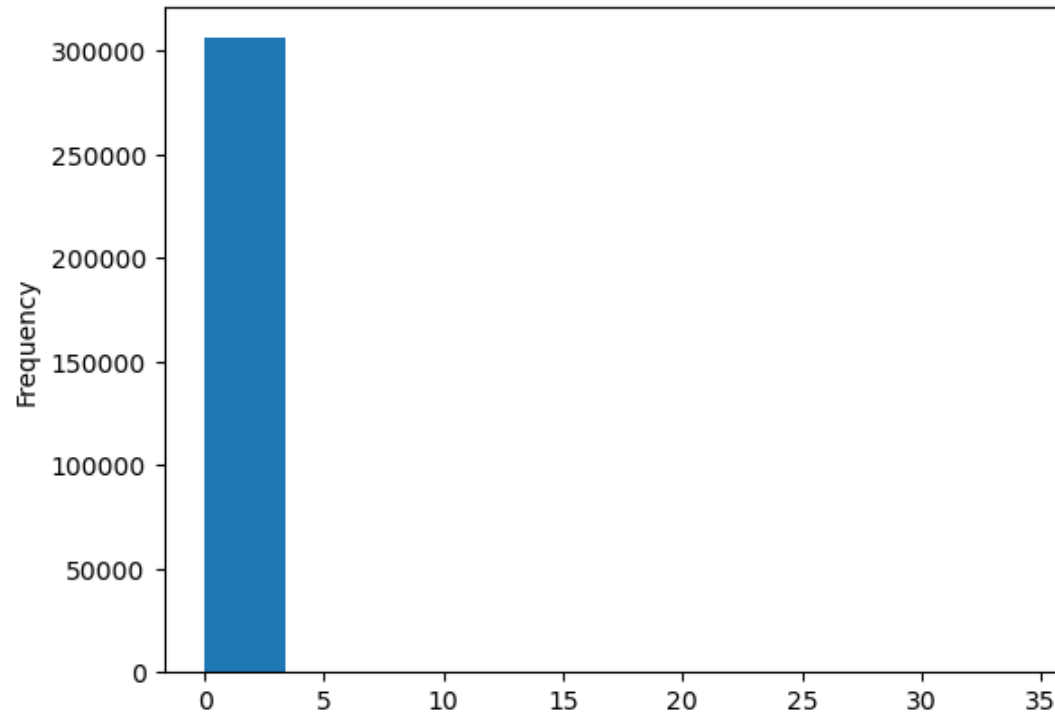
## Handling missing values

Column: OBS\_30\_CNT\_SOCIAL\_CIRCLE



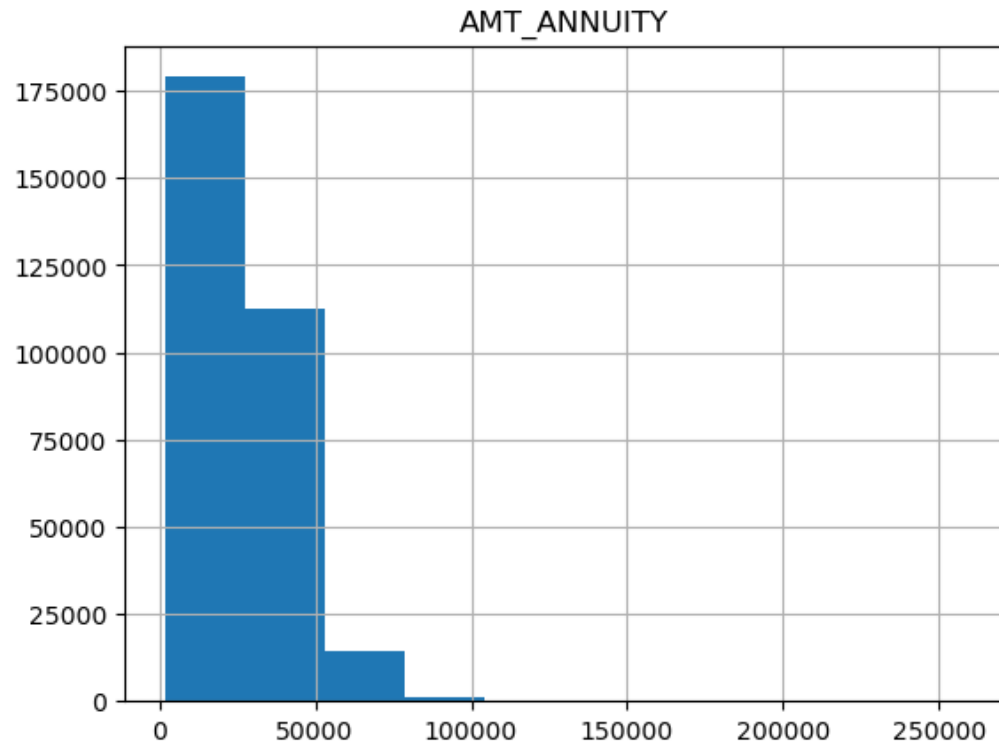
The histogram plot shows that the column has a skewed distribution and hence the suggested imputation method for the null values can be replaced by mode=0 which is the maximum frequency value.

Column:DEF\_30\_CNT\_SOCIAL\_CIRCLE



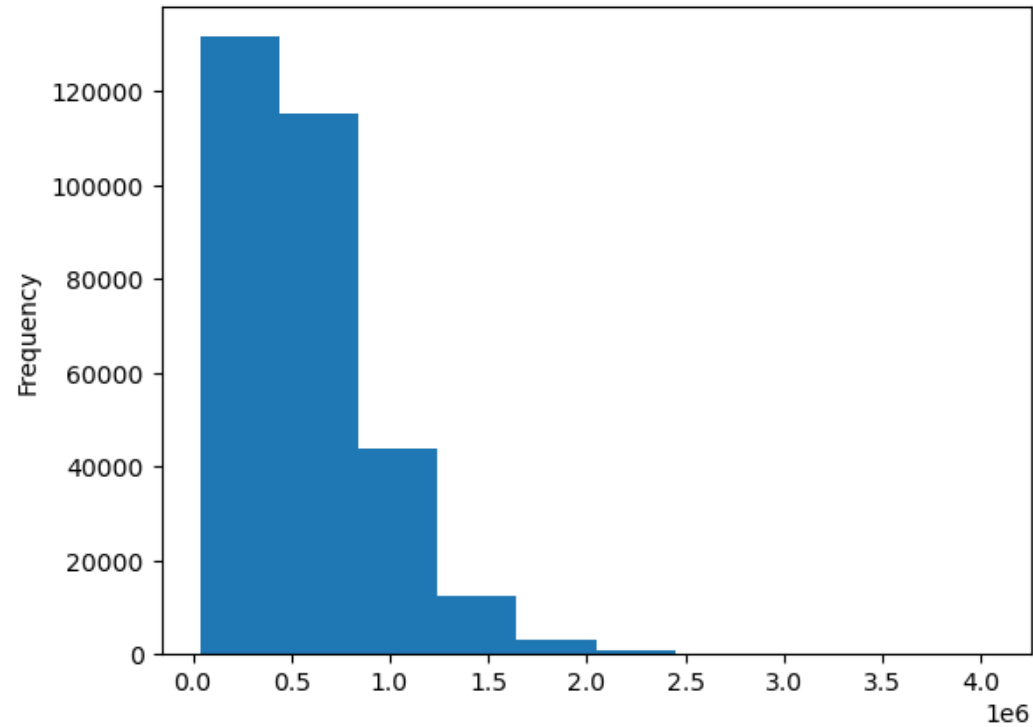
The suggested imputation method for the null values in the column DEF\_30\_CNT\_SOCIAL\_CIRCLE can be replaced by mode=0 which is the maximum frequency value.

Column: AMT\_ANNUIITY



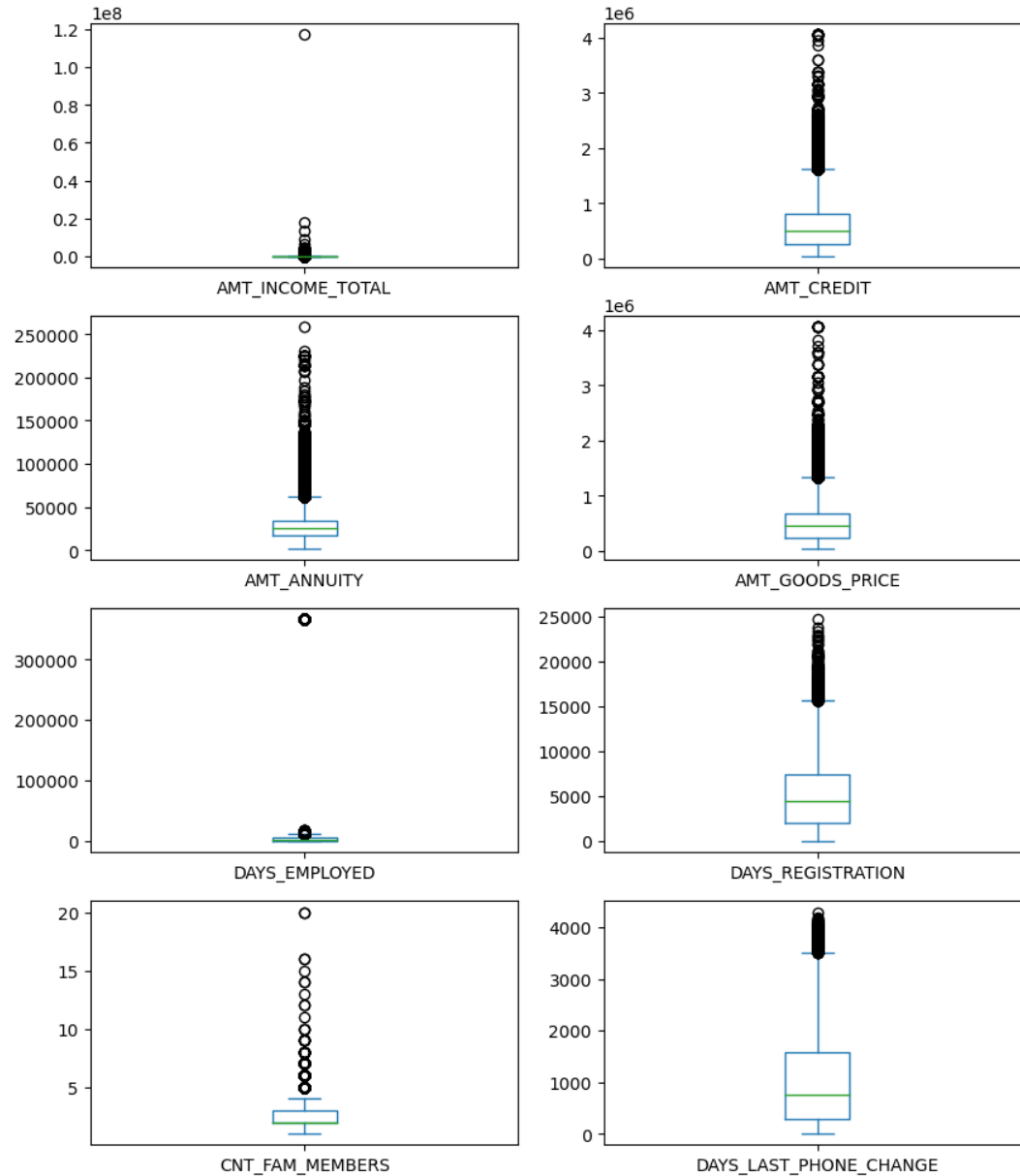
The column has an asymmetrical distribution and hence the mean values is skewed. Since the null values are very less in percentage the columns can be either dropped off or can be imputed with median value = 24903.

Column: AMT\_GOODS\_PRICE



The column AMT\_GOODS\_PRICE has an asymmetrical distribution and hence the mean values is skewed. Since the null values are very less in percentage the columns can be either dropped off or can be imputed with median value = 450000

# Handling outliers

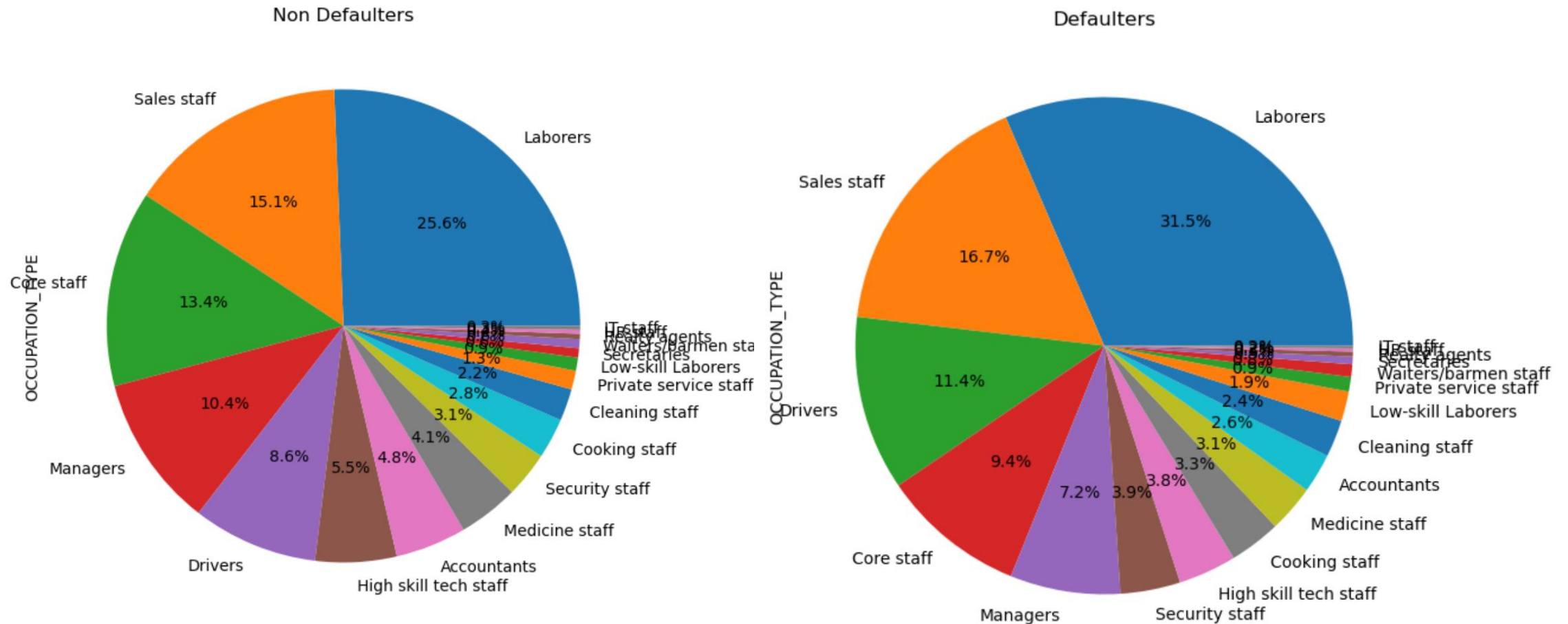


## SUGGESTIONS

- All the variables in num\_cols exhibit very low percentage of outliers which is less than 5% hence these outliers can be dropped off so that there is not much effect on the prediction of the results.
- The variable DAYS\_EMPLOYED shows 18.33% of outliers which can be capped with the upper and lower bound of the interquartile range for further analysis.

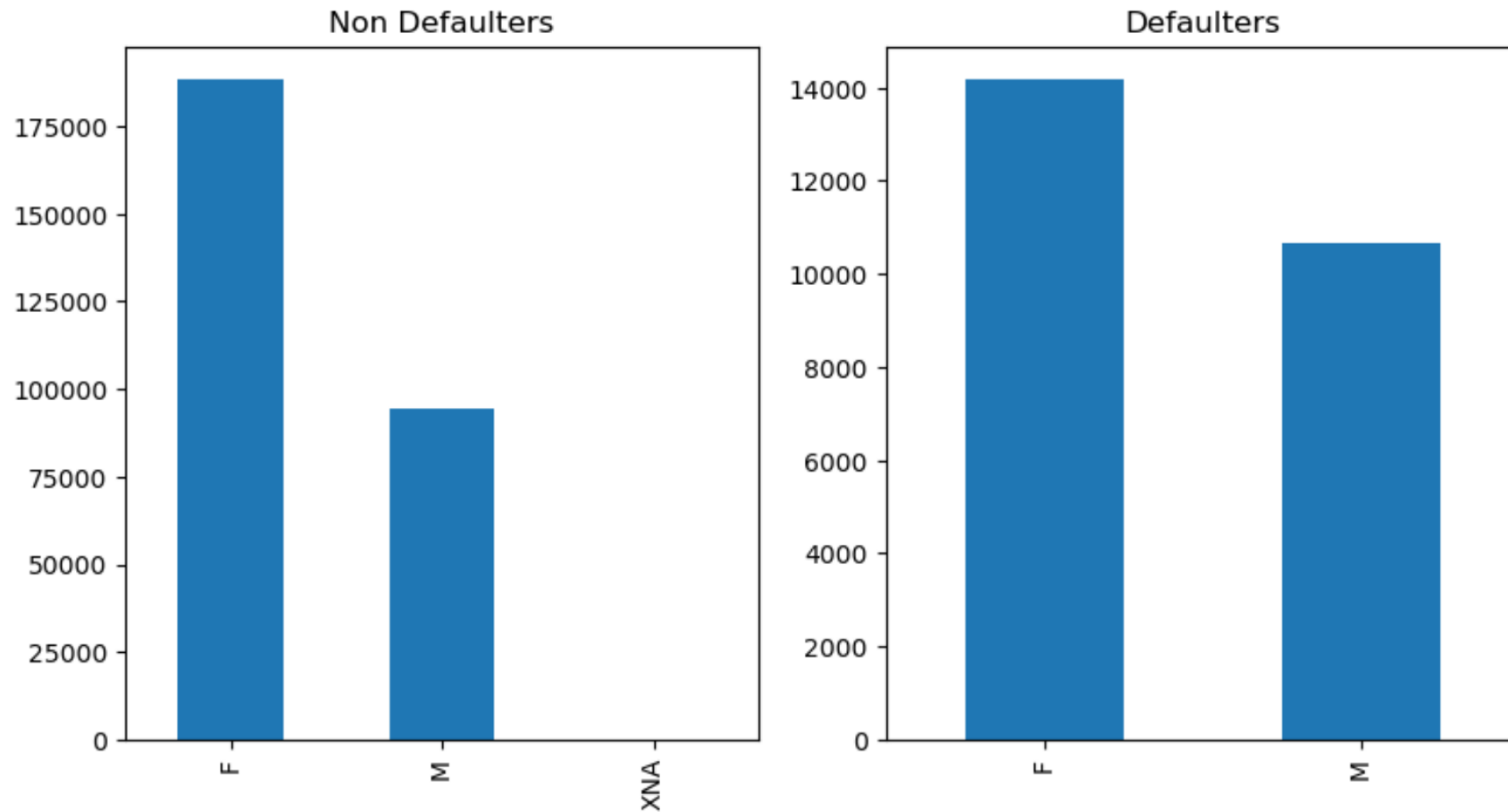
# Univariate Analysis

## OCCUPATION\_TYPE



- Laborers, Sales staff and Drivers contribute to more than 50% of defaulters list.
- Laborers, Sales staff, Drivers, Core staff and Managers contribute to more than 75% of defaulters list.
- 31.5% of defaulters are with occupation type as Laborers
- 25.6% of Non Defaulters are with occupation type as Laborers

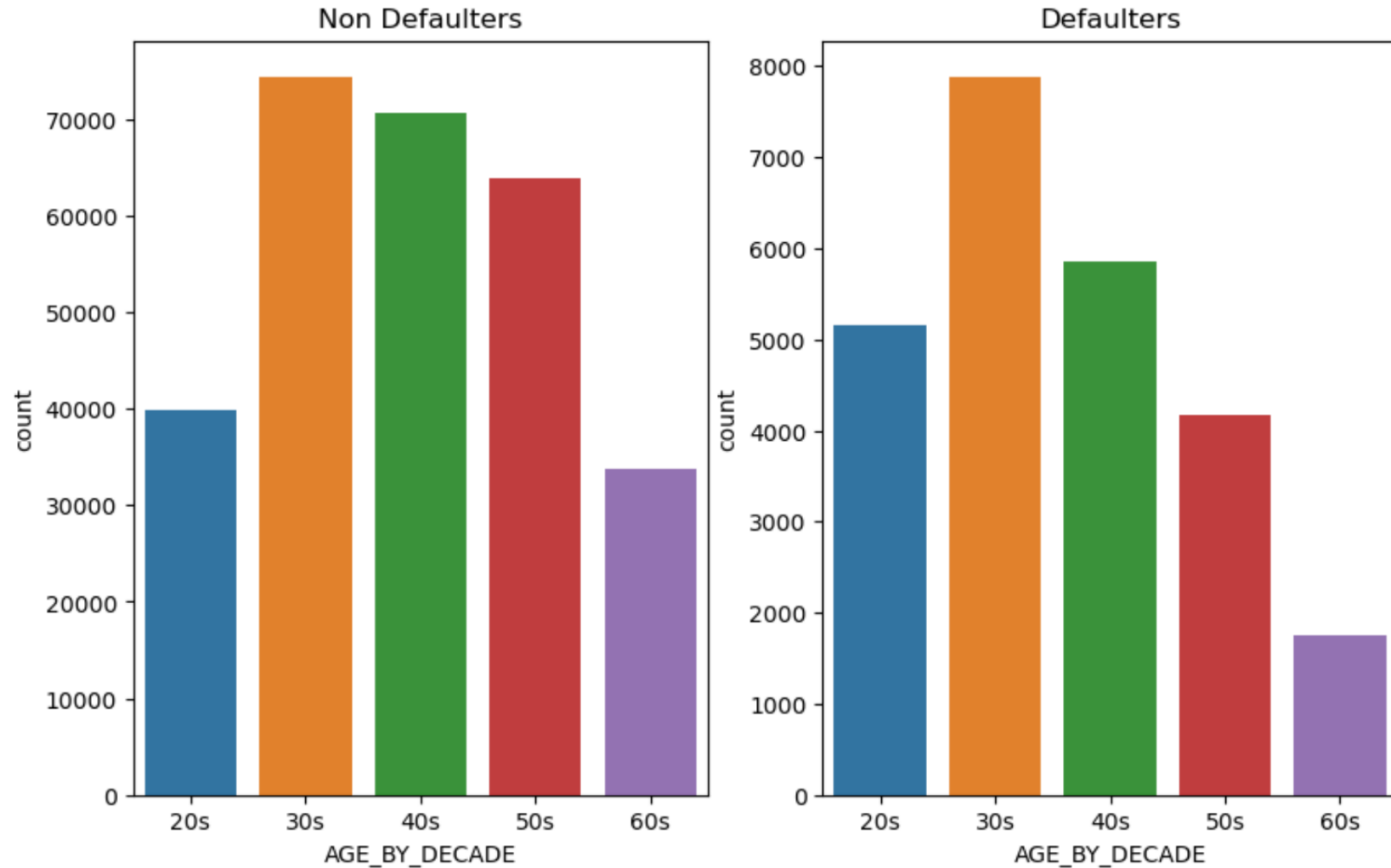
## CODE\_GENDER



- The bar plot depicts that the female gender were higher as compared to the male in applying credit loans
- The female gender percentage is 66.6% while the male percentage is 33.3% in case of Non defaulters
- The female gender percentage is 57.0% while the male percentage is 42.9% in case of defaulters The female receiving loan tend to repay in higher percentage than male receiving loans

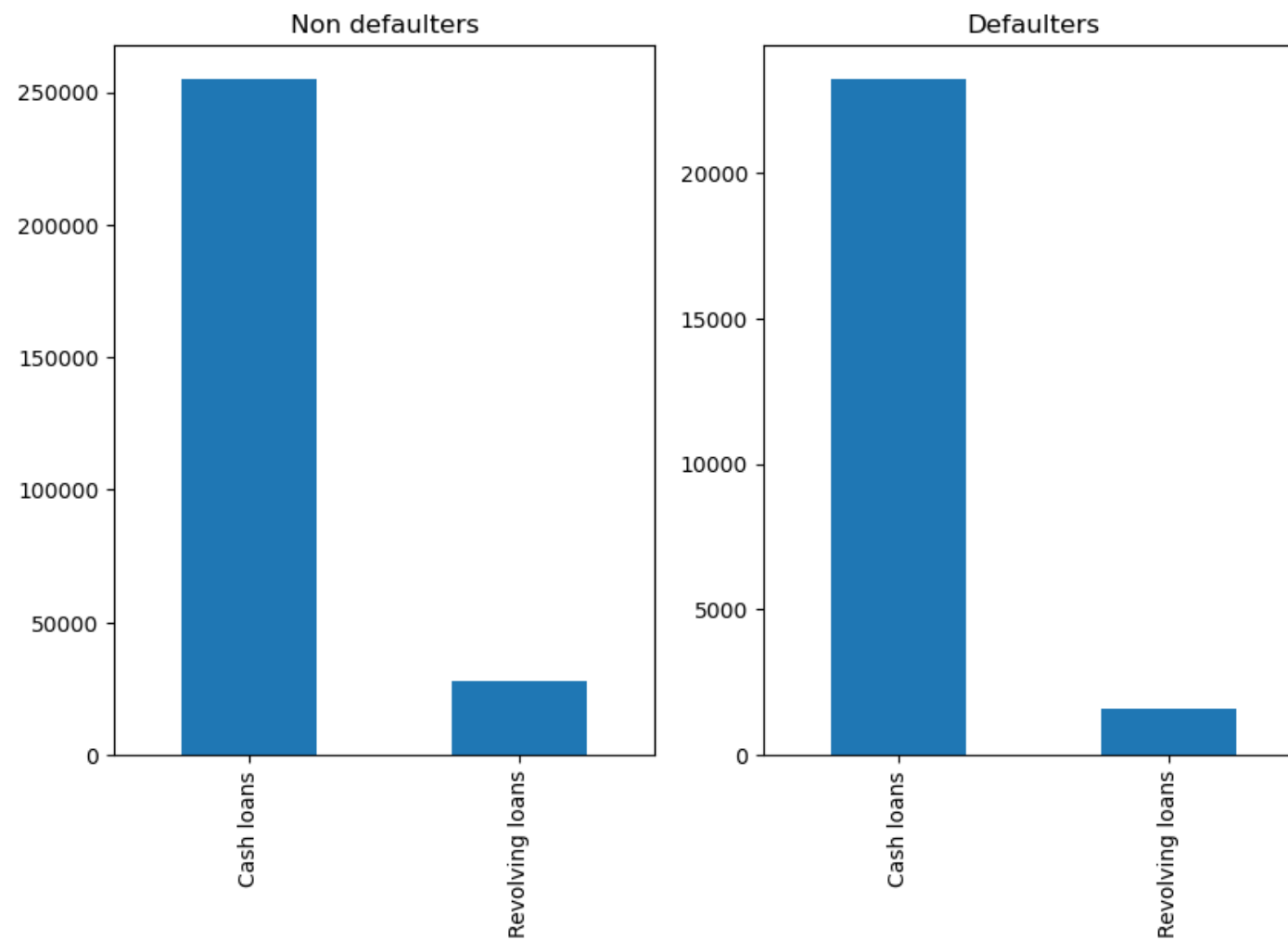


## AGE\_BY\_DEFAULT



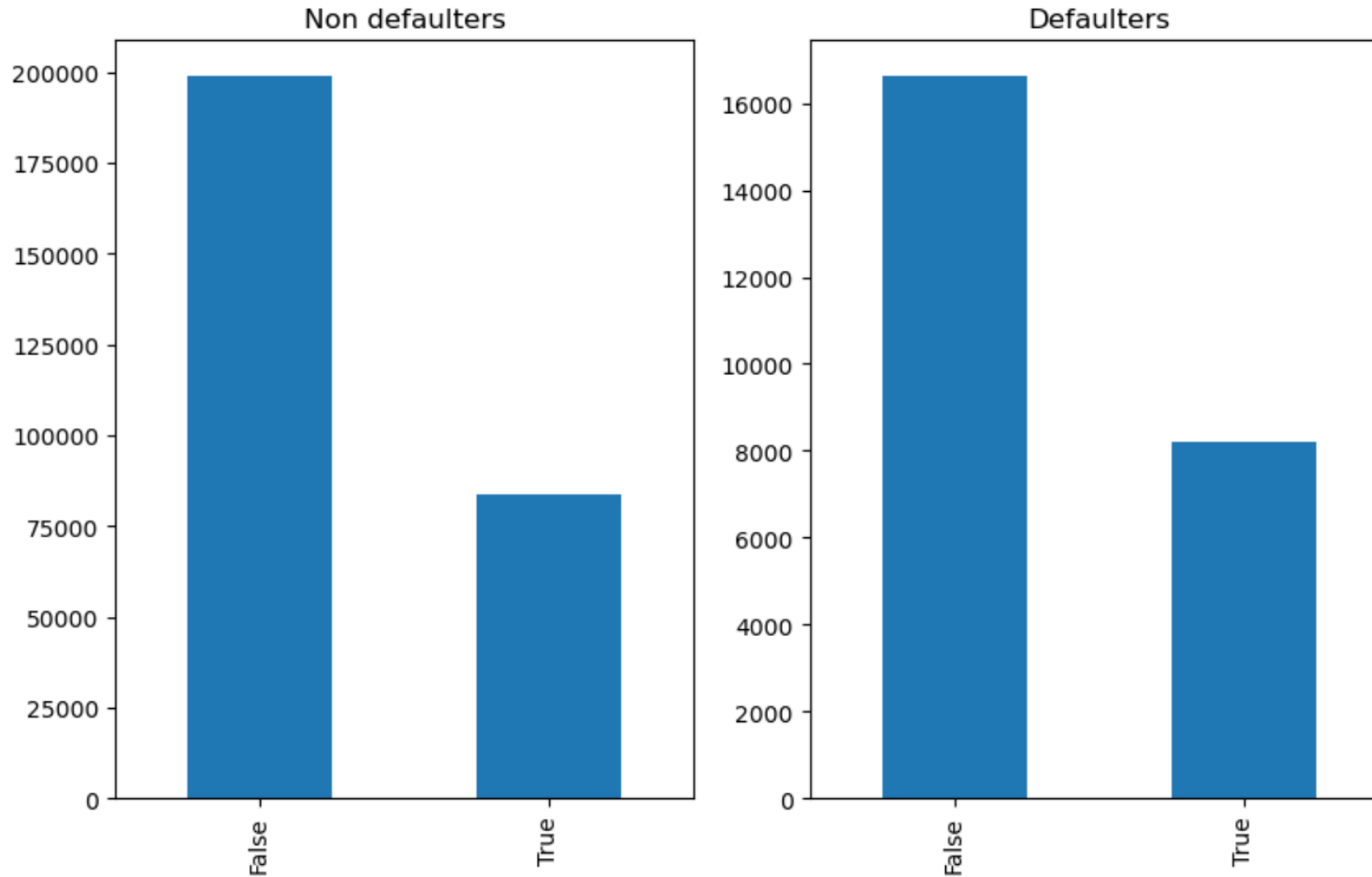
- The people in the age category of mid 30s and 40s contribute nearly 50% of the total people applying loans.
- The plot shows that 30s age category people mostly fail to repay the loan as compared to other age categories.
- The young and the older aged people are less likely to face difficulties in repaying the loan.

## NAME\_CONTRACT\_TYPE



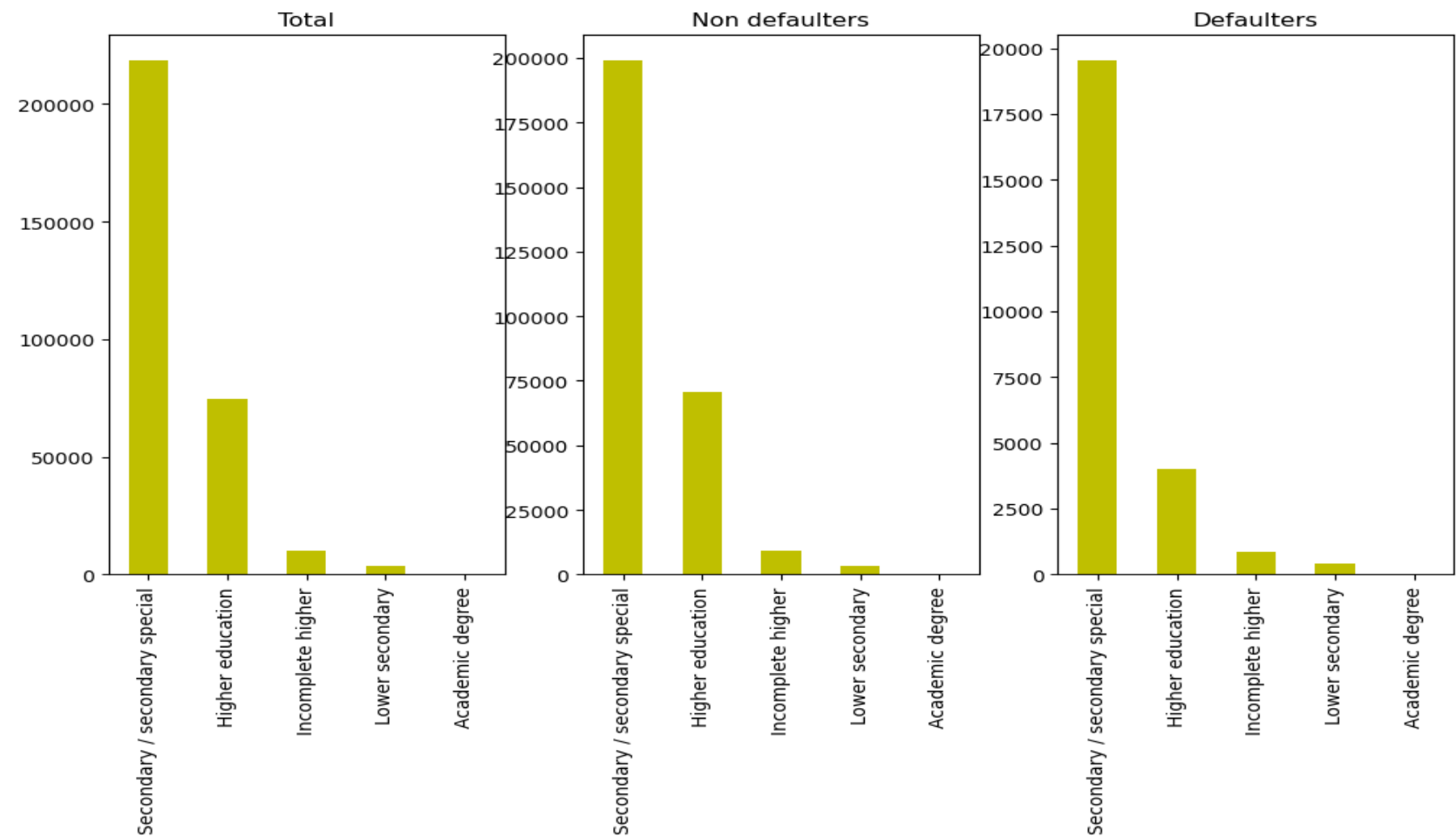
The plot shows that most of the people opted for cash loans. Only a small proportion of people opted revolving loans for both defaulters and non defaulters.

## HAS\_CHILDREN



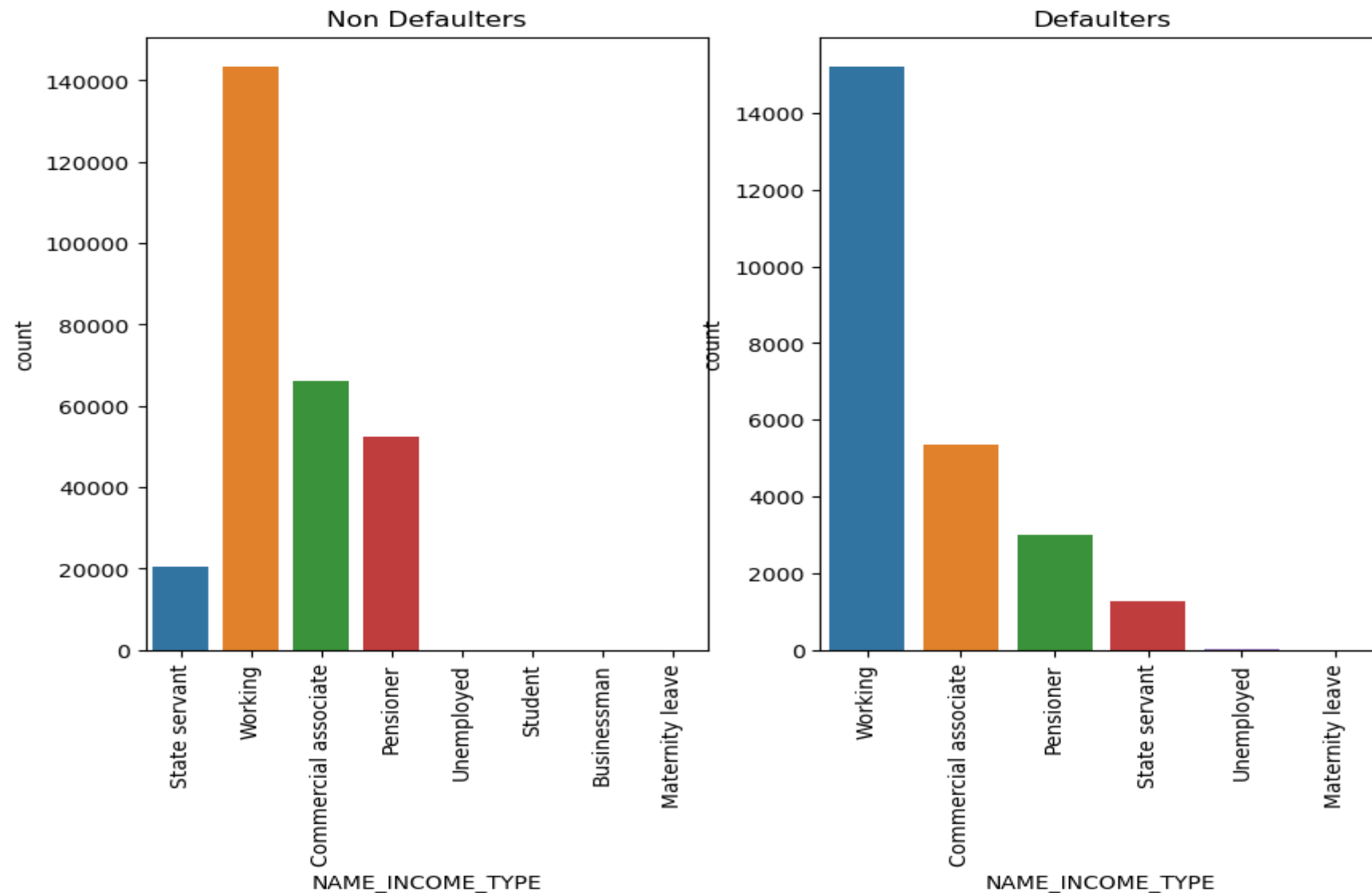
From the plot it is shown that both the defaulters and non defaulters the person who does not have children tends to apply for bank loan than people who has children.

# NAME\_EDUCATION\_TYPE



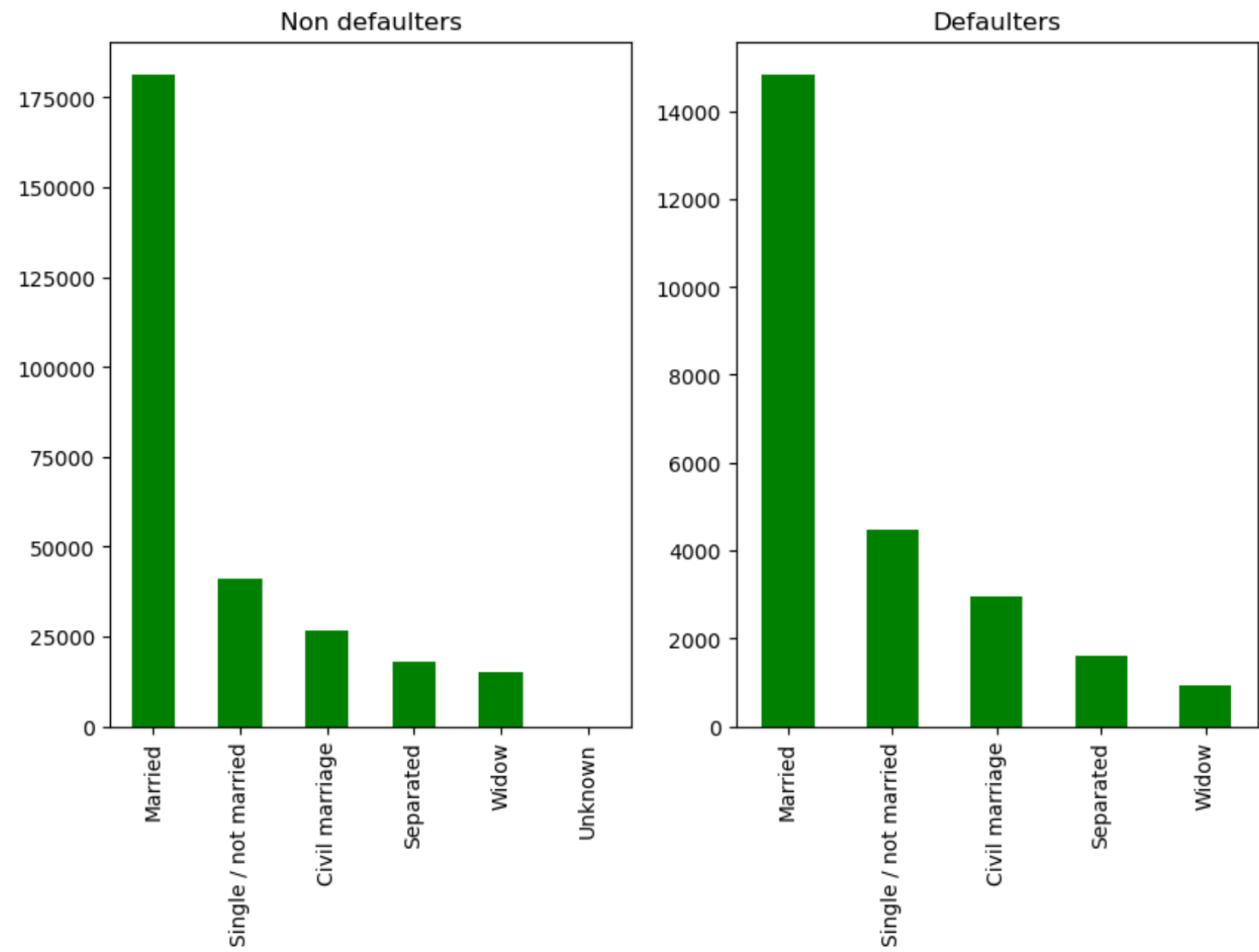
From the bar graph it is observed that people with Secondary education level apply for loans than other qualified people and they also face difficulties in repaying the loan.

## NAME\_INCOME\_TYPE



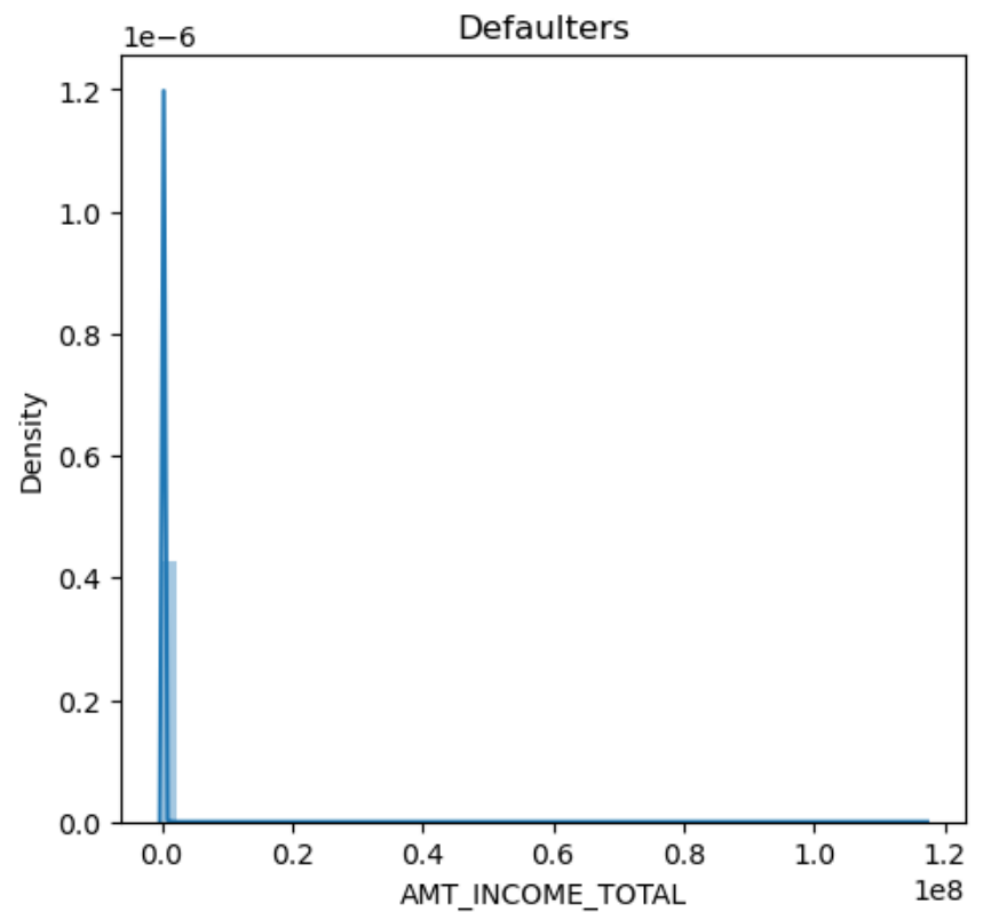
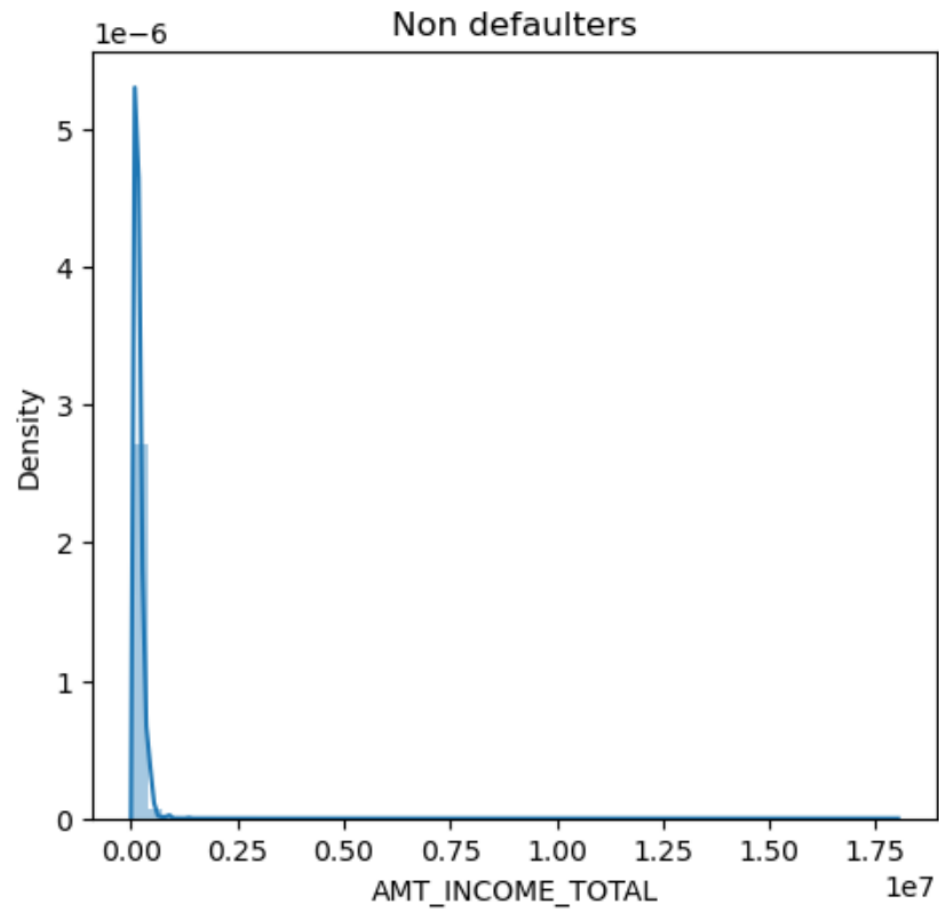
- The plot shows that the Working professionals are in higher proportions for bank loan requirements than any other professionals and there is a higher risk for loan default.
- The pensioner people and State servant applying for loan has lesser risk to default.

# NAME\_FAMILY\_STATUS



The graph shows that the Married people require bank loan than others in case of both defaulters and non defaulters and Widowed people who applies for loan has lesser risk to loan default.

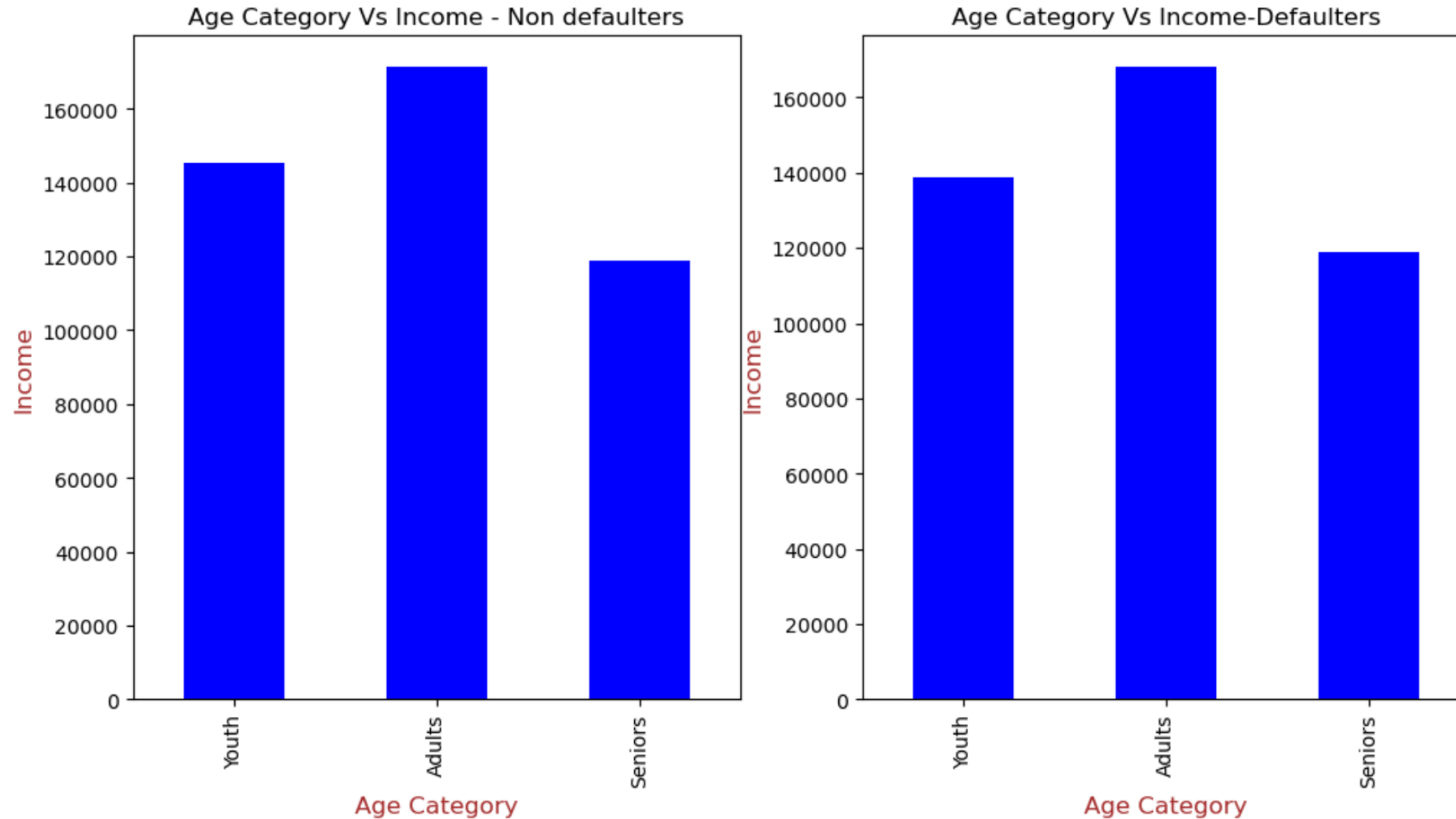
## AMT\_INCOME\_TOTAL



- In summary, the plot suggests that there is a difference in the distribution of income between non-defaulters and defaulters.
- Defaulters have a wider range of income levels, including some with high incomes, while non-defaulters have income levels more concentrated in the lower to middle range.

# Bivariate and Multivariate analysis

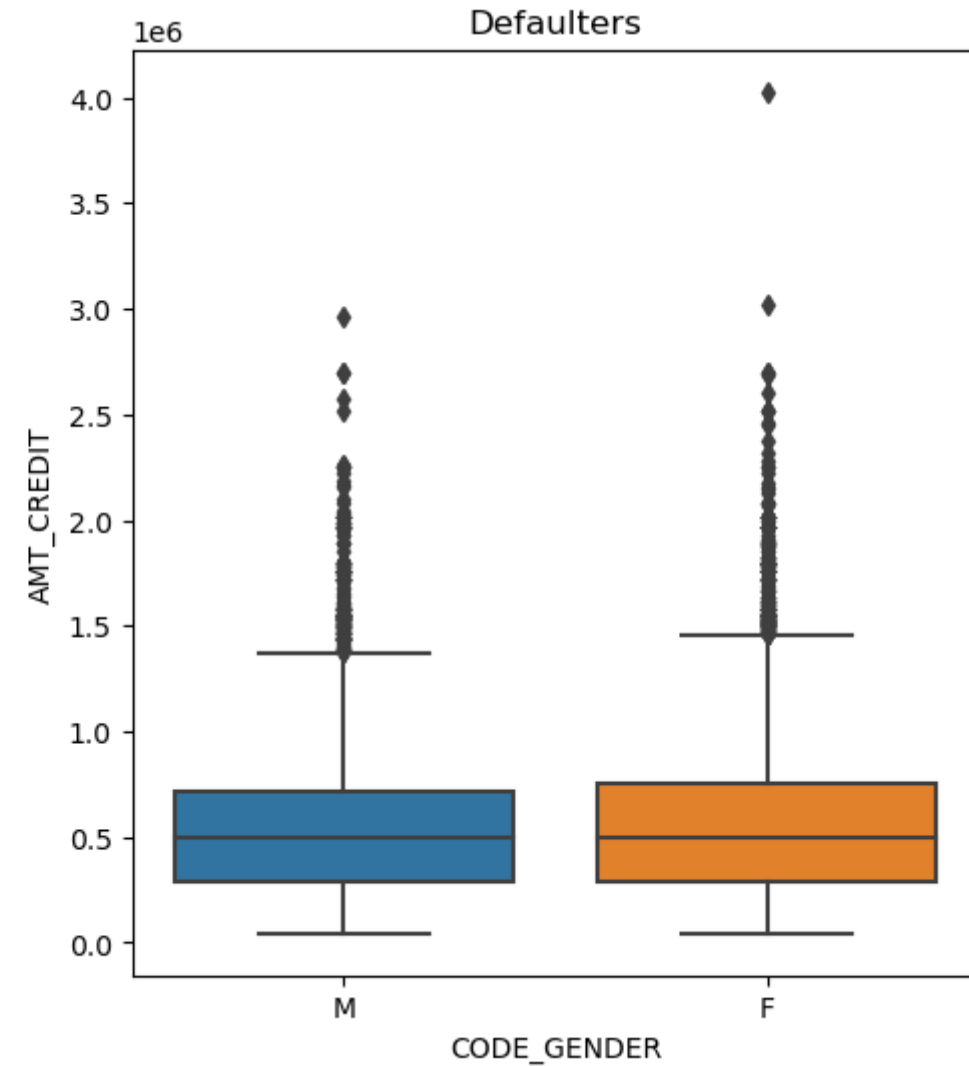
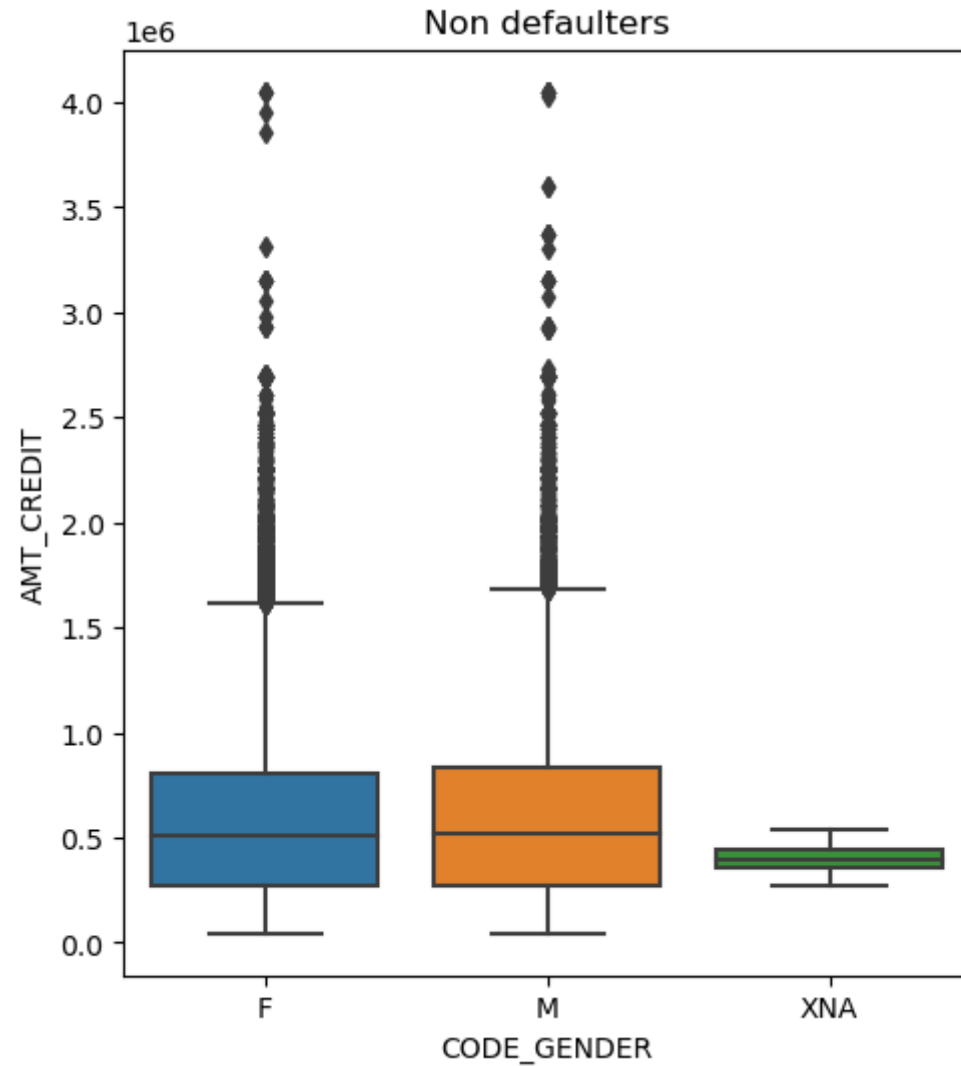
## Categorizing the Income of the persons with their Age



- In both plots, there is a clear trend where the average income tends to increase with age, up to a certain point, and then it starts to stabilize or slightly decrease in the older age categories.
- Younger age categories have lower average incomes, while older age categories generally have higher average incomes.
- For non-defaulters, the average income generally increases with age. For defaulters, the trend is similar but less pronounced.

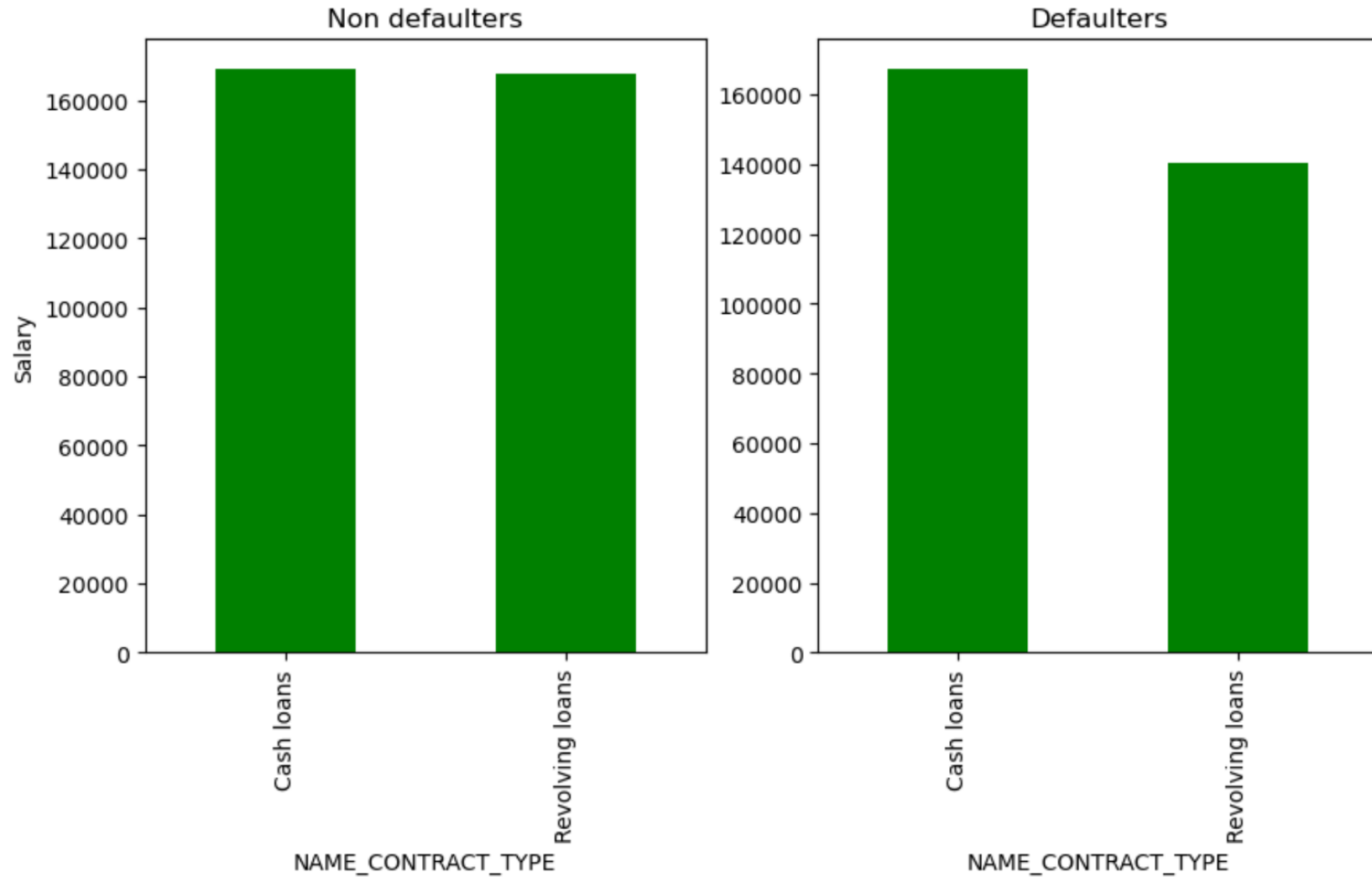


## CODE\_GENDER Vs AMT\_CREDIT



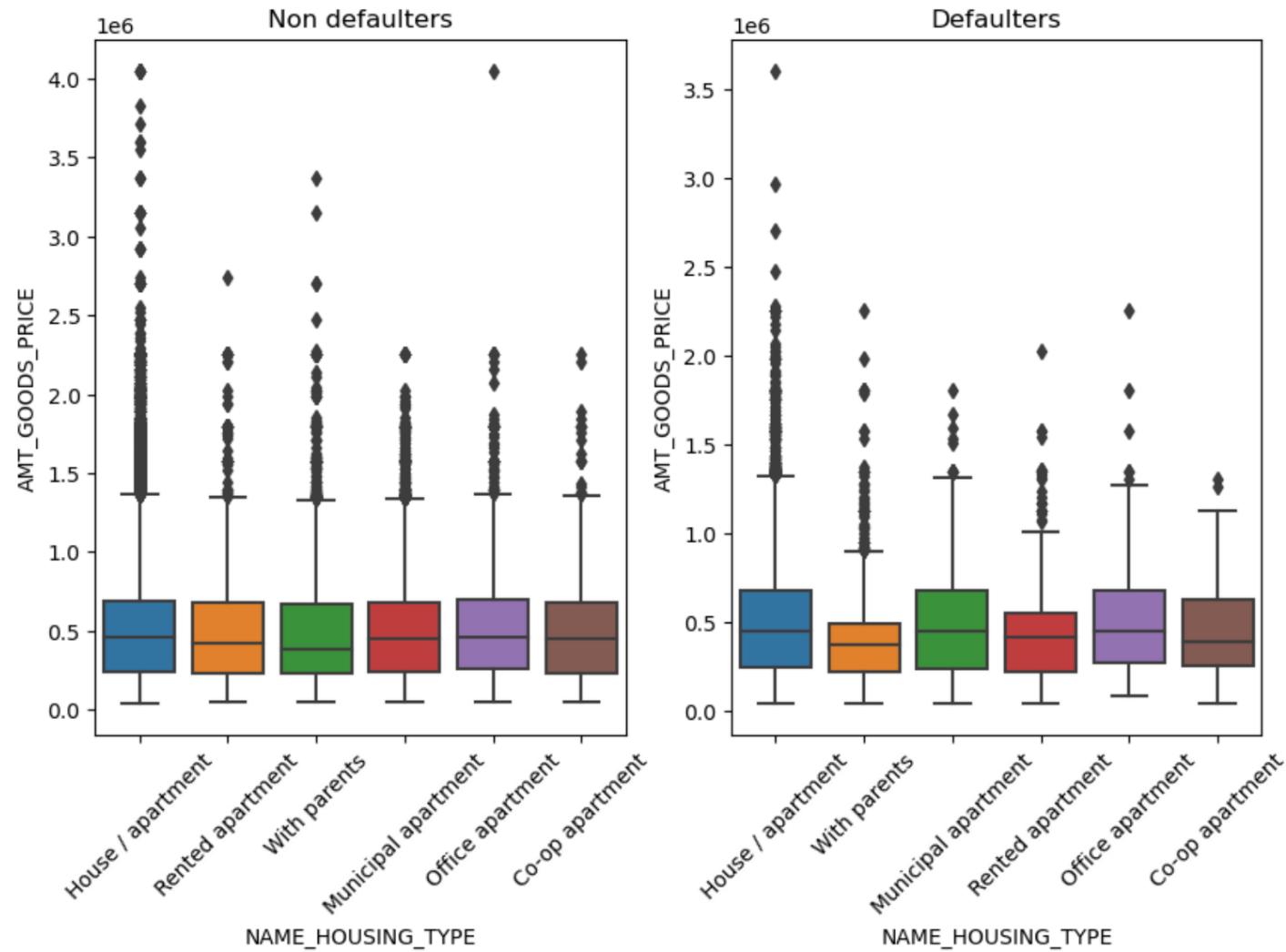
In the Non defaulters plot the credit distribution for both the genders are relatively similar. In the Defaulters plot the credit distribution for males appears to have wider spread than females. Thus there is more variation in the credit amount for male compared to female defaulters.

## NAME\_CONTRACT\_TYPE Vs AMT\_INCOME\_TOTAL



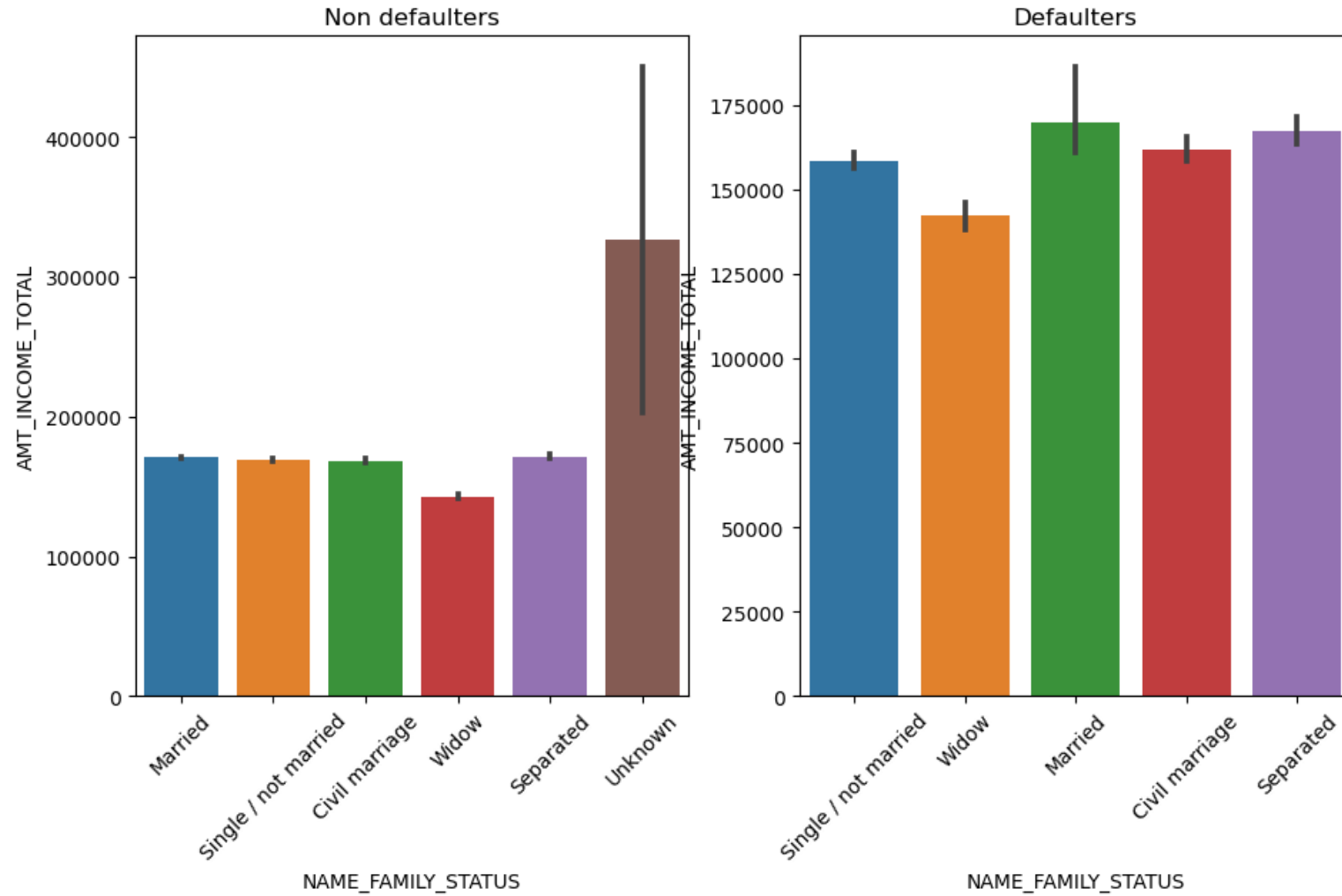
- In summary, the plot suggests that the relationship between contract type and mean income differs between non-defaulters and defaulters.
- Among non-defaulters, those with "Revolving loans" tend to have a higher mean income, while among defaulters, those with "Cash loans" tend to have a higher mean income.

## AMT\_GOODS\_PRICE Vs NAME\_HOUSING\_TYPE



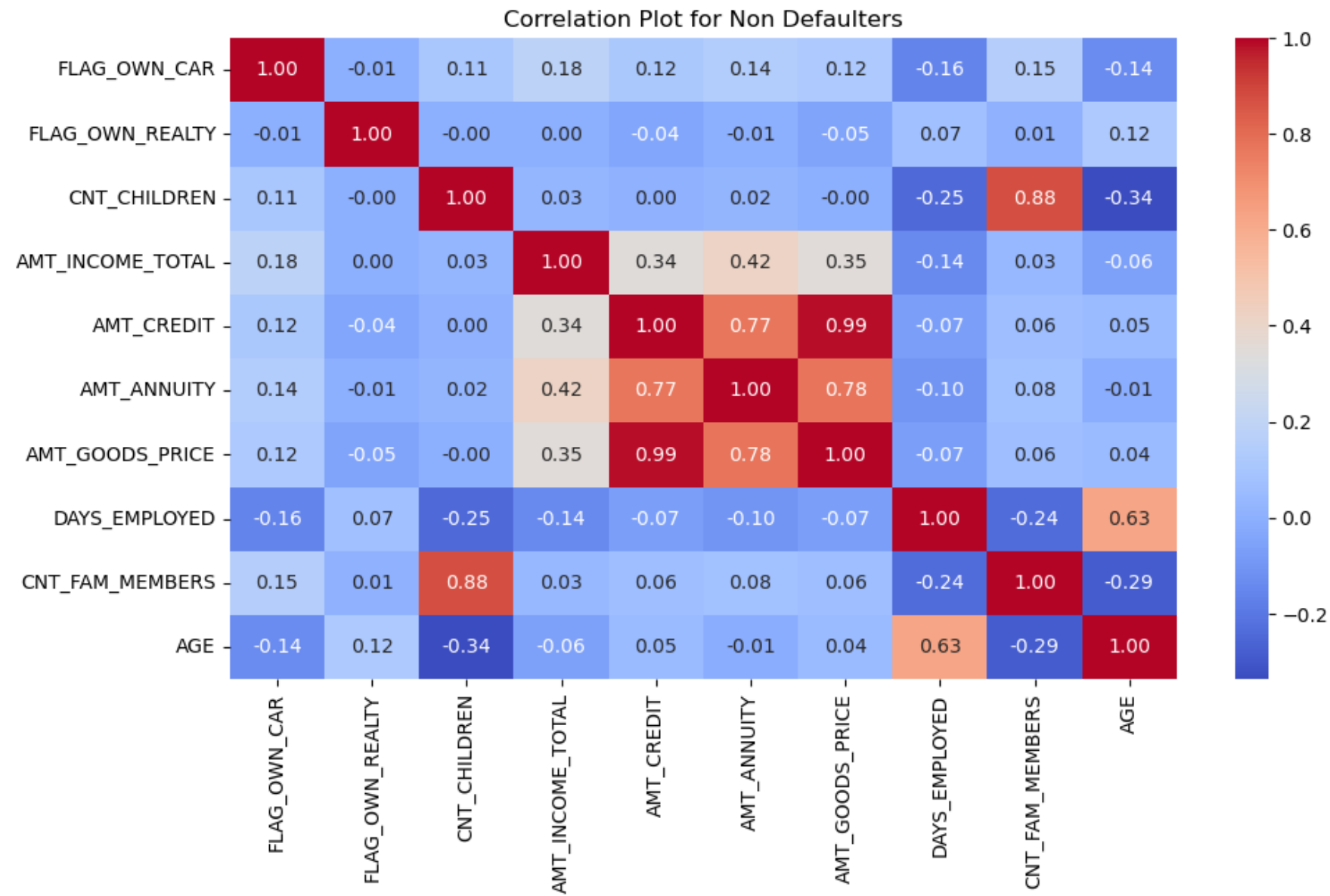
- The persons with House/apartment, Municipal apartment, Office apartment tends to have higher variability than other housing types among defaulters.
- While the non defaulters almost has similar variability of the amount of good owned for all the housing types and does not exhibit much variability in their distribution.

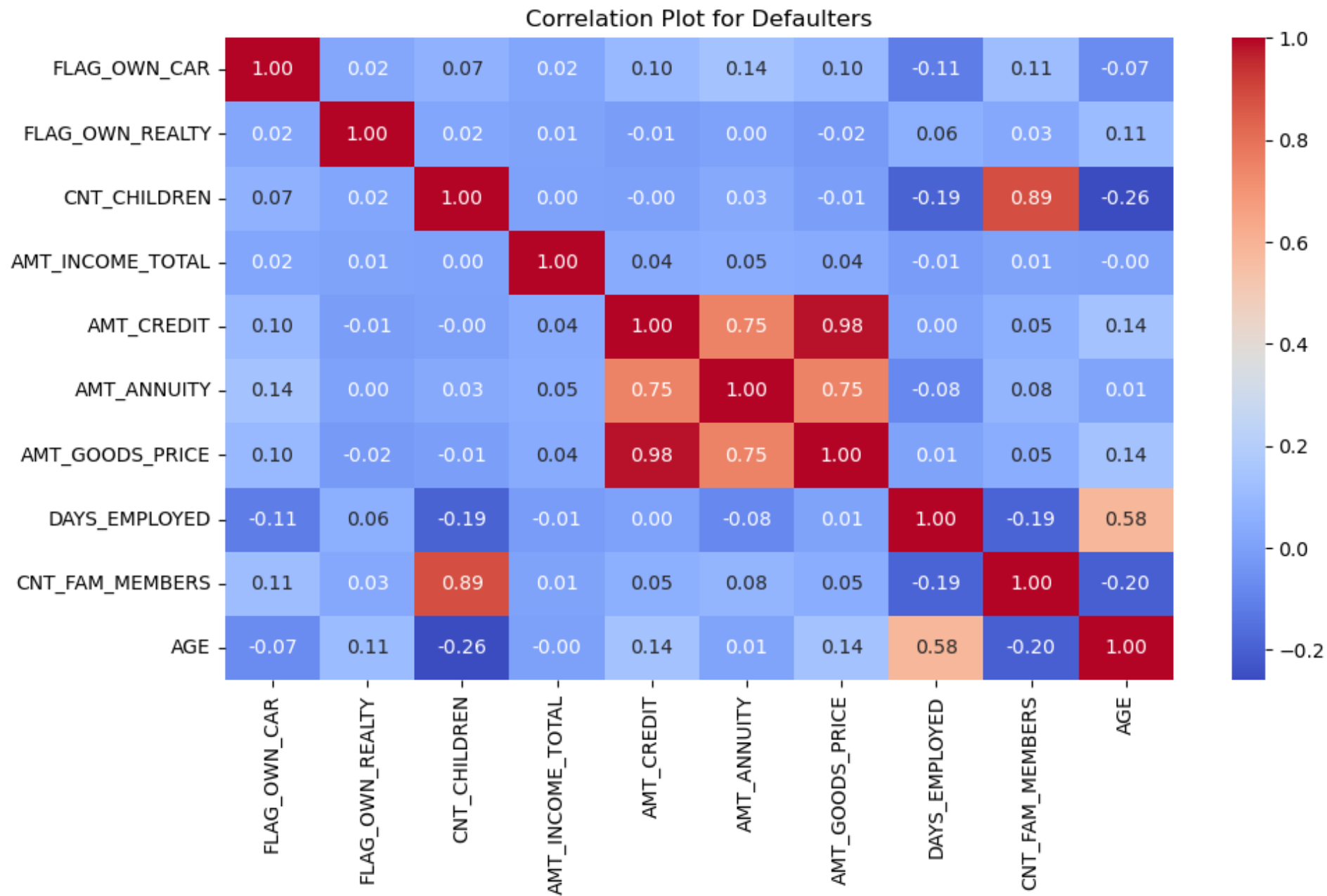
## NAME\_FAMILY\_STATUS Vs AMT\_INCOME\_TOTAL



1. Civil marriage and Married tend to have higher mean incomes compared to other family statuses.
2. Single/not married and Separated have lower mean incomes than the others for both defaulters and non defaulters.

Correlation Plots

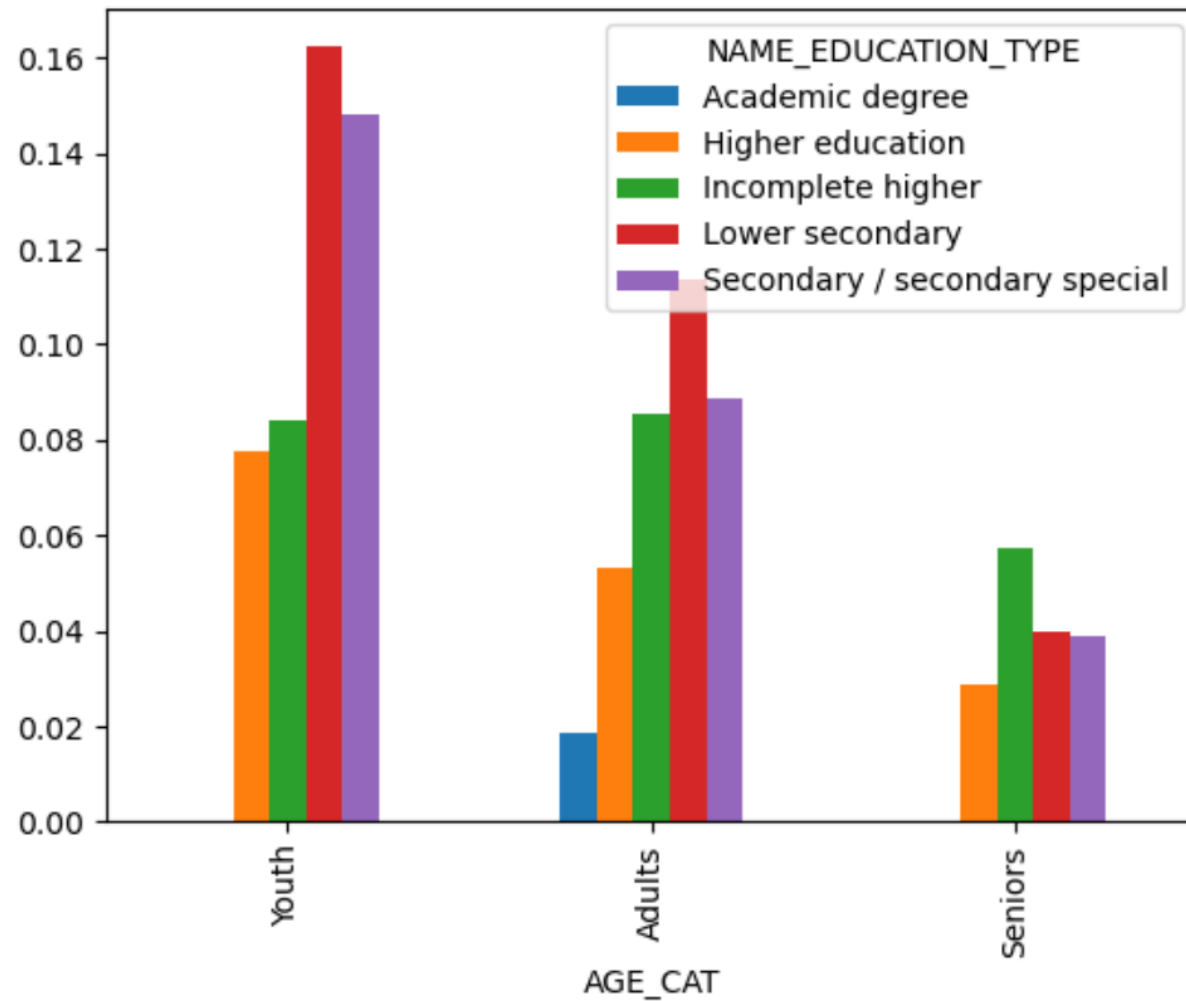




## Insights from heat maps Defaulters and Non defaulters

1. AMT\_CREDIT and AMT\_GOODS\_PRICE has strong positive correlation which implies that the loan amount is increased when the price of goods/asset increases for both Defaulters and Non defaulters
2. CNT\_CHILDREN and CNT\_FAM\_MEMBERS are positively correlated with each other. When the count of children in the family increases , the family members count also gets increased for both Defaulters and Non defaulters.
3. AMT\_CREDIT and AMT\_ANNUITY are directly proportional to each other. The total annual loan payment increases when the credit amount is high.
4. AMT\_GOODS\_PRICE, AMT\_CREDIT are negatively correlated with CNT\_CHILDREN which shows that the loan amount is higher for those who has fewer children.
5. AMT\_ANNUITY are inversely proportional to DAYS\_EMPLOYED which shows that the annual loan amount to be paid is higher for people employed for less number of days.

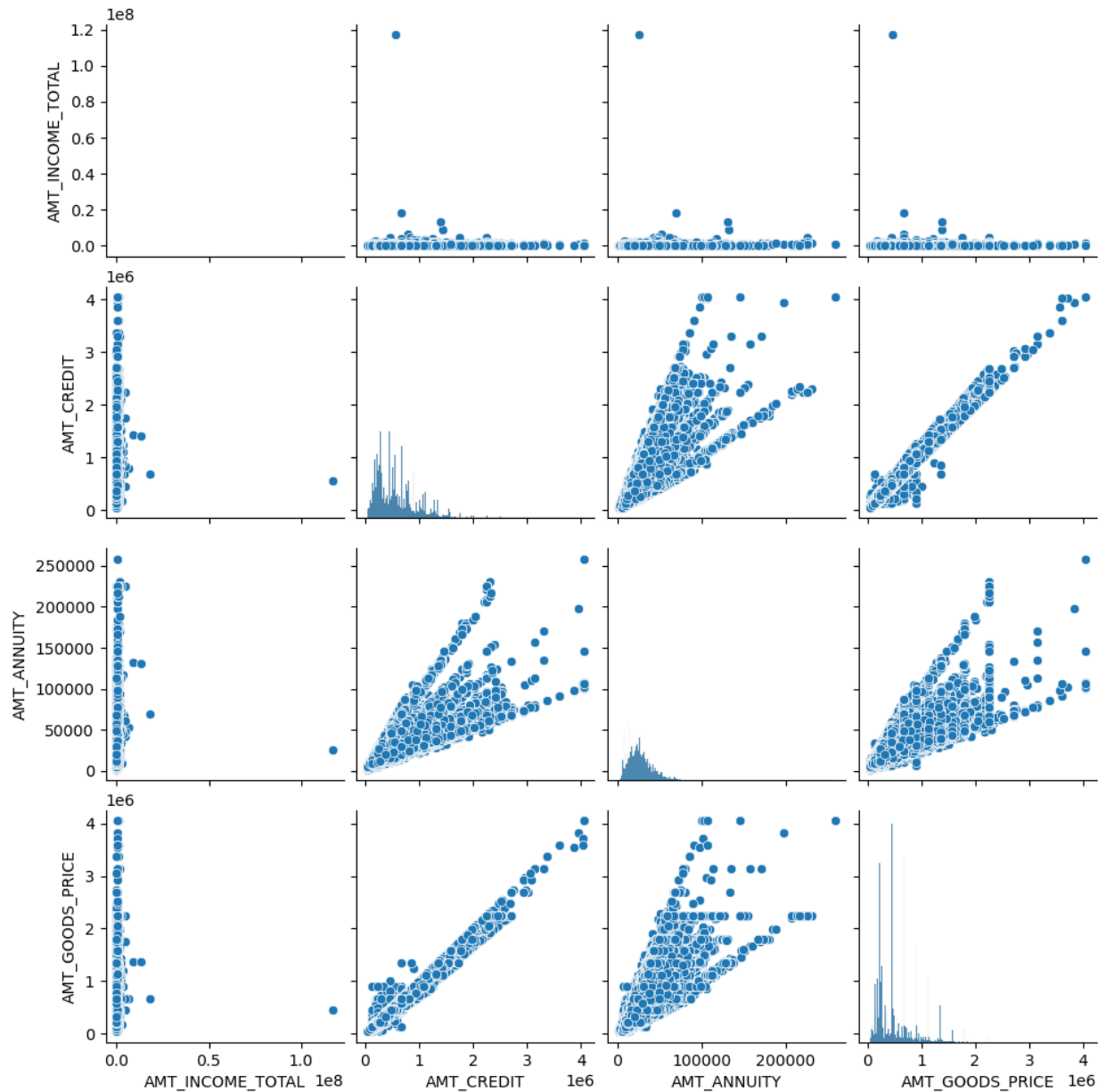
## Analysing loan defaulters by Education Type and Age Category



1. The plot shows that the youth who completed lower secondary education exposes high risk of loan default.
2. The Senior age group people who are highly educated shows minimal risk of repayment.



## Pairplots of continuous variables



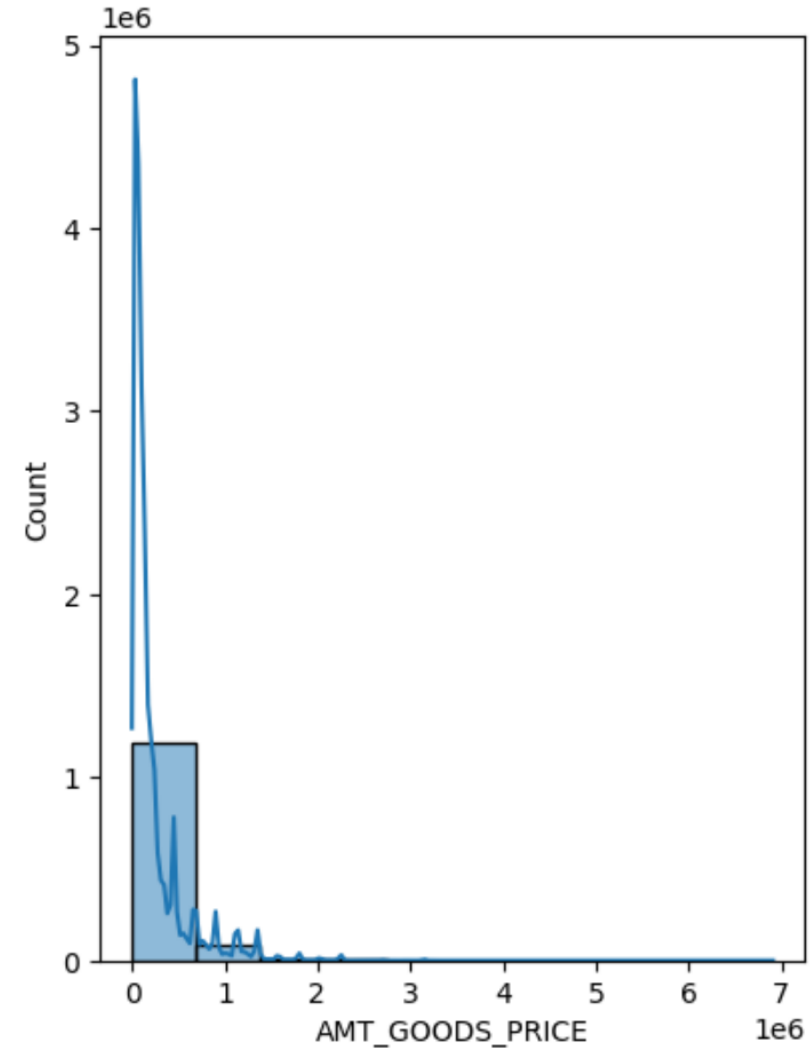
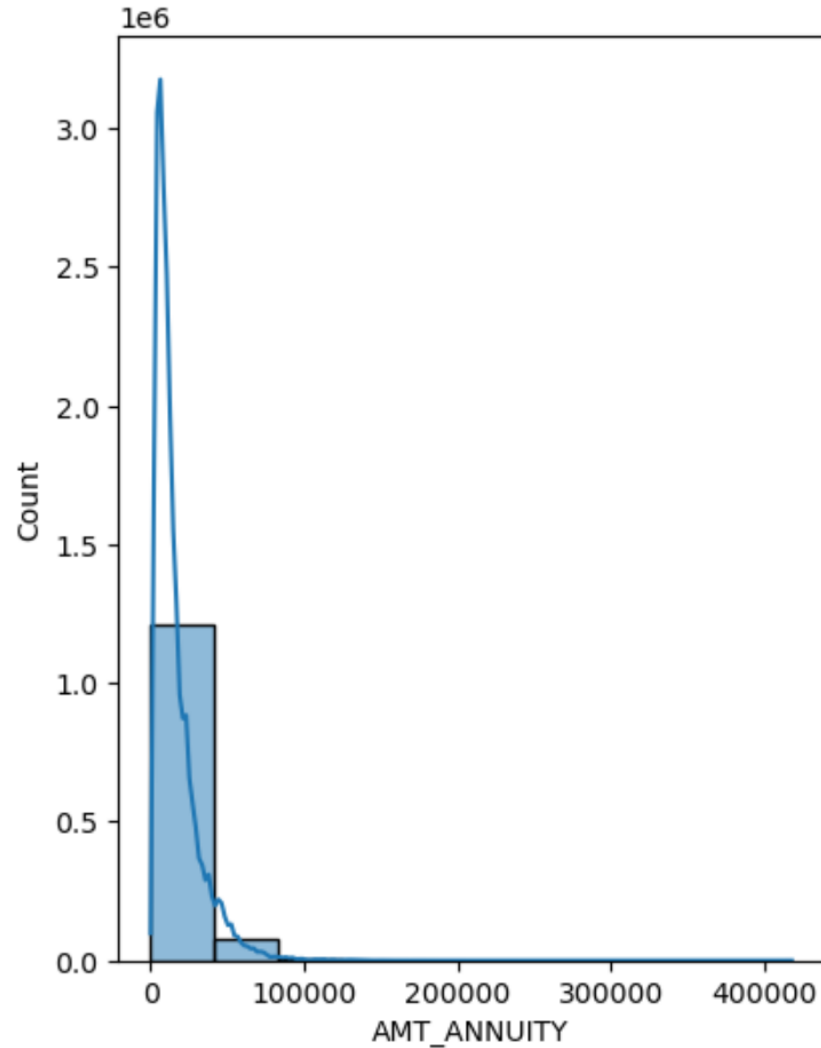
### Inference

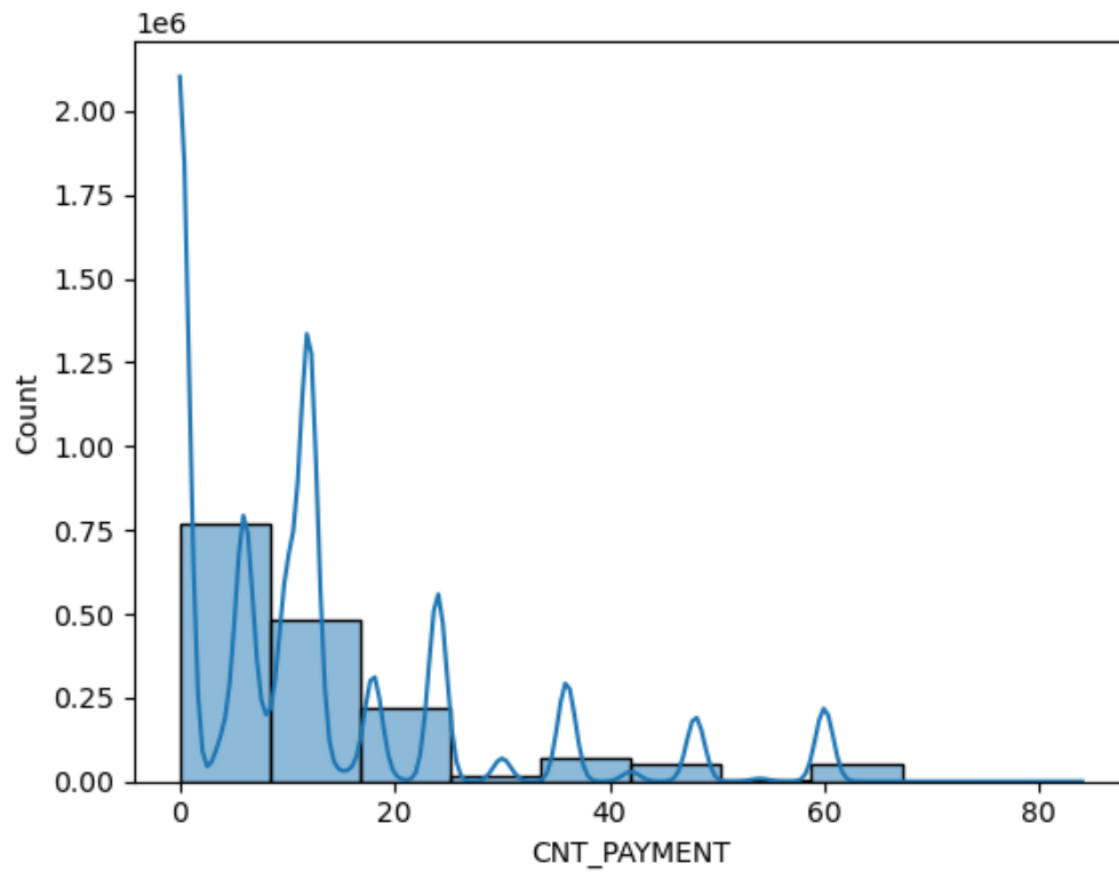
1. AMT\_CREDIT and AMT\_GOODS\_PRICE shows strong correlated with each other as an increase in the credit amount is typically associated with a higher goods price.
2. AMT\_ANNUIITY and AMT\_CREDIT may reveal a positive linear relationship, where higher credit amounts are associated with higher annuities.

## Previous\_application

### Handling missing values

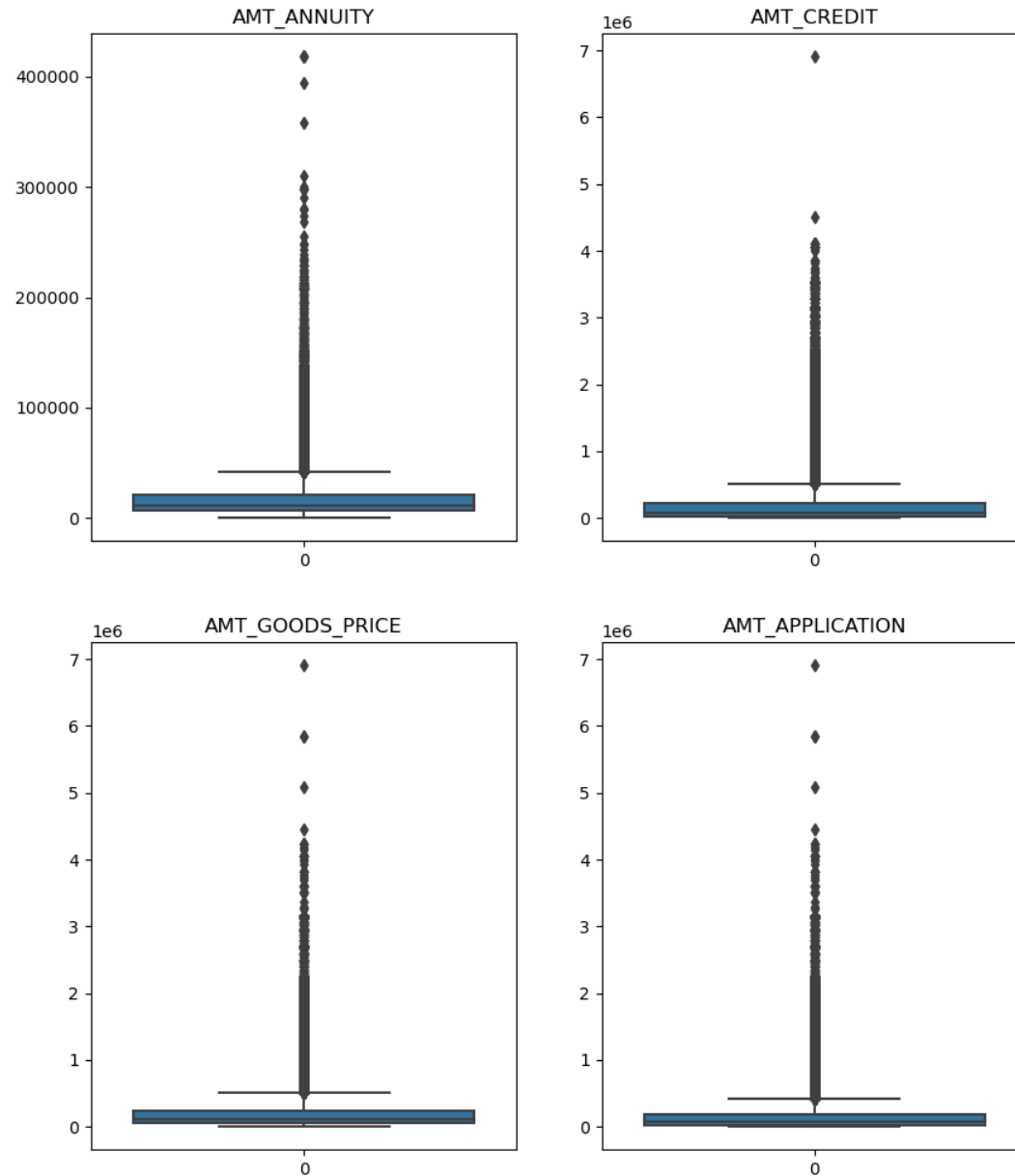
columns whose missing value is < 13% are imputed





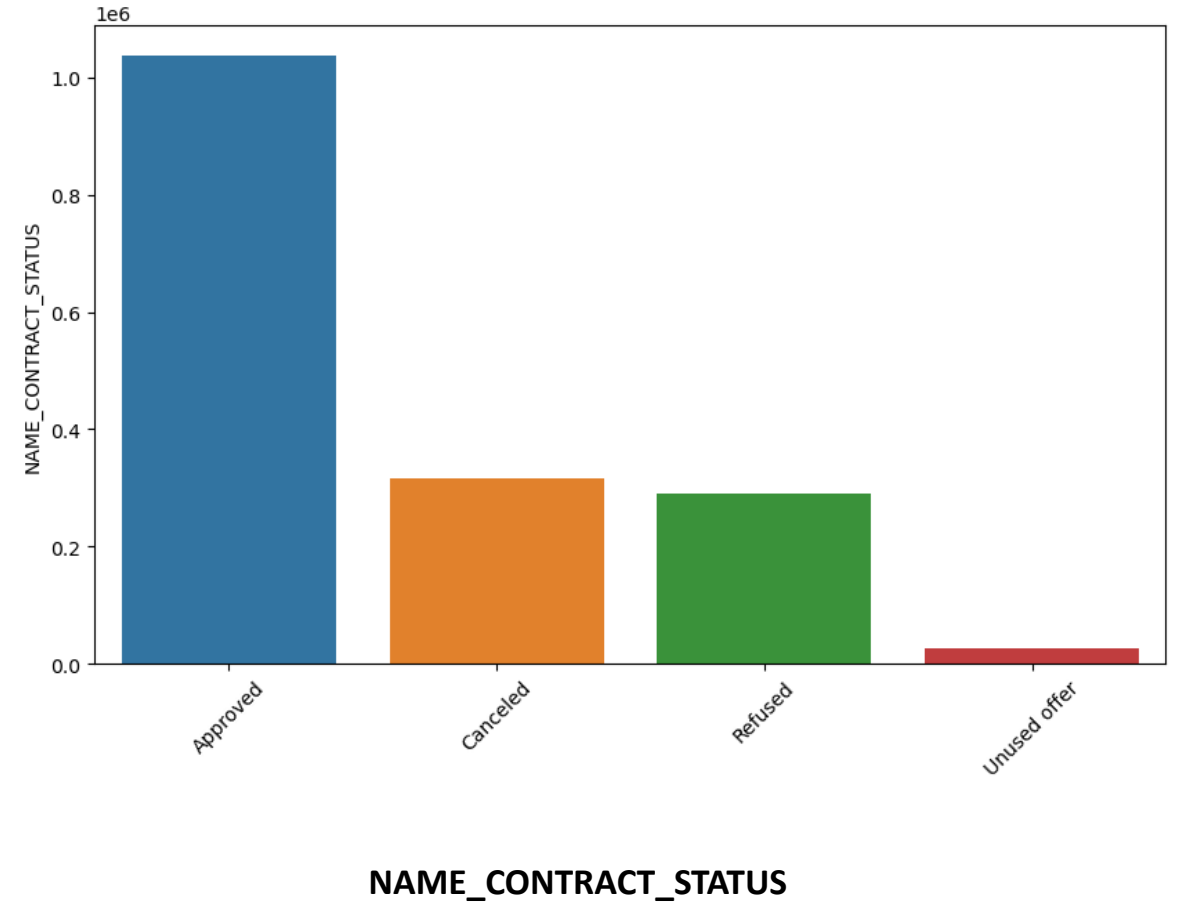
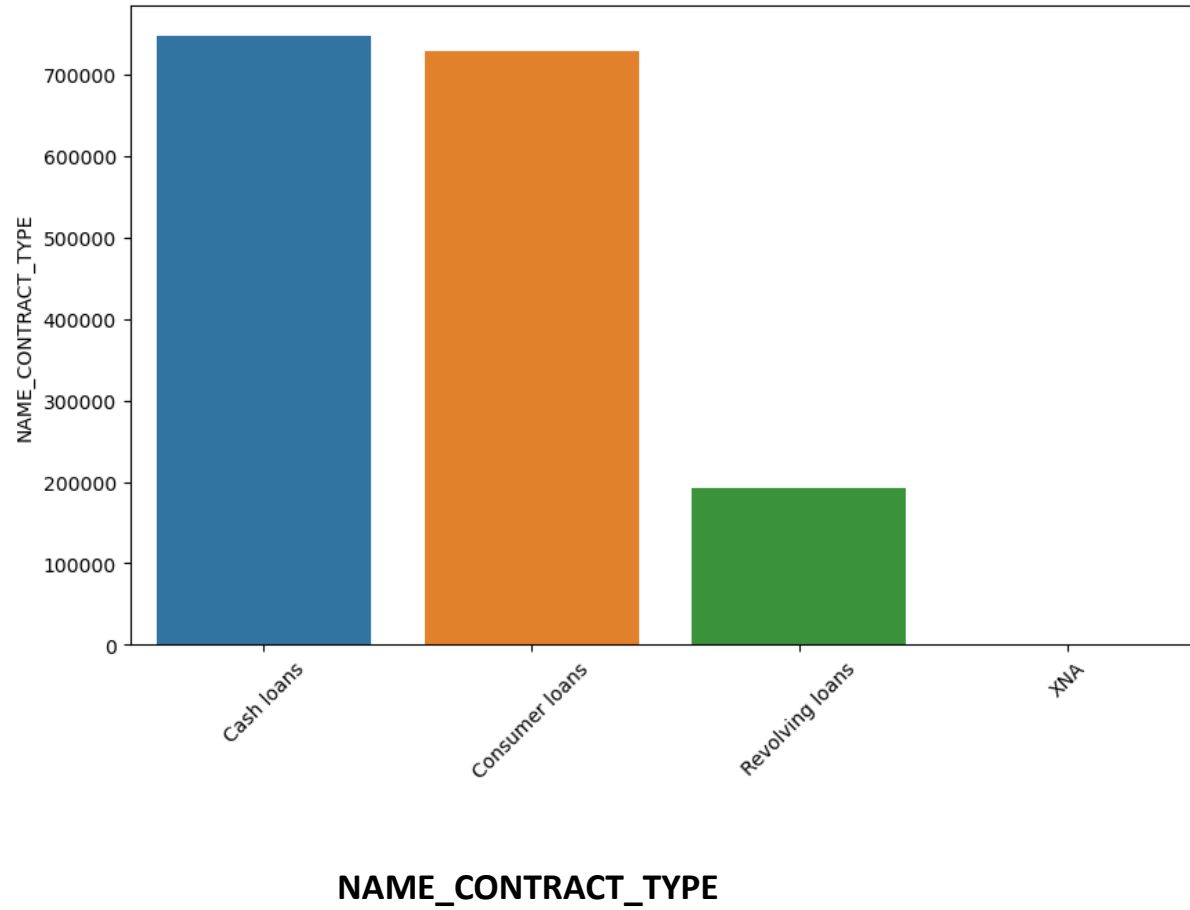
The distribution plot of the variable AMT\_ANNUITY, AMT\_GOODS\_PRICE and CNT\_PAYMENT is rightly skewed and hence the null values can be replaced with median value.

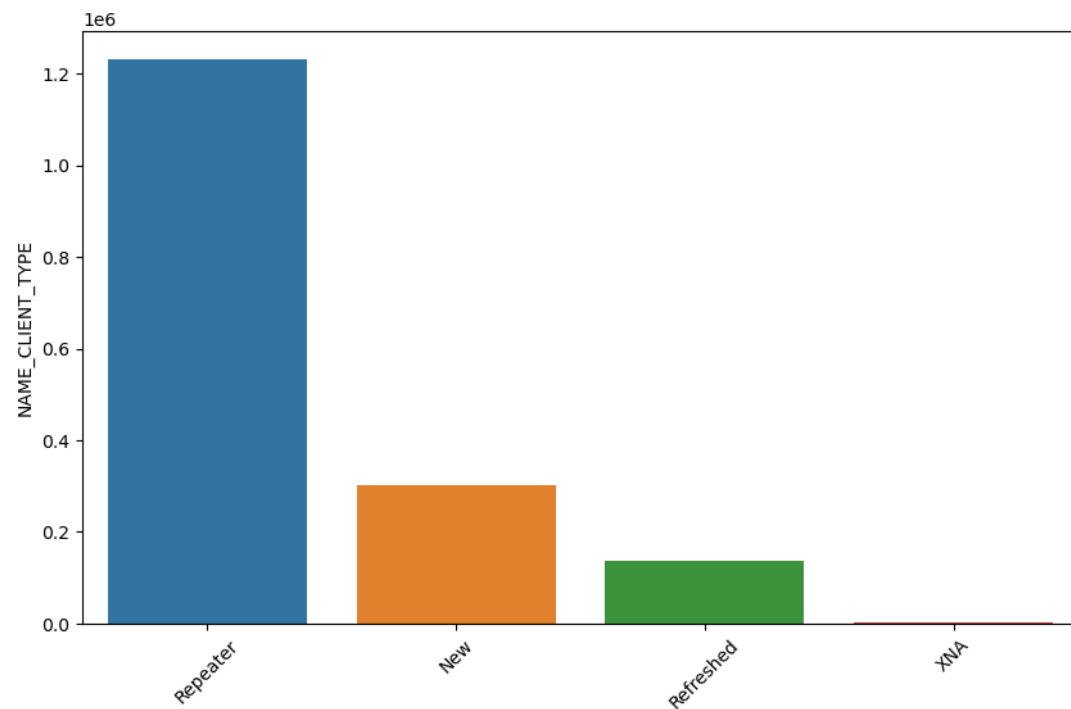
# Handling outliers



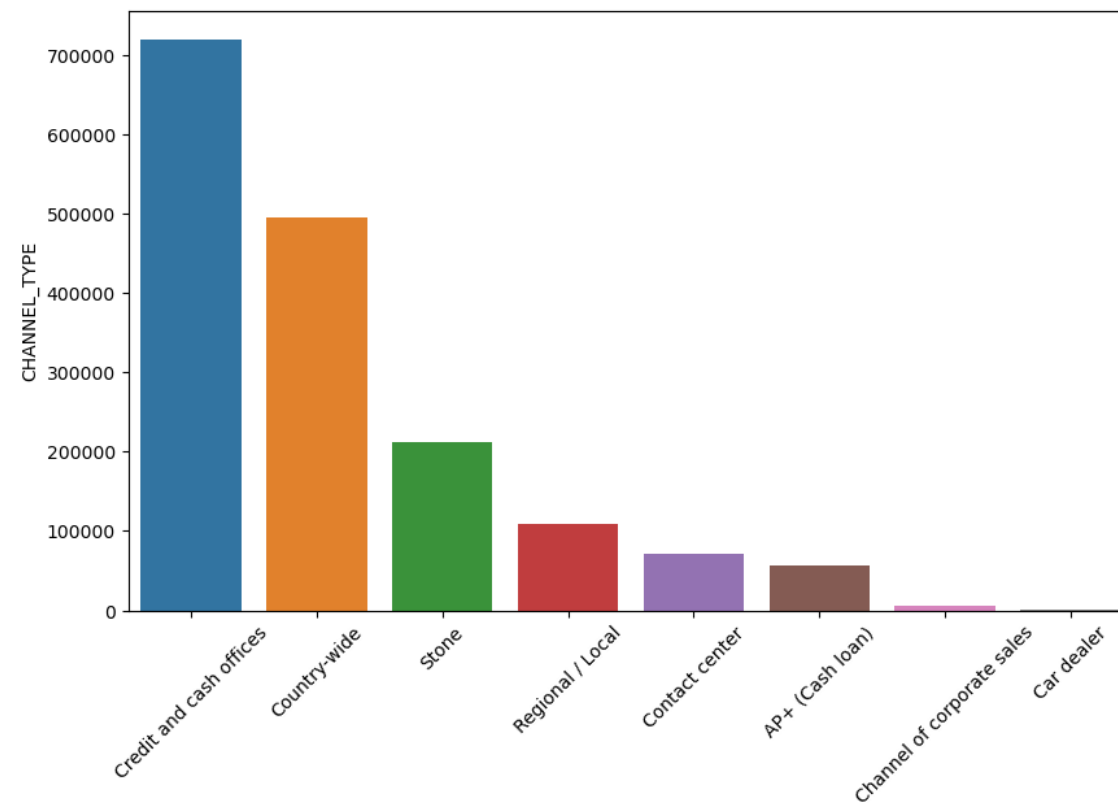
Since the percentage of outliers are very less compared to the total data, the outliers in all the variables can be dropped off, since it does not have an effect on the overall analysis.

## Categorical variables

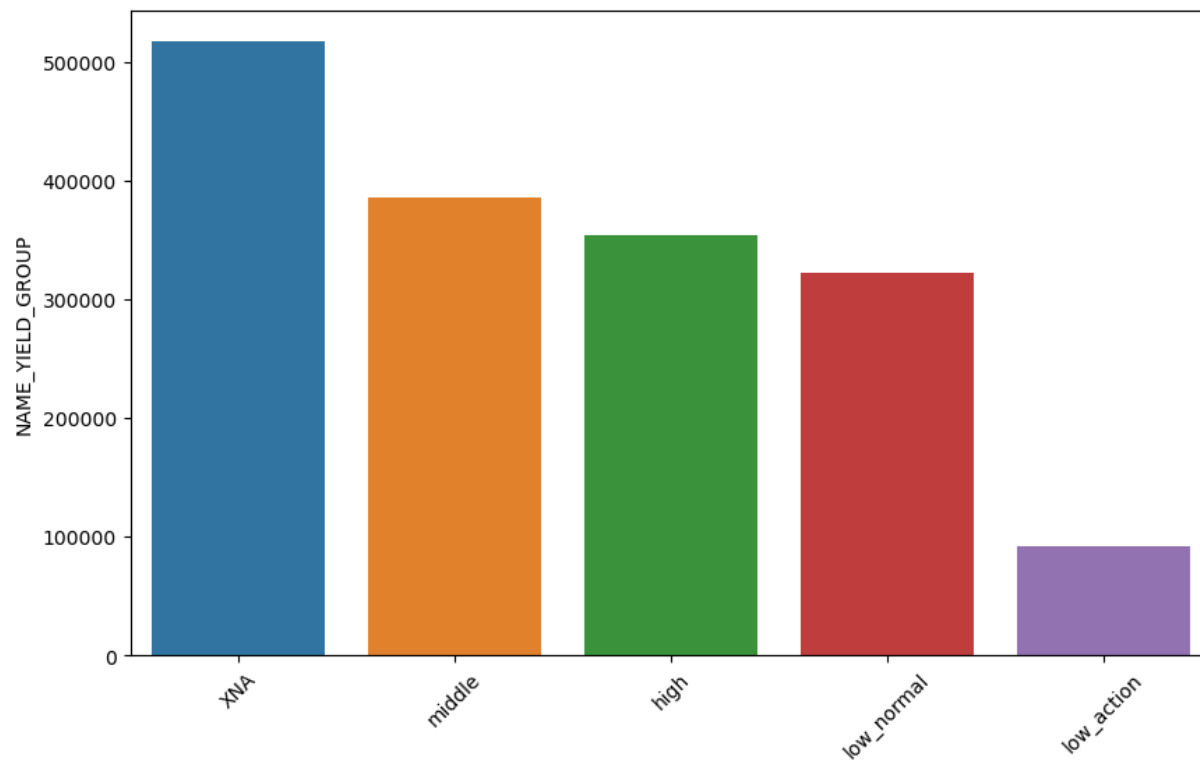




NAME\_CLIENT\_TYPE



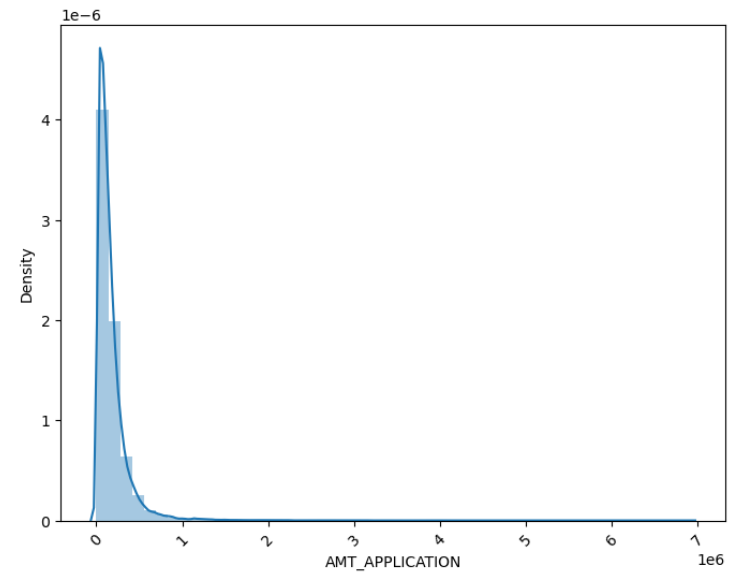
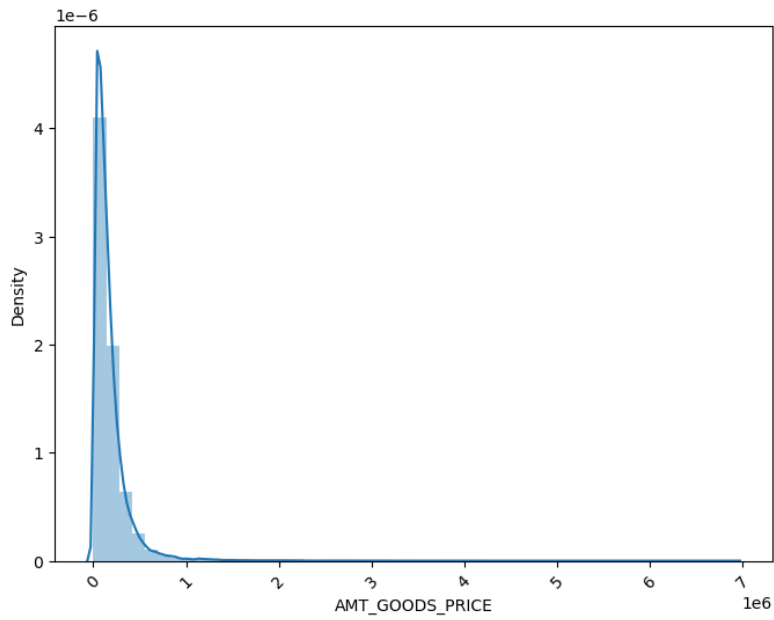
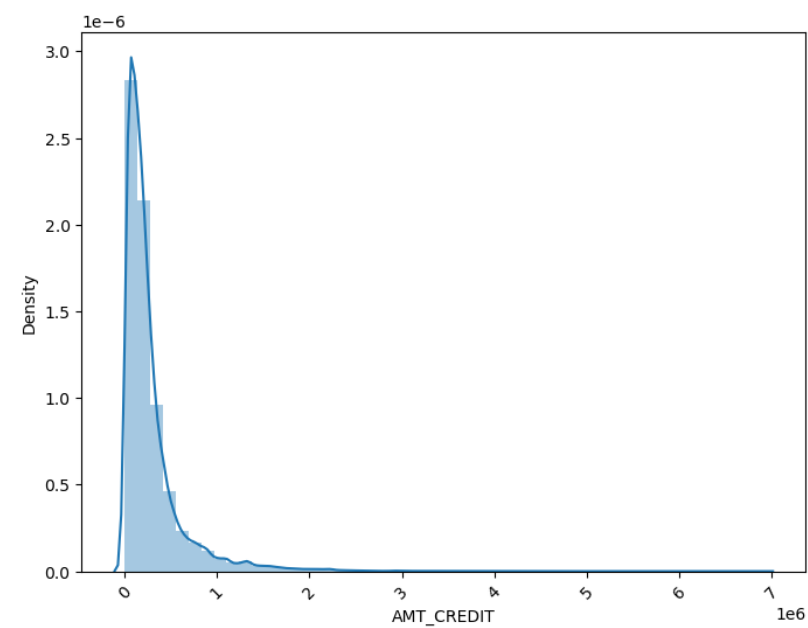
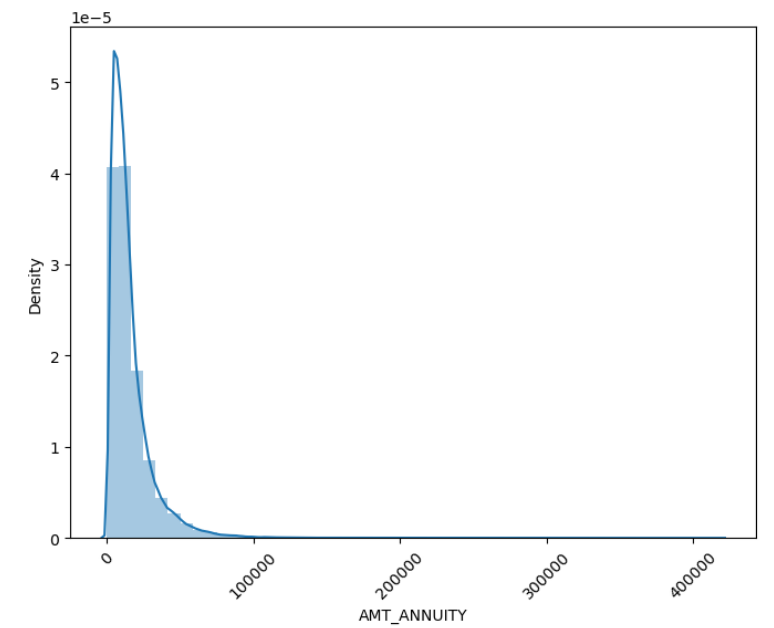
CHANNEL\_TYPE



### Insights from the Categorical variables

1. NAME\_CONTRACT\_TYPE: Cash loans are higher followed by Consumer loans
2. NAME\_CONTRACT\_STATUS: Most of the previous loan application status are approved and only a meagre amount of applicant did not use the offer.
3. NAME\_CLIENT\_TYPE: Most of the clients are already a customer of the bank
4. CHANNEL\_TYPE: The clients approached by the Credit and cash offices are higher in proportions.
5. NAME\_YIELD\_GROUP: The interest were in the middle category for most of the loan applicants.

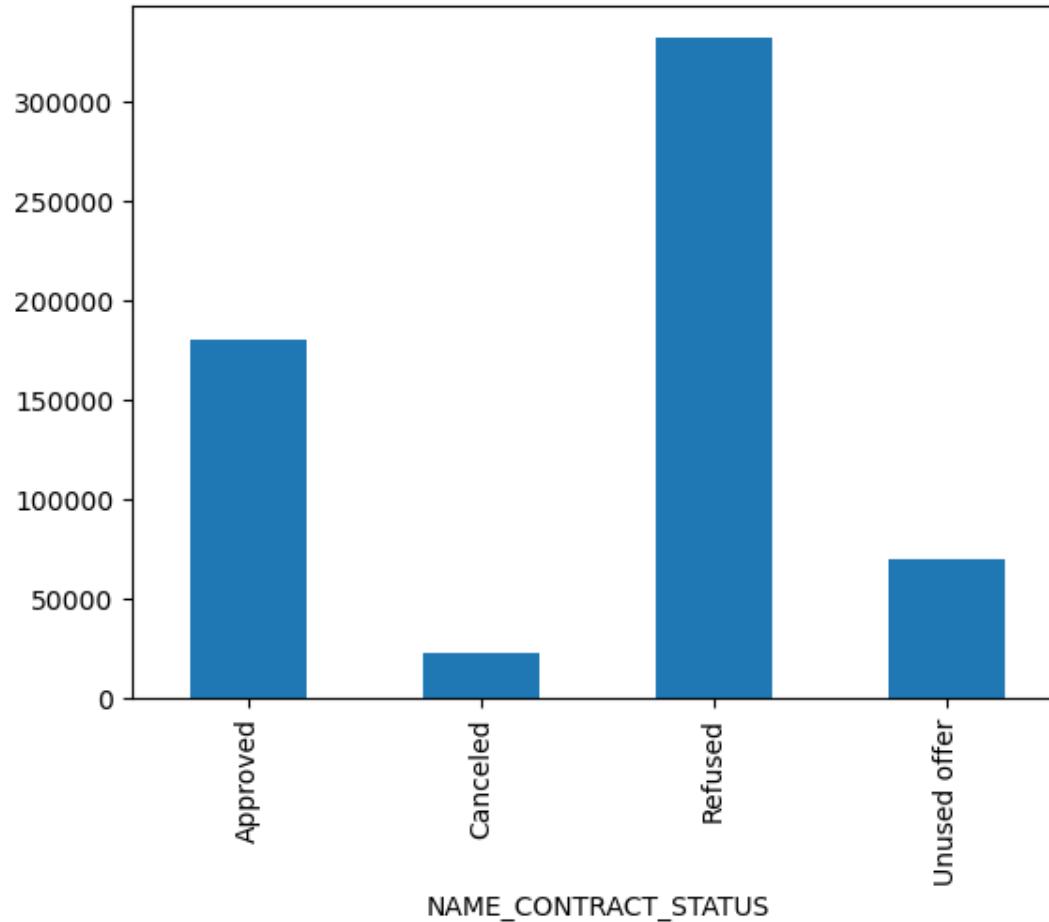
# Numerical variables





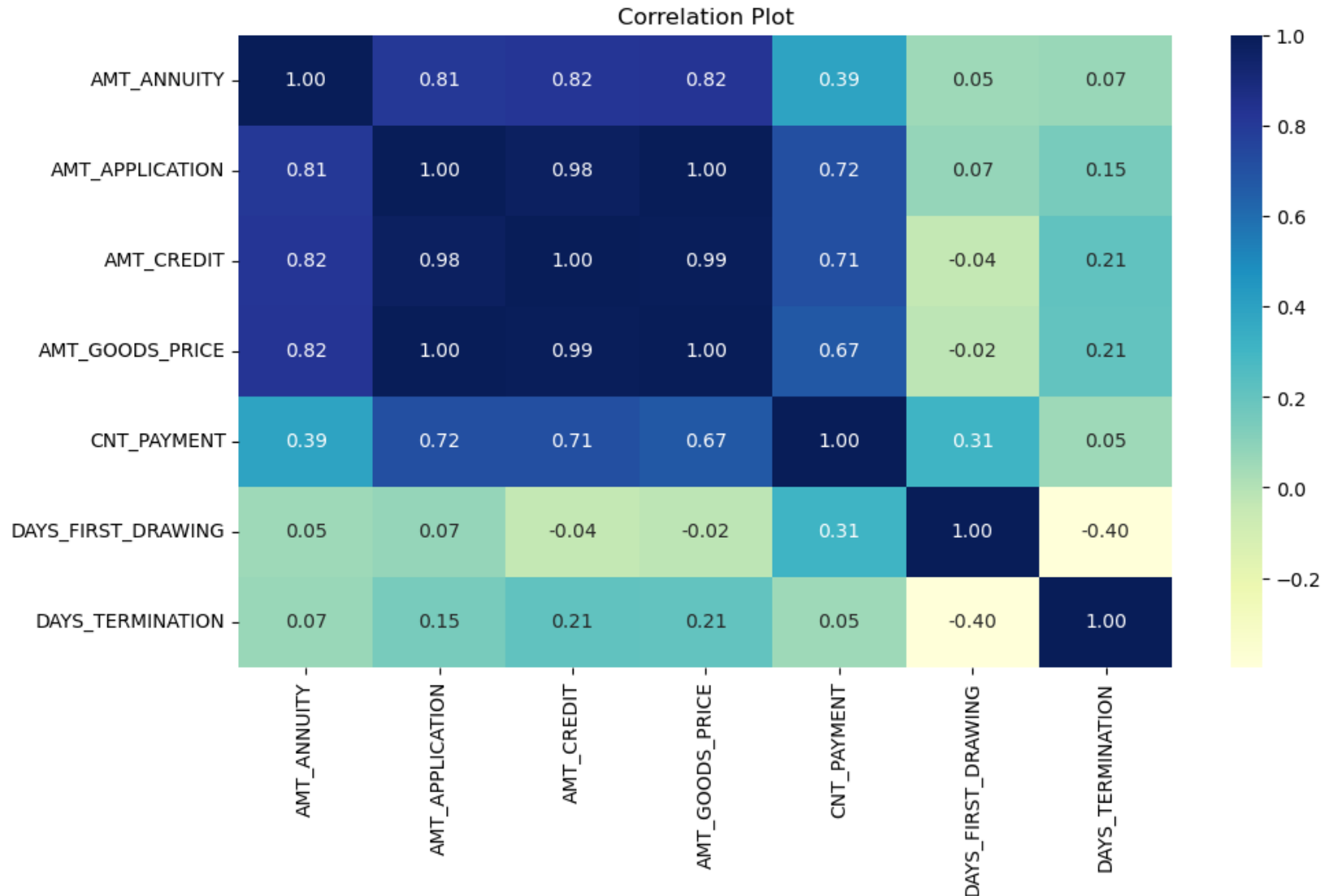
## Bivariate analysis

### Average Application amount by Contract Status



1. The plot shows that most of the clients were refused who had a very high application amount.
2. The clients who opted for average amount of application at a moderate range were accepted and those with lower amount were cancelled.

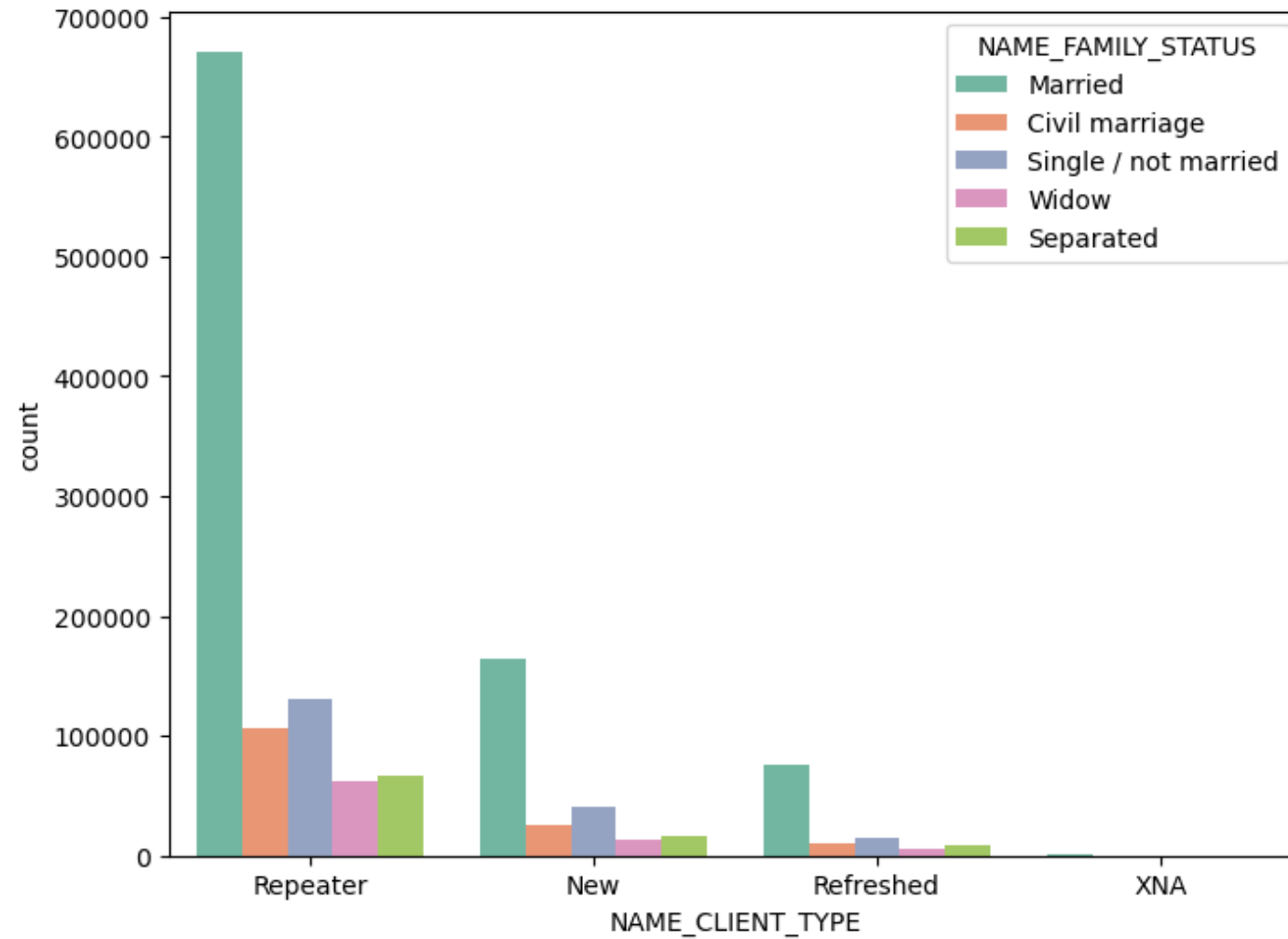
## Correlation plots for numerical variables



- The variable AMT\_ANNUITY,AMT\_APPLICATION,AMT\_CREDIT and AMT\_GOODS\_PRICE shows strong positive correlation with each other.
- While DAYS\_FIRST\_DRAWING and DAYS\_TERMINATION exhibit moderate correlation with other variables.

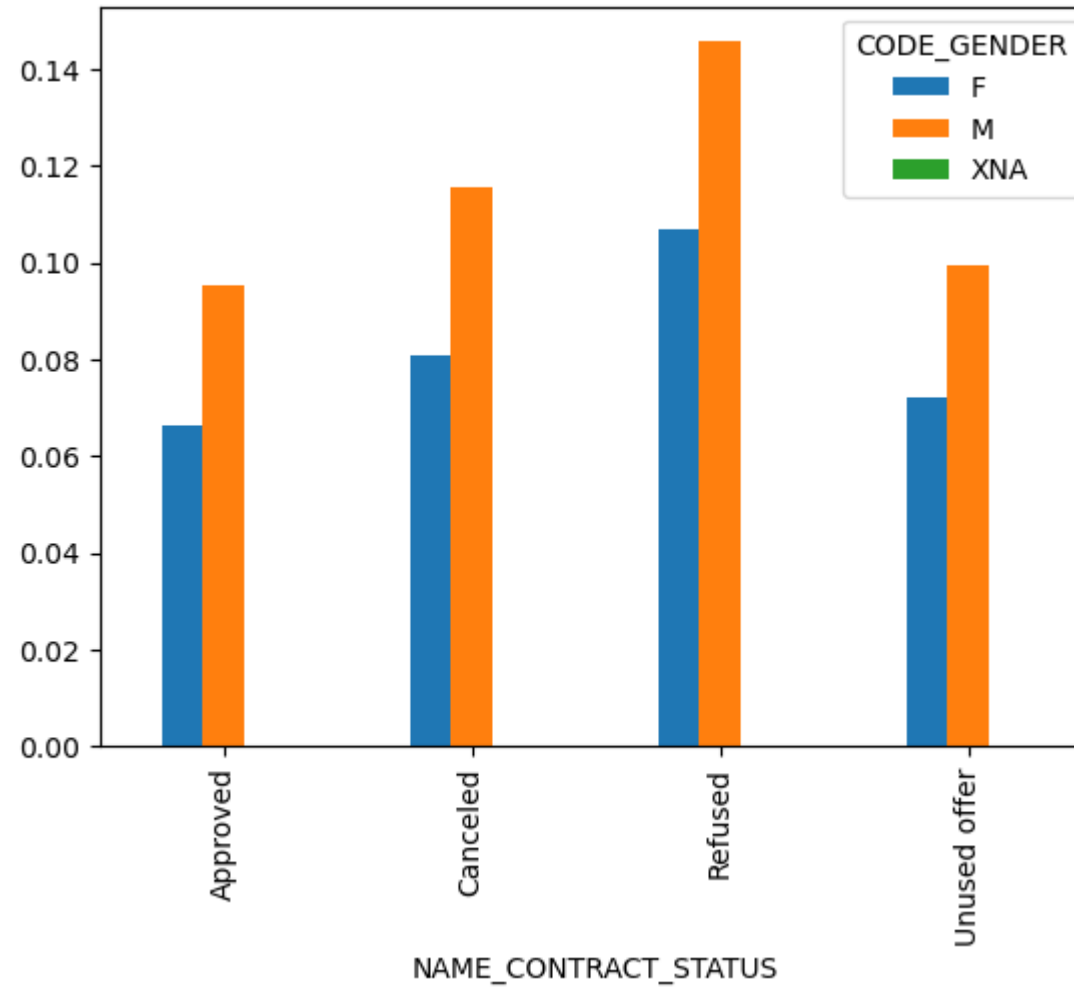
## Merged Data analysis

### Types of Client based on Family status



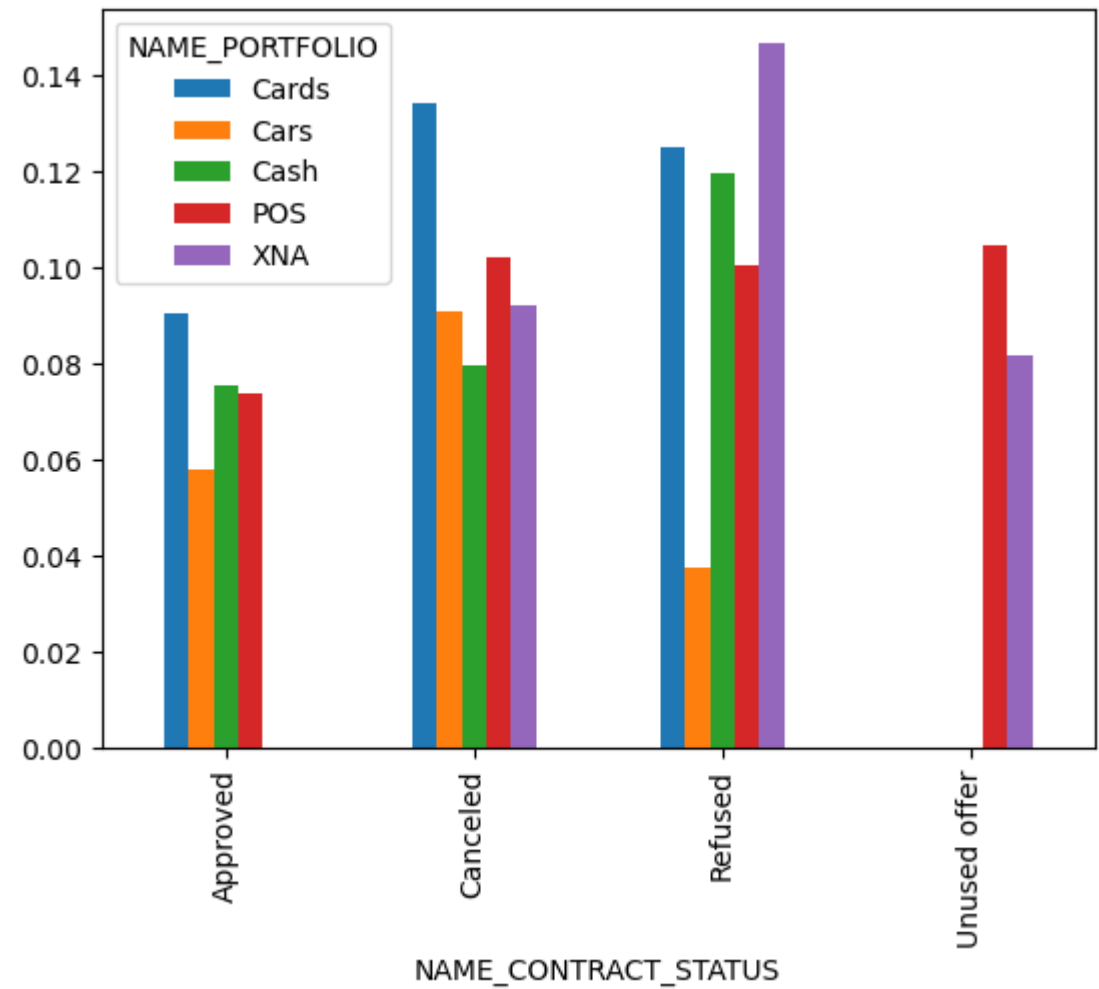
The count plot of the type of clients based on their family status shows that most of the clients approved for credit were repeated clients and married.

## Current loan defaulters based on Previous application status wrt gender



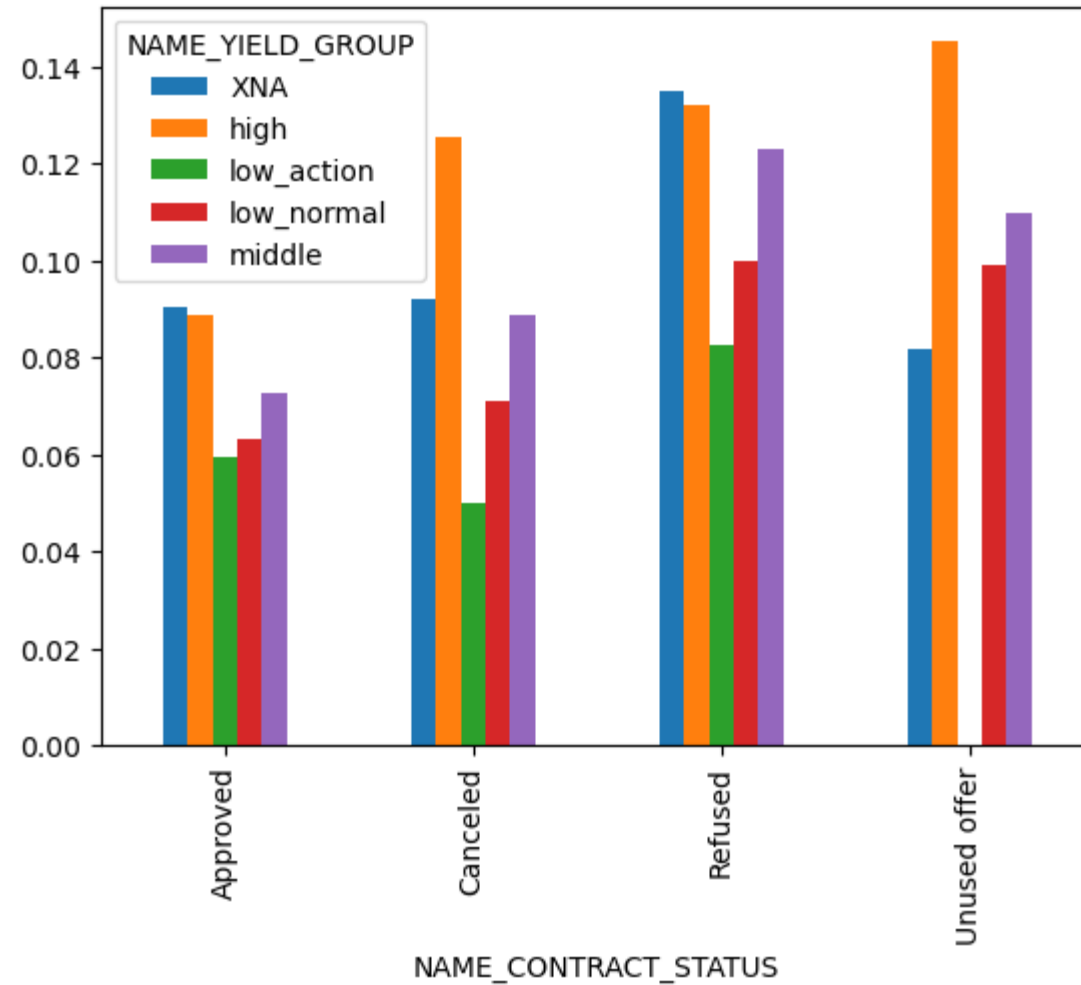
The plot shows that males are higher in proportion than female defaulters and loan to Previously refused clients has high risk of loan default.

Current loan defaulters based on Previous application status wrt Name portfolio



- 1.The loan defaulters are less for those approved for cars from previous application and exhibits high risk for cards portfolio.
- 2.For those previously refused applicants, the default loan status is less for cars portfolio.
- 3.For Previously cancelled applicants, cash exhibits low risk of loan default.

## Current loan defaulters based on Previous application status wrt Interest category



- There is a risk of loan default for those having high grouped interest amount.
- In previously approved application, low-action interest group has lesser loan default risk. The same pattern goes for both cancelled and refused loan applicants.

# Key Insights from the analysis

## Targets based on Categories

- **OCCUPATION\_TYPE** - Labourers
  - **CODE\_GENDER** - Females
  - **AGE\_BY\_DECADE** – 30s age category
  - **NAME\_CONTRACT\_TYPE** – Cash loans
  - **HAS\_CHILDREN** – Does not have children
  - **NAME\_EDUCATION\_TYPE** – Secondary/ Secondary special
  - **NAME\_INCOME\_TYPE** – Working professionals
  - **NAME\_FAMILY\_STATUS** – Married
  - **AMT\_INCOME\_TOTAL** – lower to middle range income
- 
- Most of the clients approved for credit were repeated clients and married.
  - Loan to Previously refused clients has high risk of loan default.
  - The loan defaulters are less for those approved for cars from previous application
  - Low-action interest group has lesser loan default risk

THANK YOU