

# CREDIT EXPLORATORY DATA ANALYSIS CASE STUDY

**Executive Post Graduate Programme in Data Science**  
**By Dharshak Chandra P**  
**DS C52 Batch**

# Business Objective

- This case study aims to identify patterns which indicate if a client has difficulty paying their instalments which may be used for taking actions such as denying the loan, reducing the amount of loan, lending (to risky applicants) at a higher interest rate, etc. This will ensure that the consumers capable of repaying the loan are not rejected. Identification of such applicants using EDA is the aim of this case study.

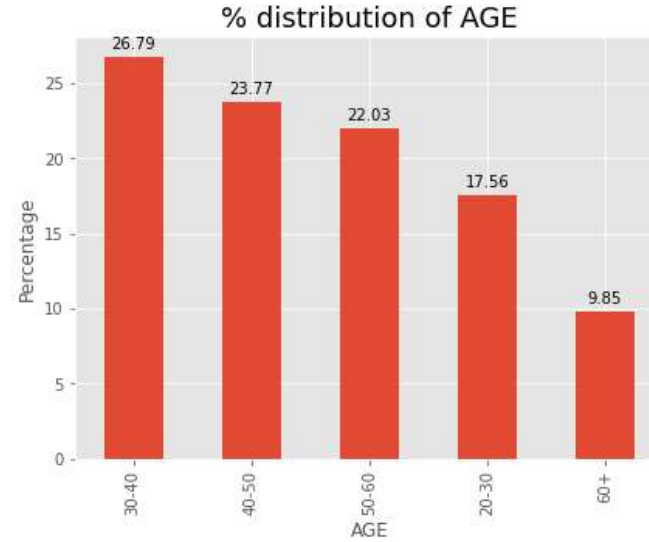
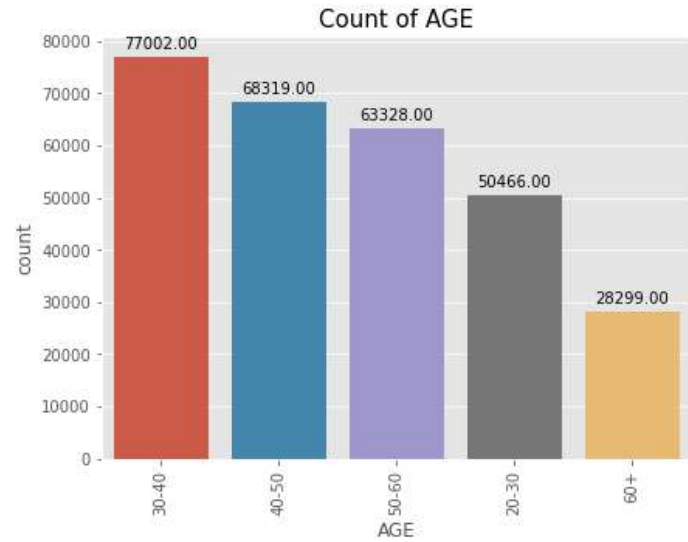
# Function created to do Univariate Analysis on Categorical columns

```
def cat_uni_analysis(x):  
    print('Categorical Univariate Analysis of {}'.format(x))  
  
    plt.figure(figsize=[15,5])  
    plt.subplot(1,2,1)  
    plots = sns.countplot(app0[x],order = app0[x].value_counts().index)  
    plt.xticks(rotation = 90)  
    plt.title('Count of {}'.format(x), fontdict={'fontsize':15})  
    for bar in plots.patches:  
        plots.annotate(format(bar.get_height(), '.2f'),  
                        (bar.get_x() + bar.get_width() / 2,  
                         bar.get_height()), ha='center', va='center',  
                        size=10, xytext=(0, 8),  
                        textcoords='offset points')  
  
    plt.subplot(1,2,2)  
    plt.title('% distribution of {}'.format(x), fontdict={'fontsize':18})  
    plt.xlabel(x)  
    plt.ylabel('Percentage')  
    plots = (app0[x].value_counts()*100/len(app0[x])).plot.bar()  
    for bar in plots.patches:  
        plots.annotate(format(bar.get_height(), '.2f'),  
                        (bar.get_x() + bar.get_width() / 2,  
                         bar.get_height()), ha='center', va='center',  
                        size=10, xytext=(0, 8),  
                        textcoords='offset points')  
  
    fig.tight_layout()  
    plt.show()
```

# Function created to do Univariate Analysis on Numerical columns

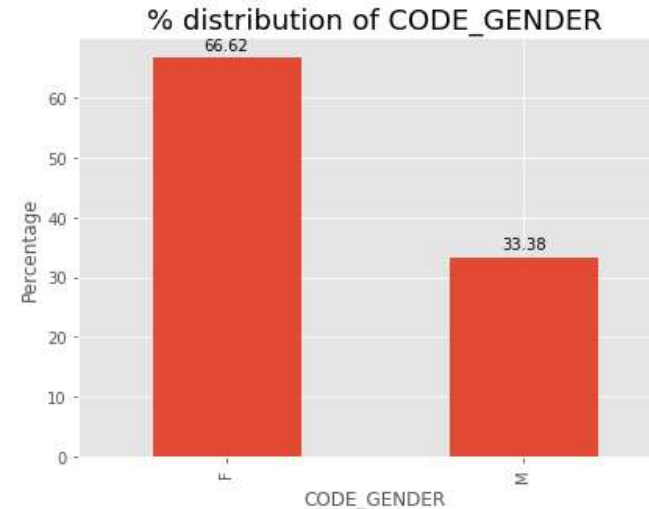
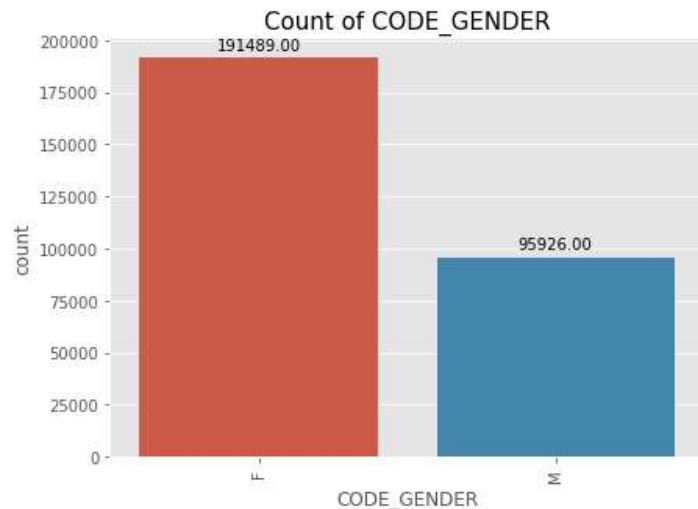
```
def num_uni_analysis(x):  
    print('Numerical Univariate Analysis of {}'.format(x))  
    print(app0[x].describe())  
    plt.figure(figsize=[15,7])  
    plt.subplot(1,2,1)  
    sns.boxplot(app0[x])  
    plt.title('Box Plot of {}'.format(x),fontdict={'fontsize':18})  
  
    plt.subplot(1,2,2)  
    sns.distplot(app0[x],hist=False)  
    plt.title('Distplot of {}'.format(x),fontdict={'fontsize':18})  
  
    fig.tight_layout()  
    plt.show()
```

Categorical Univariate Analysis of AGE



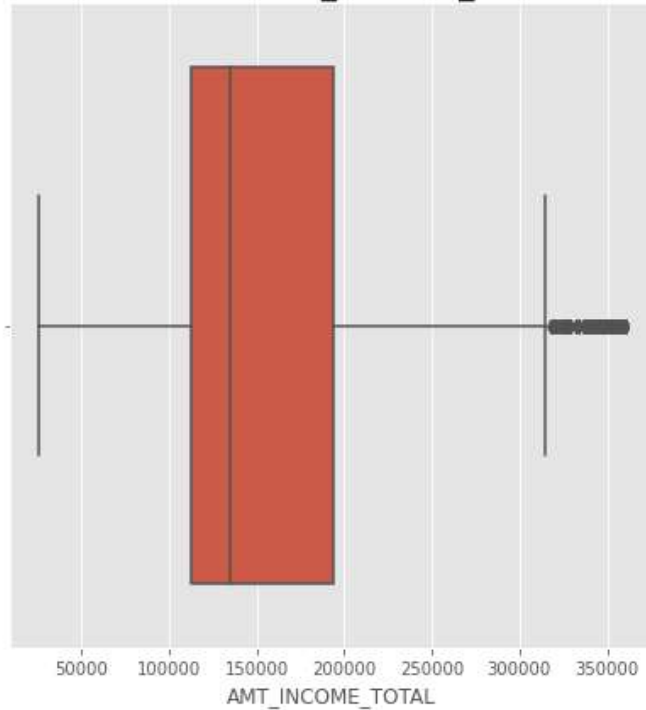
Examples of  
Visualization  
after applying  
the Function to  
Categorical  
columns

Categorical Univariate Analysis of CODE\_GENDER

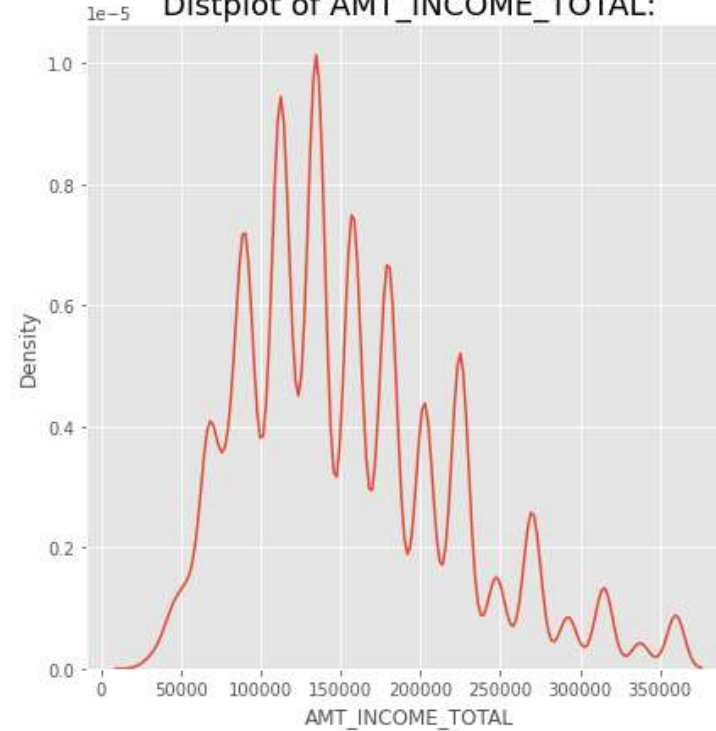


```
Numerical Univariate Analysis of AMT_INCOME_TOTAL:  
count    287415.000000  
mean     154470.901144  
std      66518.029770  
min      25650.000000  
25%     112500.000000  
50%     135000.000000  
75%     193500.000000  
max      360000.000000  
Name: AMT_INCOME_TOTAL, dtype: float64
```

Box Plot of AMT\_INCOME\_TOTAL:



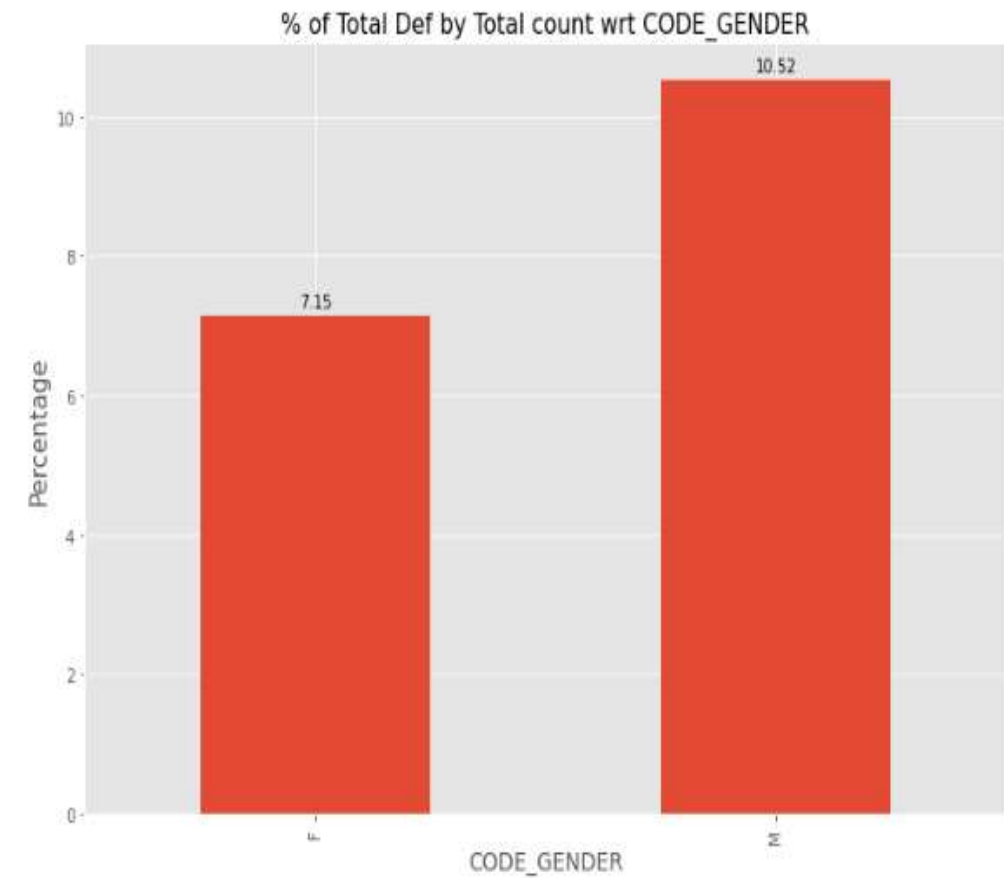
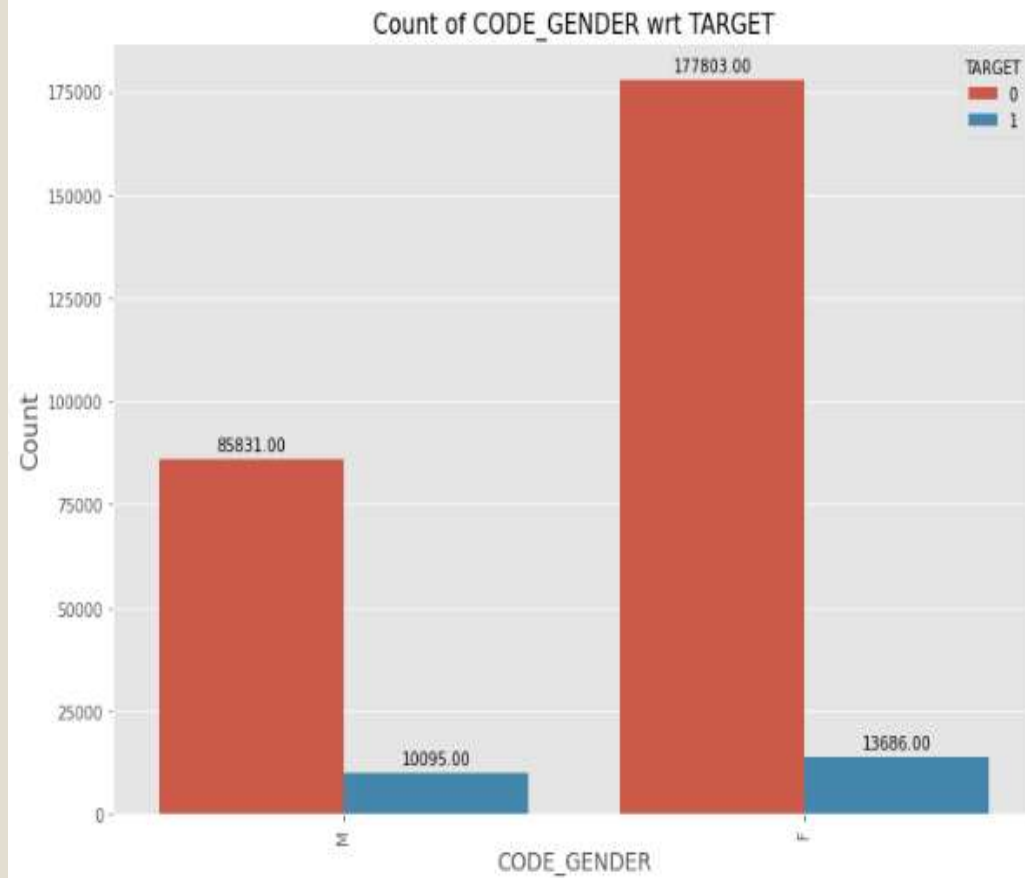
Distplot of AMT\_INCOME\_TOTAL:



**Examples of  
Visualization  
after applying  
the Function to  
Numerical  
columns**

# **Segmented Univariate Analysis and Bivariate Analysis**

## Segmented Categorical Univariate Analysis of CODE\_GENDER

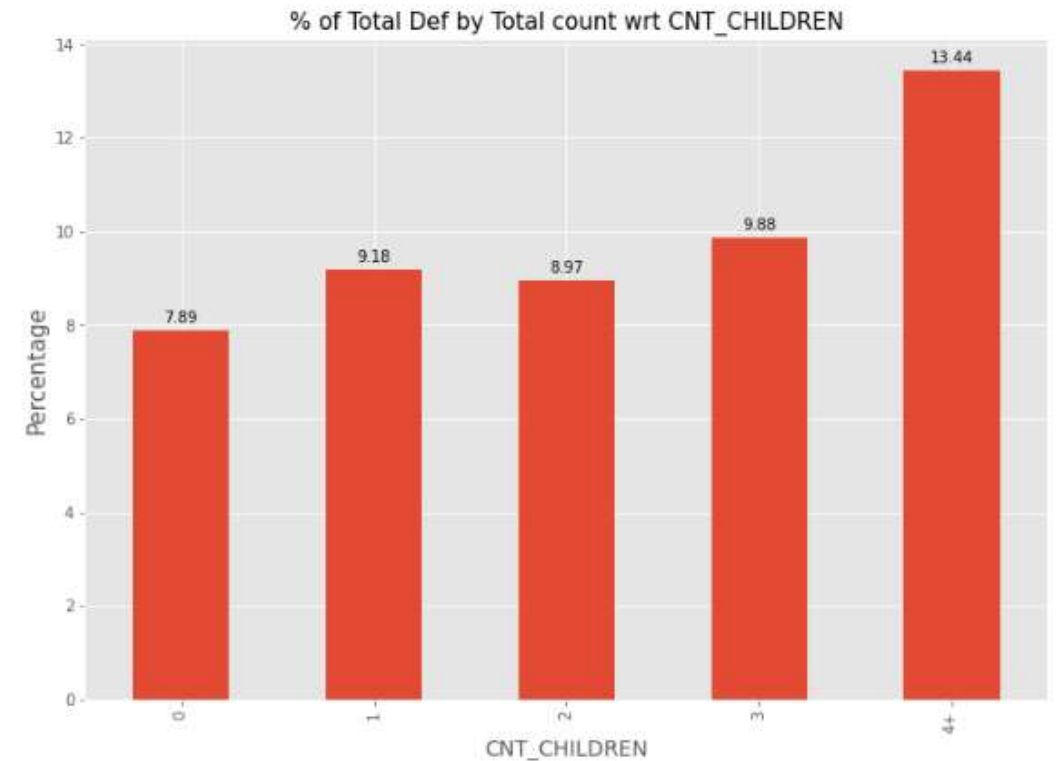
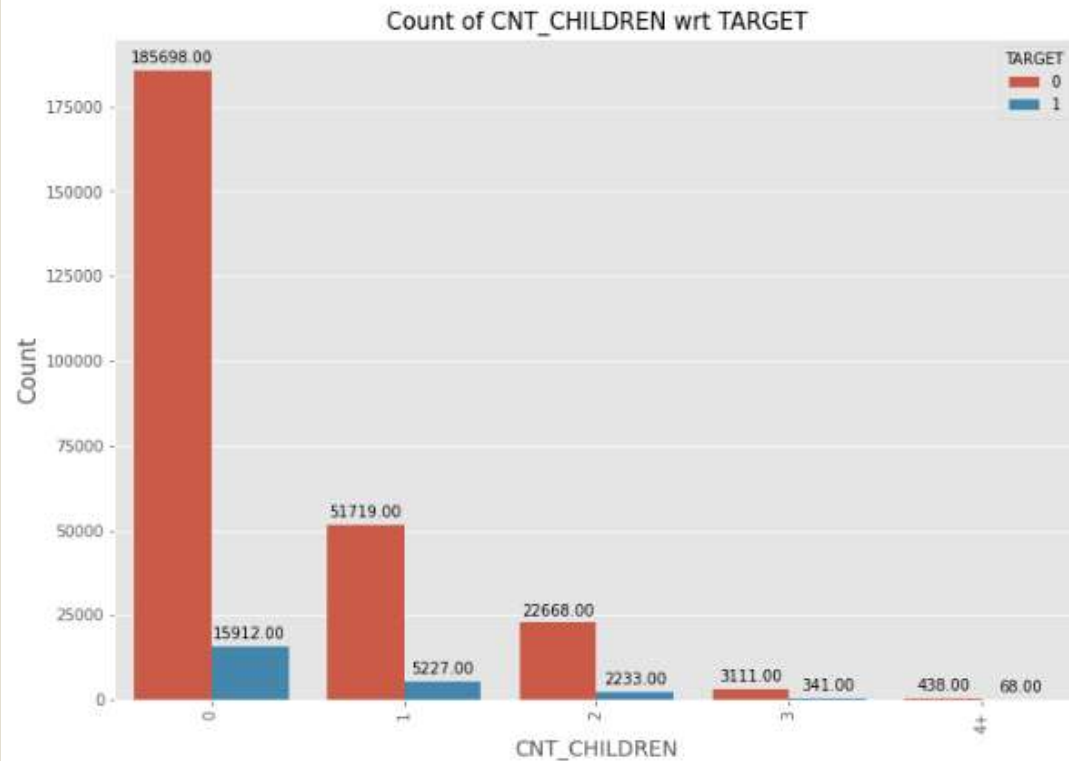


Inference:

- Male are high Defaulters compared to Female



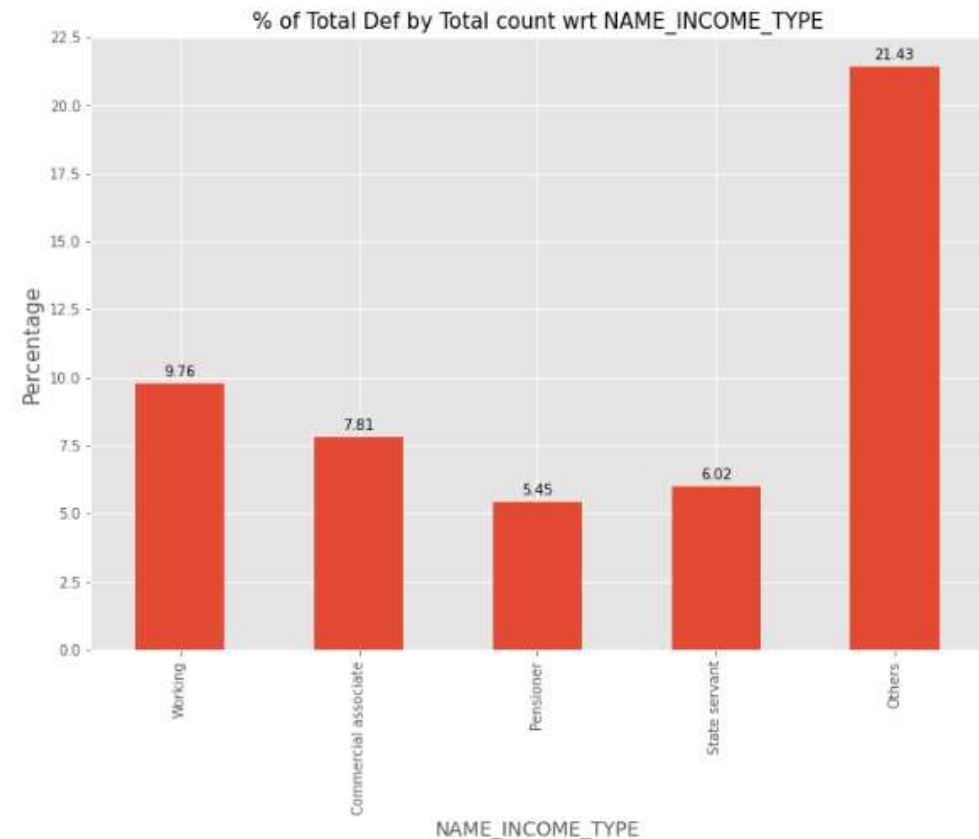
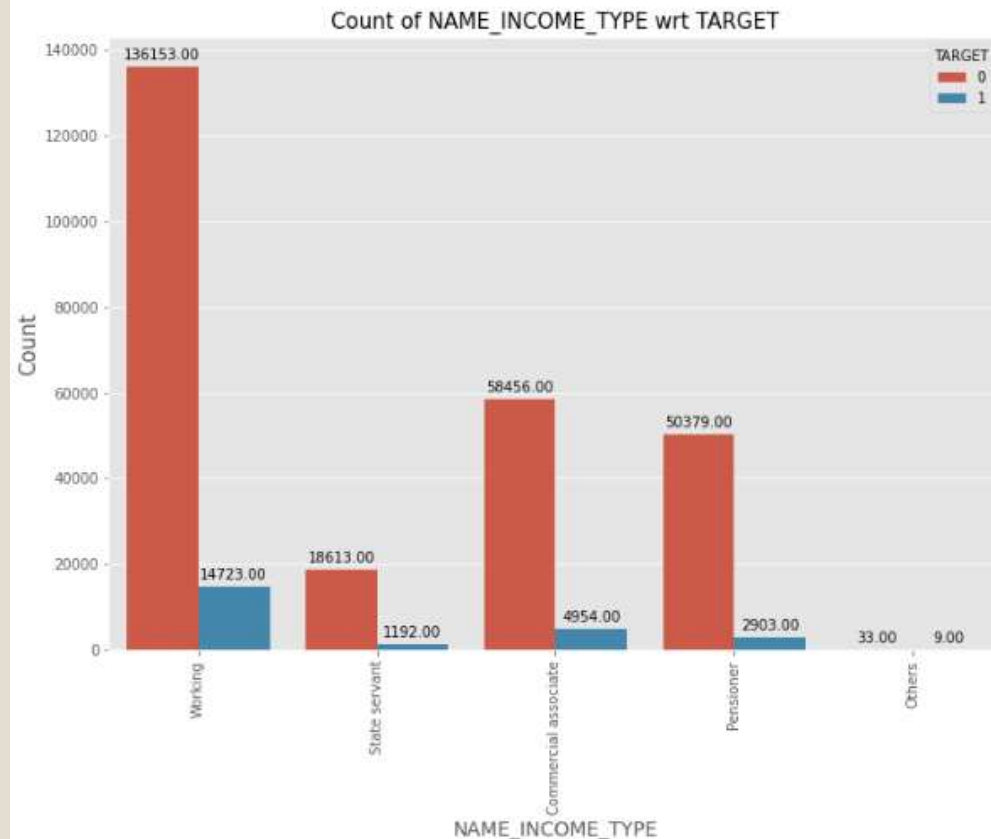
## Segmented Categorical Univariate Analysis of CNT\_CHILDREN



Inference:

- As the number of children increase, the rate of defaulters also increases.

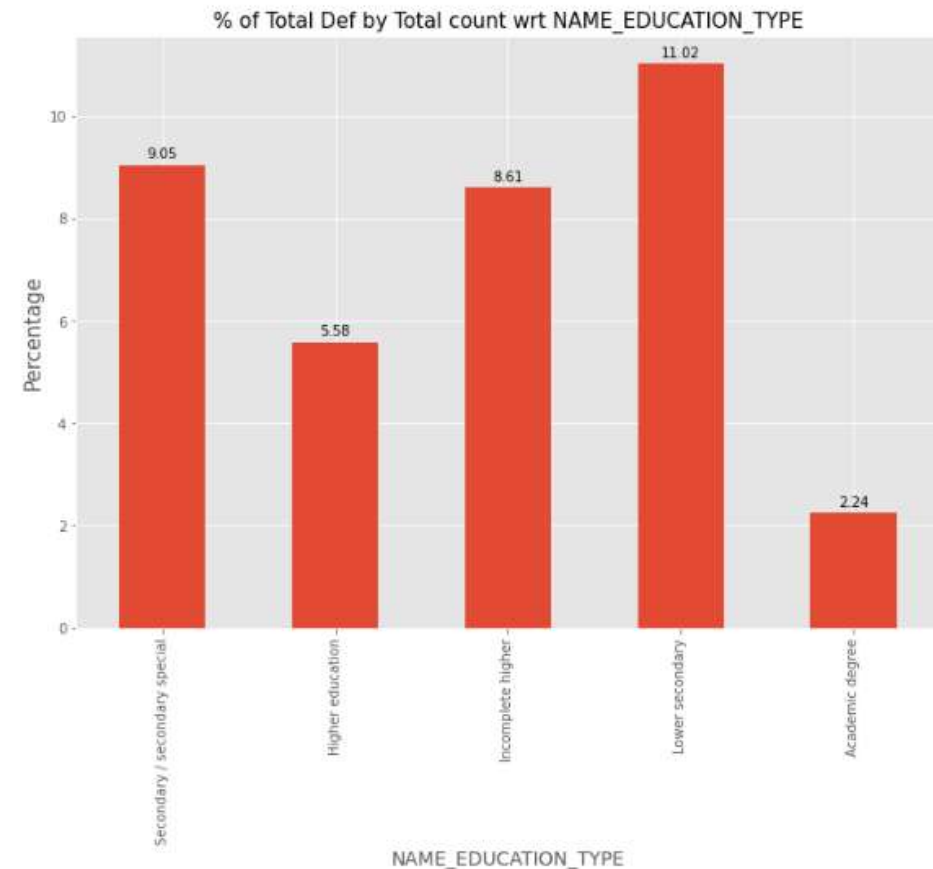
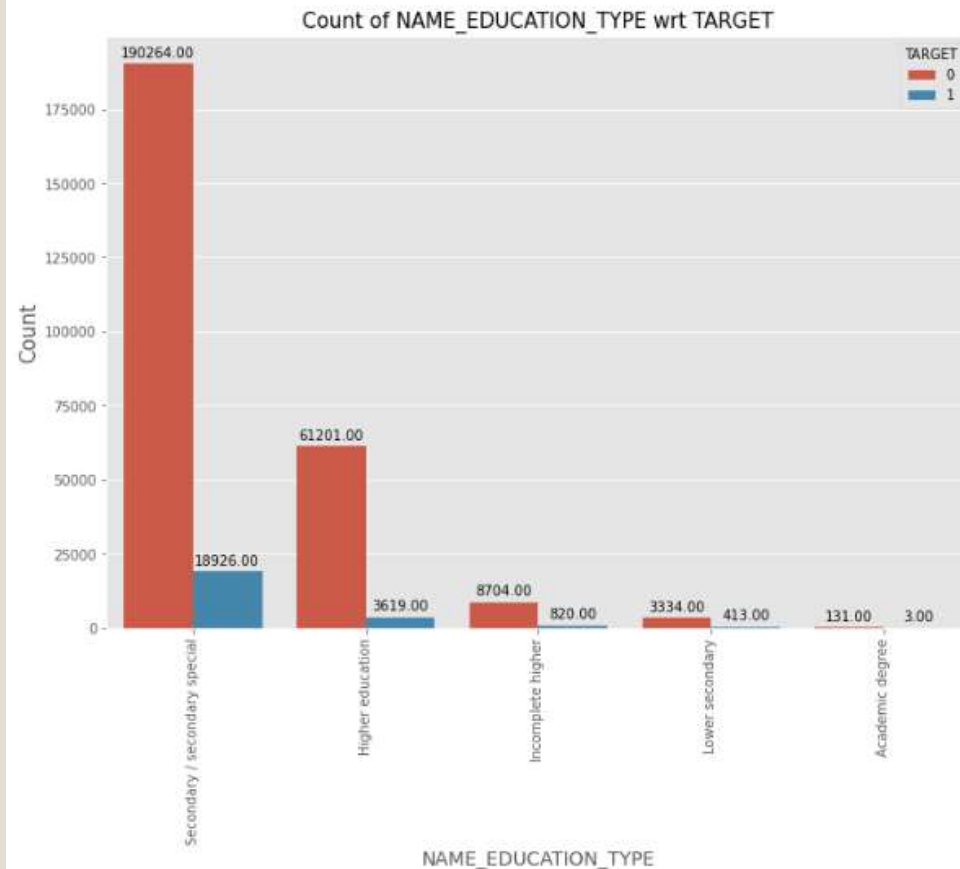
## Segmented Categorical Univariate Analysis of NAME\_INCOME\_TYPE



### Inference:

- In the NAME\_INCOME\_TYPE, Others are the high defaulter. This others contain ['Unemployed', 'Student', 'Businessman', 'Maternity leave']. This is followed by Working and Commercial Associate as high defaulters. Pensioners and State servant are the least defaulters.

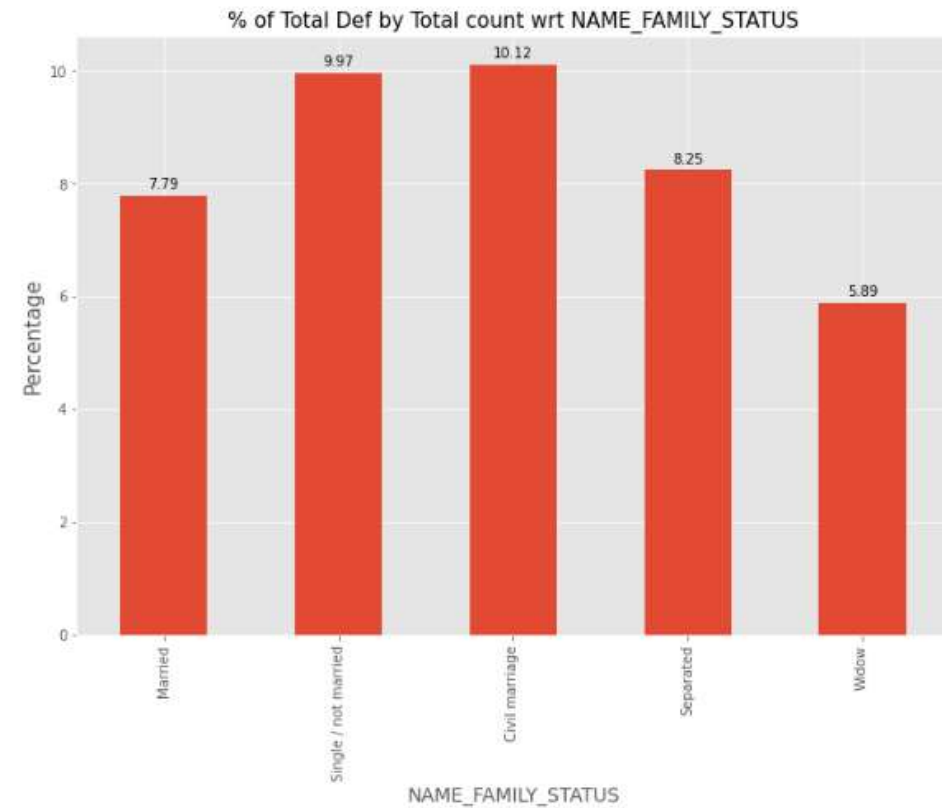
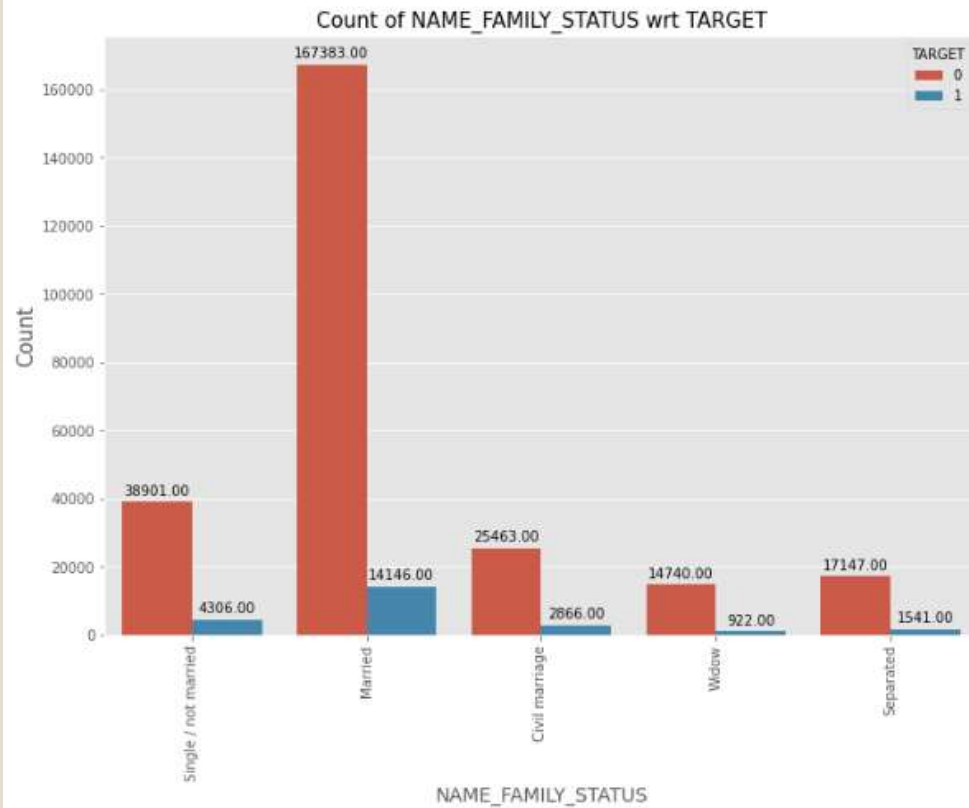
## Segmented Categorical Univariate Analysis of NAME\_EDUCATION\_TYPE



### Inference:

- Under NAME\_EDUCATION\_TYPE we see that Lower Secondary and Secondary special are the group of people who are high defaulters. People with Academic degree amount to the least amount of defaulters.

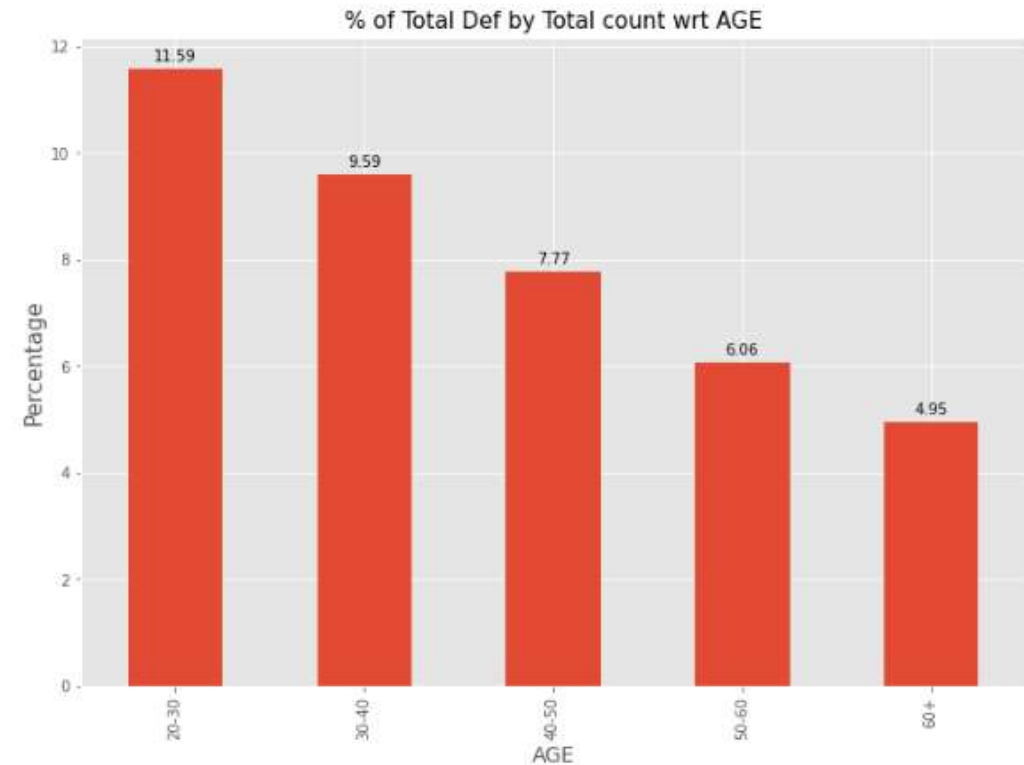
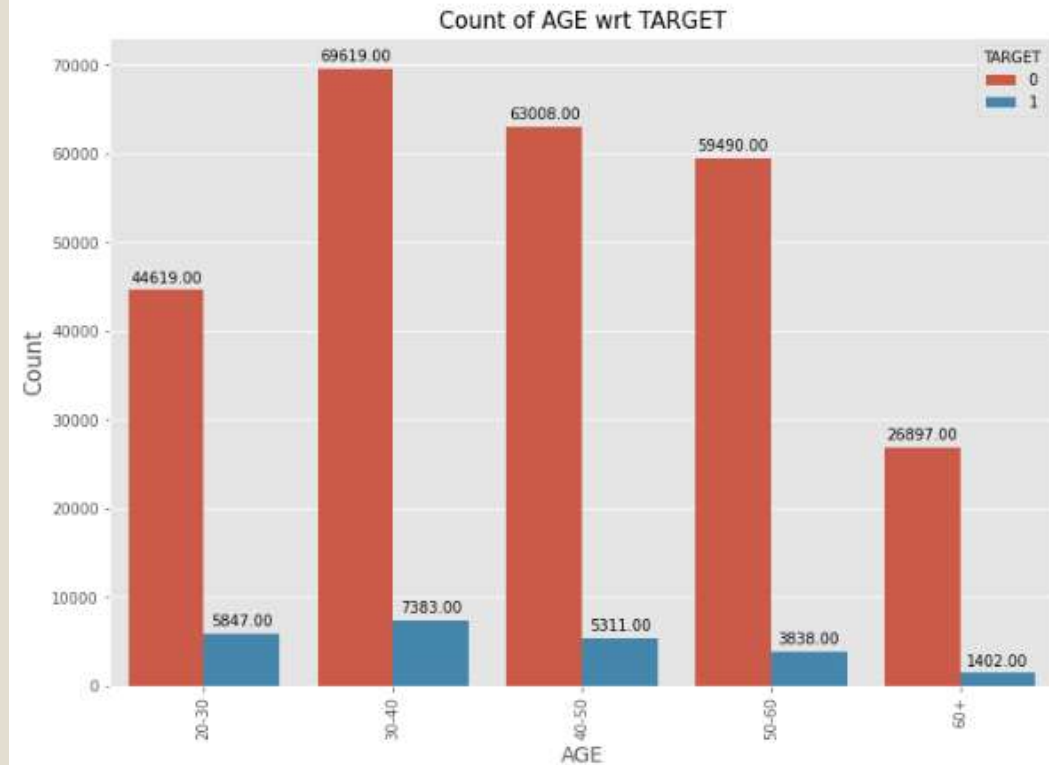
### Segmented Categorical Univariate Analysis of NAME\_FAMILY\_STATUS



#### Inference:

- For NAME\_FAMILY\_STATUS we see that Civil Marriage and Single/ not married have higher defaulting rates. Widows are the least amount of defaulters.

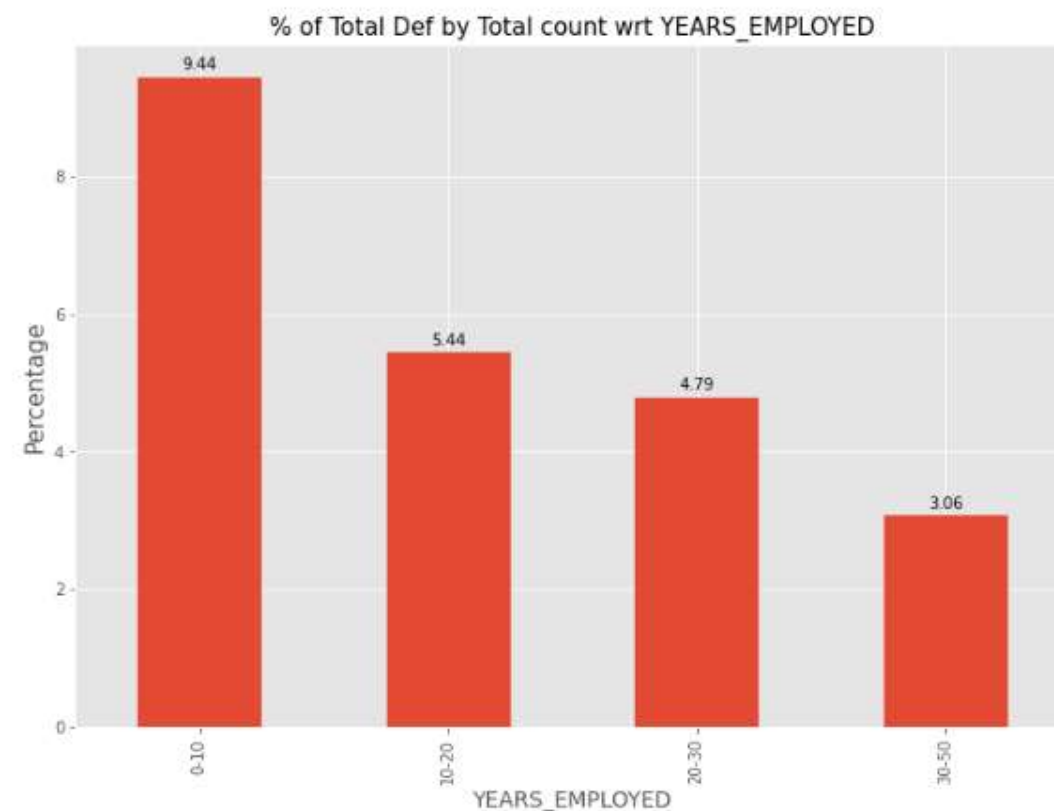
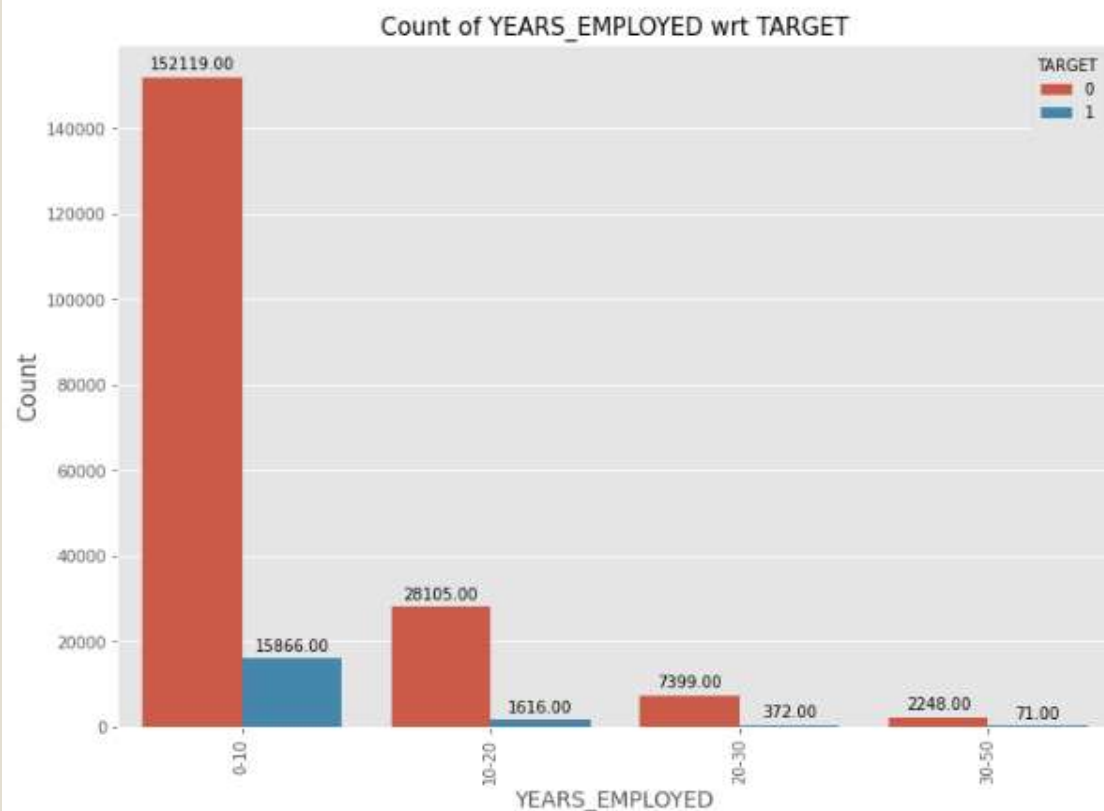
## Segmented Categorical Univariate Analysis of AGE



### Inference:

- For AGE column there is a clear trend which shows that young people have difficulties to pay back the loans compared to the old people. Therefore the rate of defaulting decreases as the age increases.

## Segmented Categorical Univariate Analysis of YEARS\_EMPLOYED



### Inference:

- For YEARS\_EMPLOYED column there is a clear trend which shows that people with less experience have difficulties to pay back the loans compared to the highly experienced people. Therefore the rate of defaulting decreases as the amount of experience increases. This could be possible as people who are more experienced could be paid higher compared to fresher.

# List of all the Inference made on these Categorical columns

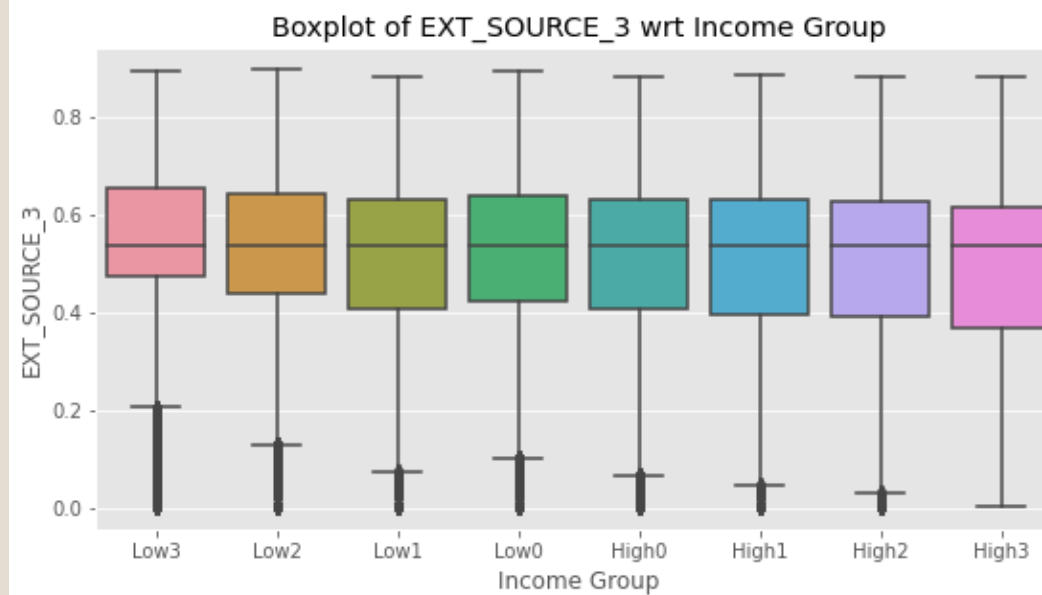
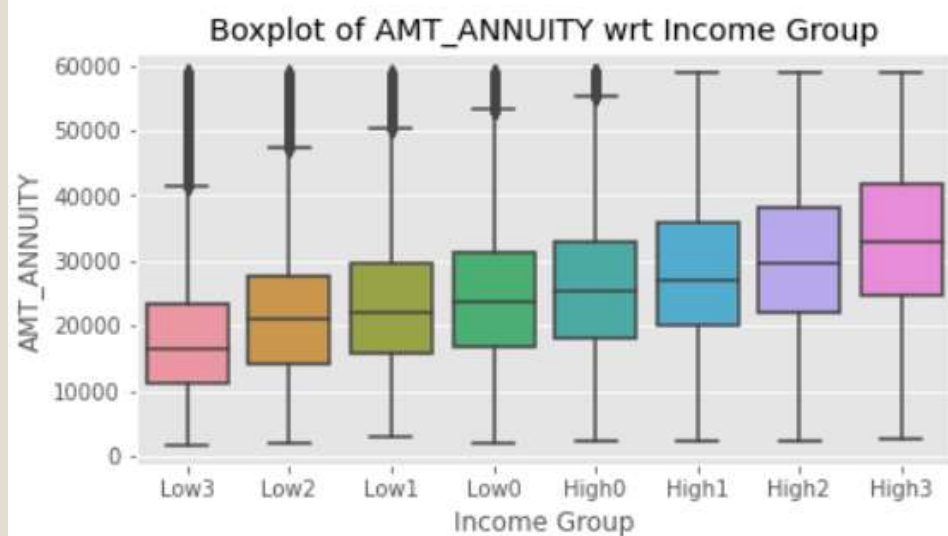
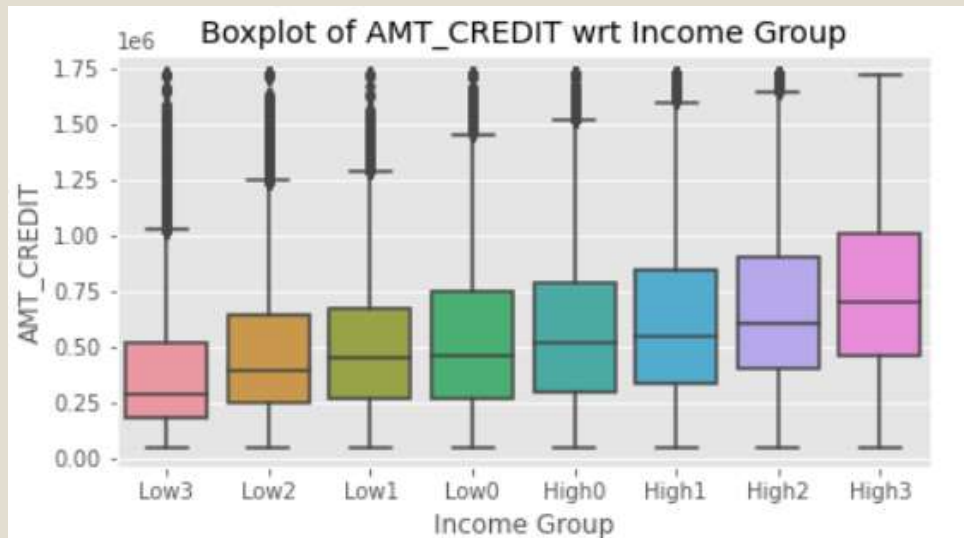
- Male are high Defaulters compared to Female
- As the number of children increase, the percentage of defaulters also increases.
- In the NAME\_INCOME\_TYPE, Others are the high defaulter. This others contain ['Unemployed', 'Student', 'Businessman', 'Maternity leave']. This is followed by Working and Commercial Associate as high defaulters. Pensioners and State servant are the least defaulters.
- Under NAME\_EDUCATION\_TYPE we see that Lower Secondary and Secondary special are the group of people who are high defaulters. People with Academic degree amount to the least amount of defaulters.
- For NAME\_FAMILY\_STATUS we see that Civil Marriage and Single/ not married have higher defaulting rates. Widows are the least amount of defaulters.
- For NAME\_HOUSING\_TYPE we see that people living in rented apartments or people living with their parents have higher defaulting rates. Whereas people living in Office Apartments and House have lesser defaulting rates.
- For REGION\_POPULATION\_RELATIVE there is a clear trend which shows that as the population of the place increase, the defaulting percentage decrease. This could be because since there are more people there is a possibility to have more jobs available and more opportunity to earn money.
- For AGE column there is a clear trend which shows that young people have difficulties to pay back the loans compared to the old people. Therefore the rate of defaulting decreases as the age increases.

# List of all the Inference made on these Categorical columns

- For YEARS\_EMPLOYED column there is a clear trend which shows that people with less experience have difficulties to pay back the loans compared to the highly experienced people. Therefore the rate of defaulting decreases as the amount of experience increases. This could be possible as people who are more experienced could be paid higher compared to fresher.
- YEARS\_REGISTRATION also follows the trend of AGE and YEARS\_EMPLOYED.
- For OWN\_CAR\_AGE we see that people who have very old cars, that is over 12 years have higher defaulting rates compared to people who have cars of age 0 to 9. We also see that people who don't own cars also have higher defaulting rates.
- For OCCUPATION\_TYPE we see that people who are low skilled like- Low skilled labours, Drivers, Waitress have higher defaulting percentage. Whereas people who are in jobs that require high educations like- Accountants, HR staff, IT staff, Managers have the least defaulting percentage.
- From Total\_documents we see that people who have submitted 3 documents to have high defaulting percentage.(Can't really tell the reason for it).
- From Income Group we see that people having the least amount of income are finding it hard to pay back compared to people who have highest income. It is expected, so need to be careful while giving loans to the people who have low income.
- AMT\_CREDIT\_RANGE follows trend of a normal distribution curve where the extremes that is 0-200000 and 800000 and above have the least defaulting percentage and the middle 400000-600000 has highest defaulting rates.



## Plotting Boxplot for Income Groups against other Numerical columns



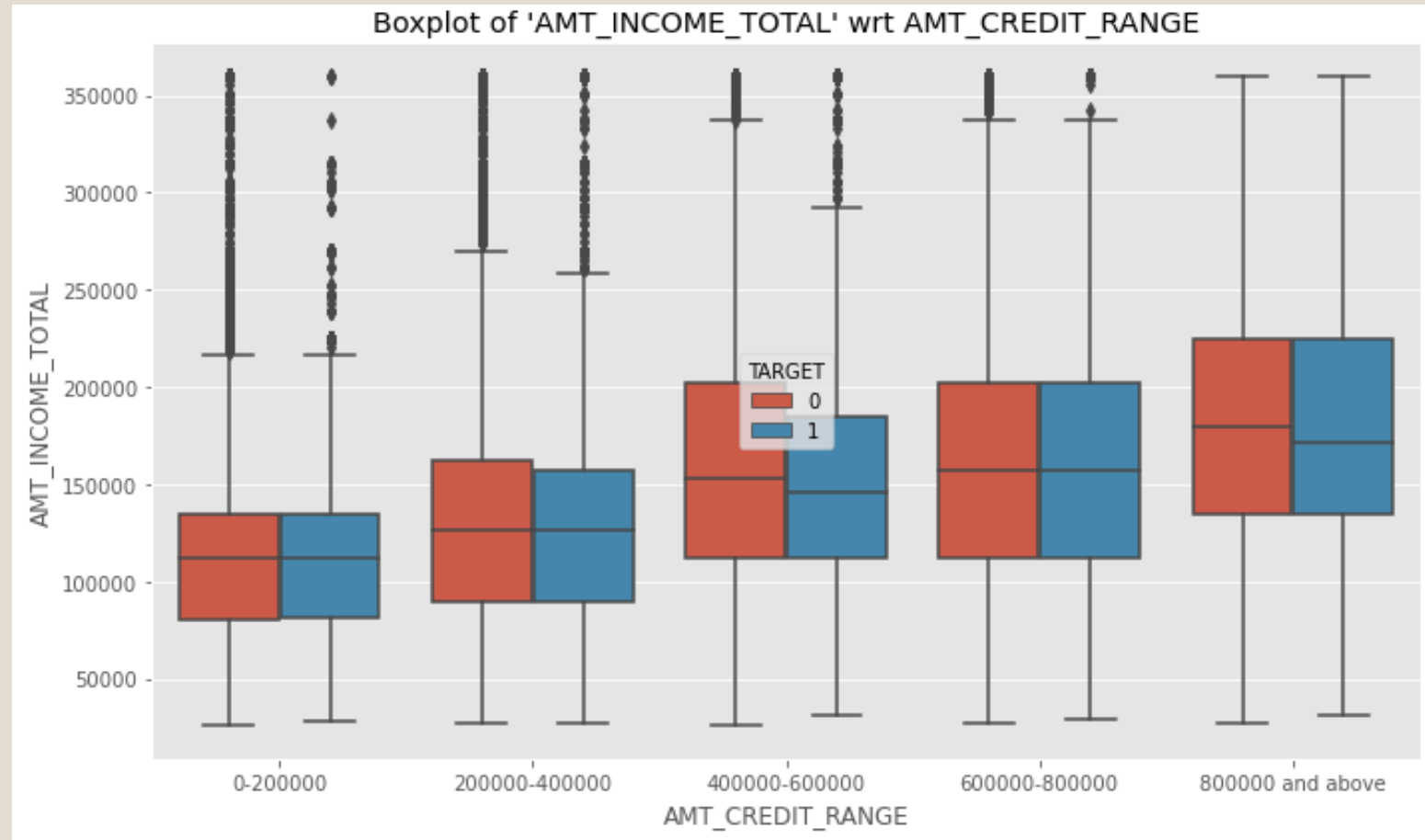
# Inference

- As people start earning more, people take higher amount of credit.
- As people start earning more, people pay higher amount of annuity. This is because Credit and Annuity have high correlation.
- As people start earning more, people take credit for higher amount of goods price. This is because Credit and goods price have very high correlation as expected.
- As Income increases, EXT\_SOURCE\_2 also increases slightly.
- As Income increases, EXT\_SOURCE\_3 decreases slightly. That is, they are inversely proportional.

## Plotting a box plot for 'AMT\_INCOME\_TOTAL' wrt 'AMT\_CREDIT\_RANGE' to understand why defaulters are struggling to pay back the loan.

### Inference:

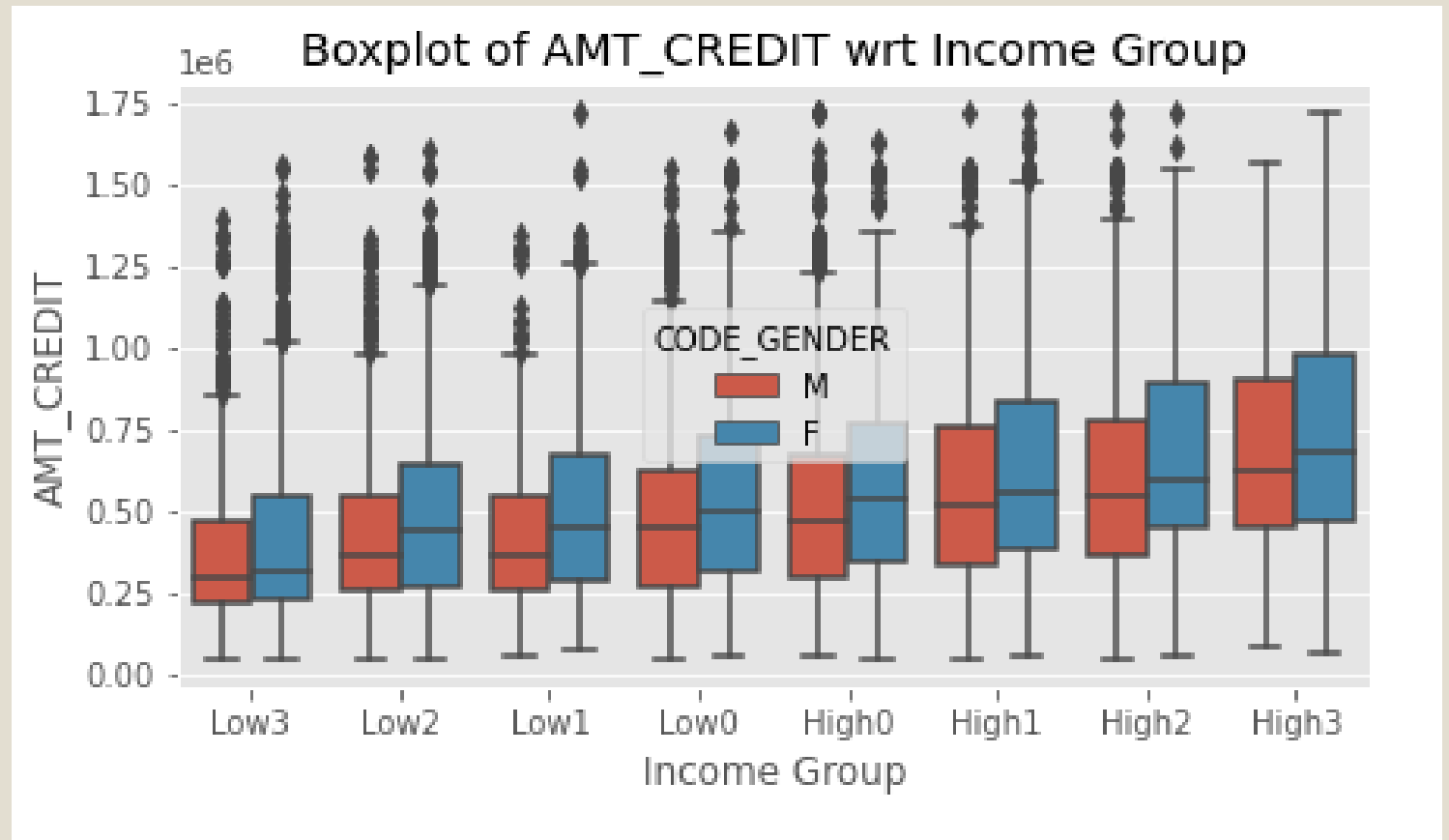
- From the given plot we can clearly understand that from the high defaulting range that is 200000-400000 and 400000-600000, the defaulters have less amount of income compared to non-defaulters and they have taken loans which are slightly above their range to pay it back on time. Therefore they are defaulting. To stop this we can increase the interest rates for people having lesser income but wanting higher credit. This way we can make sure that people who really believe they can pay on time will take the loan at those rates and also bank can minimize the loss in case of defaulting.



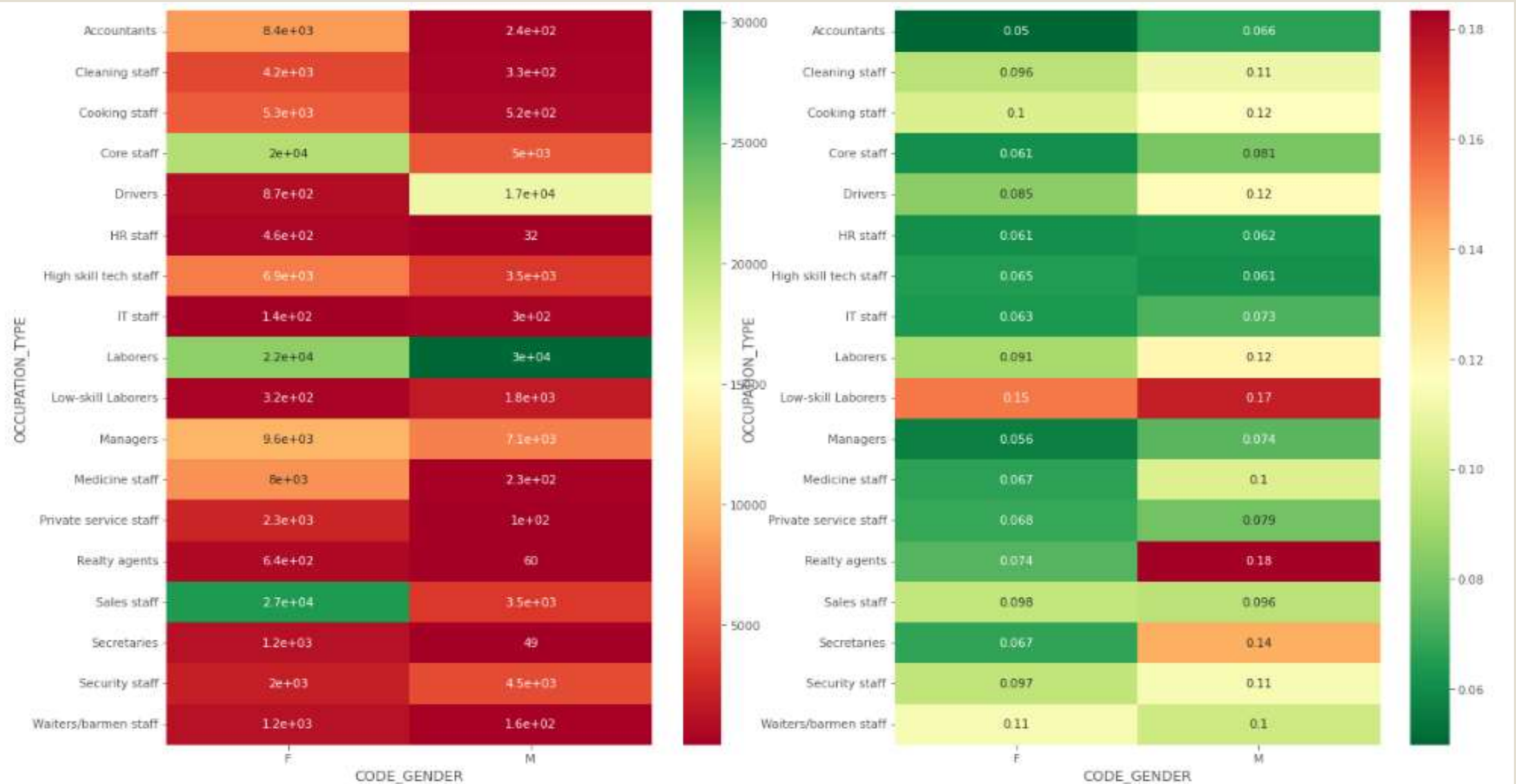
## Plotting a box plot for 'AMT\_CREDIT' wrt 'Income Group' wrt 'CODE\_GENDER'

### Inference:

- From this plots we see that Men in spite of taking lower credit compared to women for same income range are struggling to pay back the loan on time.

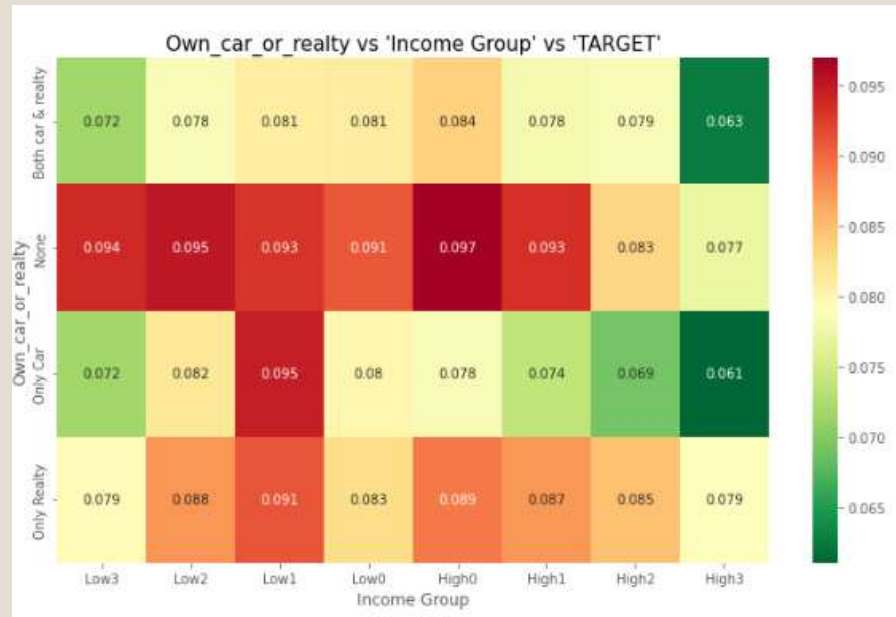


## Plotting Heatmap to understand why Men are higher Defaulters than Women

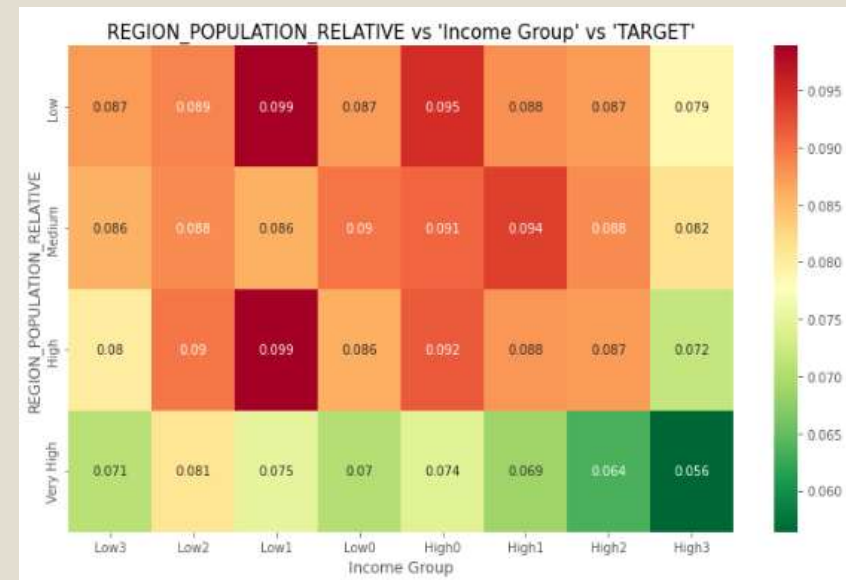
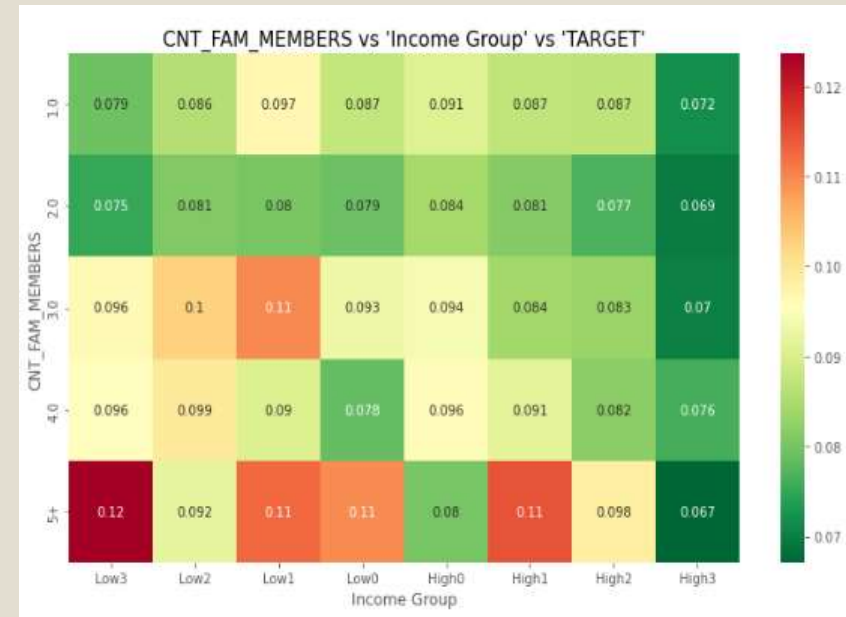


# Inference:

- As we can see here, there are more men in the job roles like "Labourers, Low-skill Labourers and Driver" categories compared to women. These categories make up the highest default making categories. Therefore, men have more default rates than women.
- We also see that the count of women in high paying jobs like "Accountants, HR Staff, Managers, Medicine Staff" are higher than Men in these categories.
- Most loans were provided to Men Laborers who are high defaulters. Thus, increasing the default rate for Men.
- Banks need to be careful while providing loans to these high defaulting categories of people.



**Few Examples of  
Heatmap plotted to  
find the relation  
between  
Categorical  
columns vs  
'Income Group' vs  
'TARGET'**





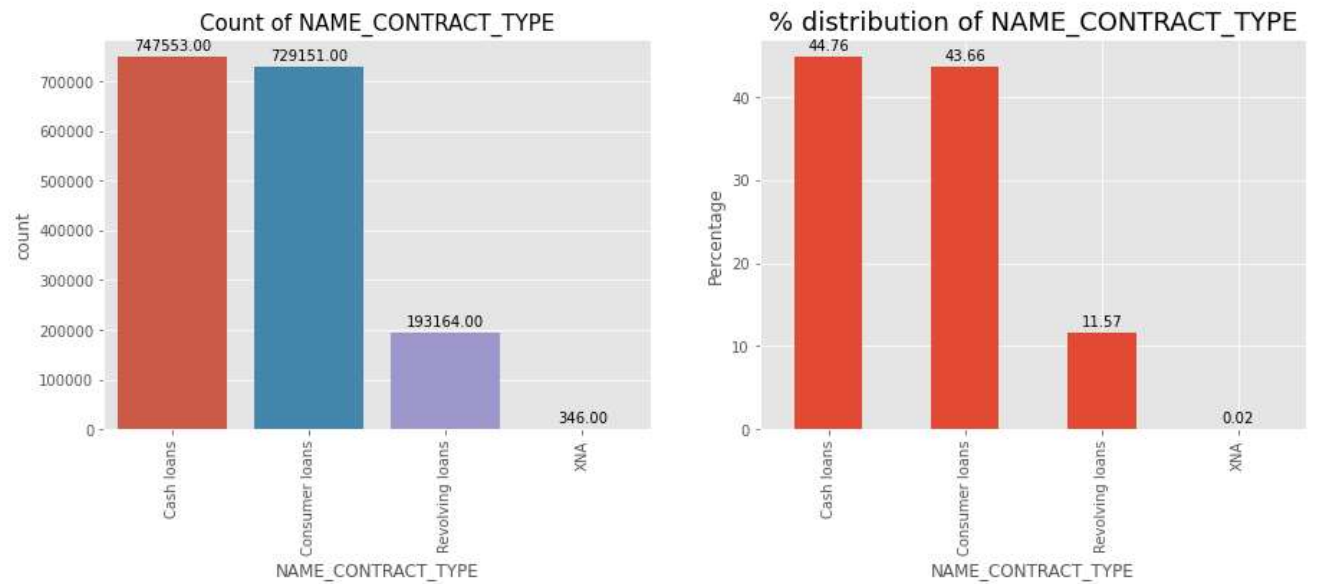
## Inference made after Plotting Heatmap to find the relation between Categorical columns vs 'Income Group' vs 'TARGET'

- From heatmap of CODE\_GENDER, we observe Men are high defaulters in spite of having income up to the range High1
- From heatmap of CNT\_CHILDREN, we observe that as the number of children increase, the rate of defaulting also increases.
- From heatmap of NAME\_INCOME\_TYPE, we learn that category of Others have less income and are high defaulters.
- From heatmap of NAME\_EDUCATION\_TYPE, we observe that People with lower education in spite of earning well tend to default.
- From heatmap of NAME\_FAMILY\_STATUS, we can see that Widow are the safest category to provide loans as the rate of default is very less even for people having low income. Also people belonging to the category of Civil marriage and Single are high defaulters.
- From heatmap of NAME\_HOUSING\_TYPE, we see that people who live with their parents and living in Rented apartment have high default rates.
- From heatmap of REGION\_POPULATION\_RELATIVE, we observe that default rate decreases with increase in population of city. Therefore they are inversely proportional.
- From heatmap of AGE, we can see that people ranging from 20-40 have higher rate of default. Whereas people who are above 50 are the best range of people to provide loans as their rate of default is very low.
- From heatmap of YEARS\_EMPLOYED, similar to AGE people with less experience tend to default more than people with higher experience.
- From heatmap of OWN\_CAR\_AGE, we observe that people who have very old cars, that is above 12 years of age and have income in the range of low2 to high2 default a lot compared to others.
- From heatmap of OCCUPATION\_TYPE, we see that Low-Skill labourers are high defaulters. This is true to people who are earning very well as well. People with high paying jobs that have high education requirements tend to default very less.
- From heatmap of CNT\_FAM\_MEMBERS, we see that people with family of more than or equal to 3 and salary less than Low0 tend to default a lot. So, it is better to provide loans to people less than 3 members in their family and for people who have 3 or more members but with Income greater than High1.
- From heatmap of Own\_car\_or\_realty, we see that people with No car or Realty and Only realty tend to have higher defaulting rates.

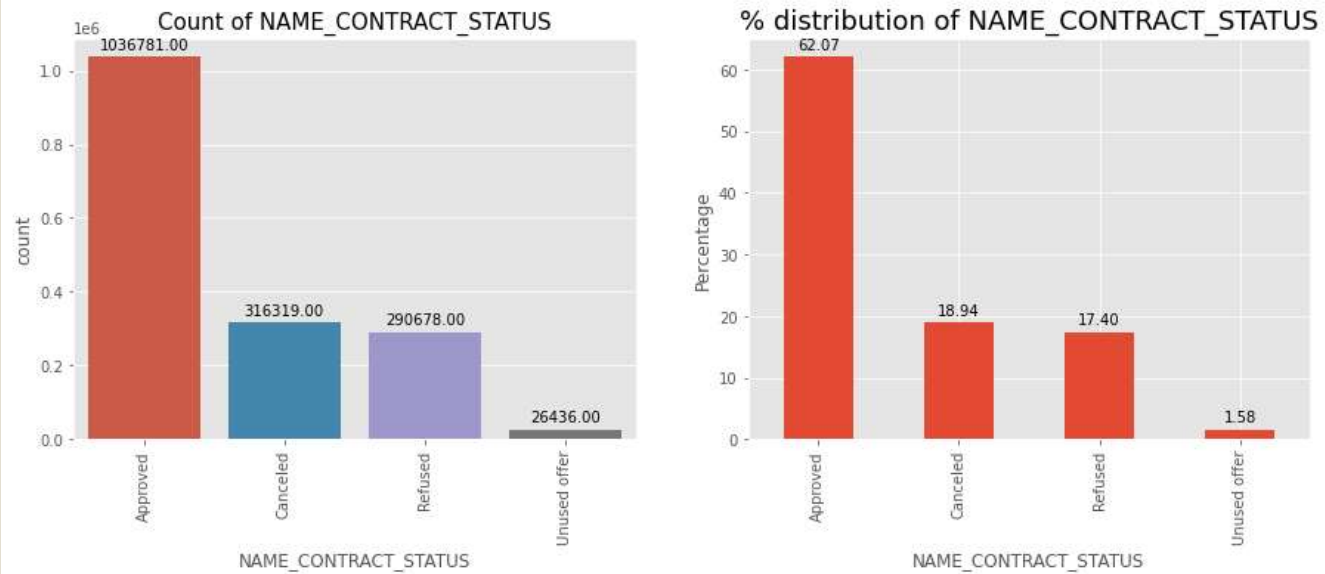


# Univariate Analysis for Previous\_Application Dataset

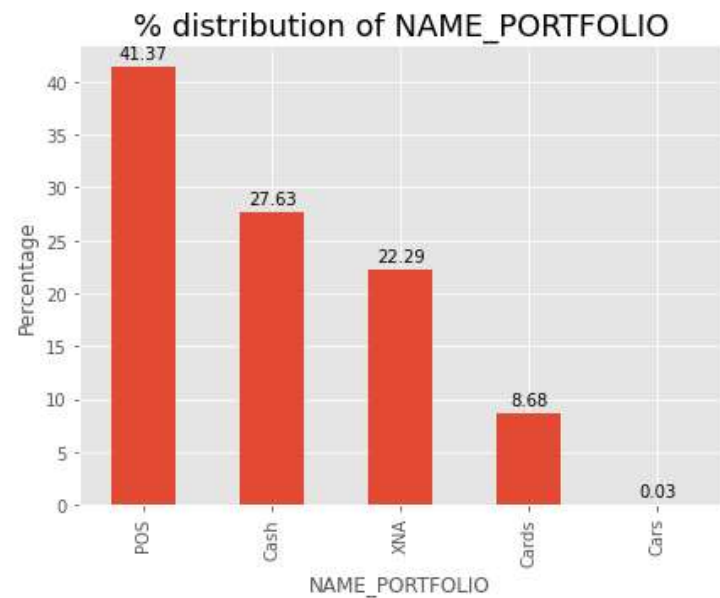
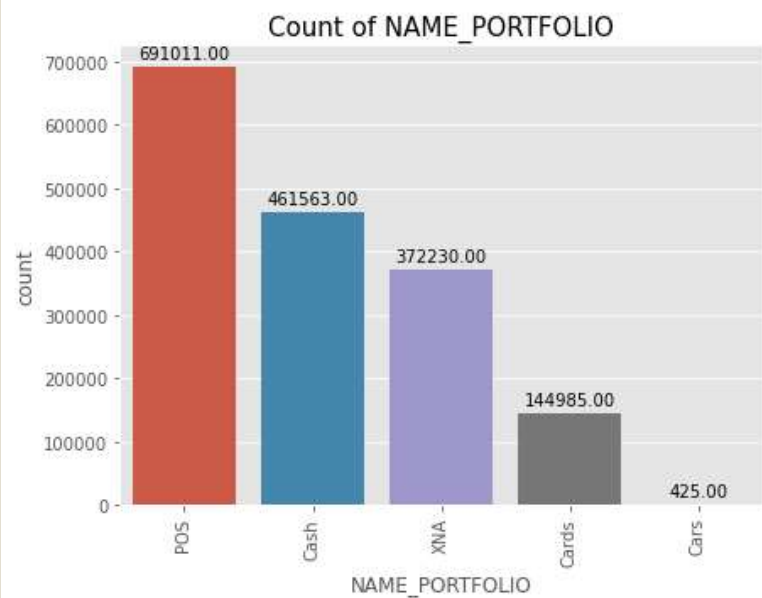
Categorical Univariate Analysis of NAME\_CONTRACT\_TYPE



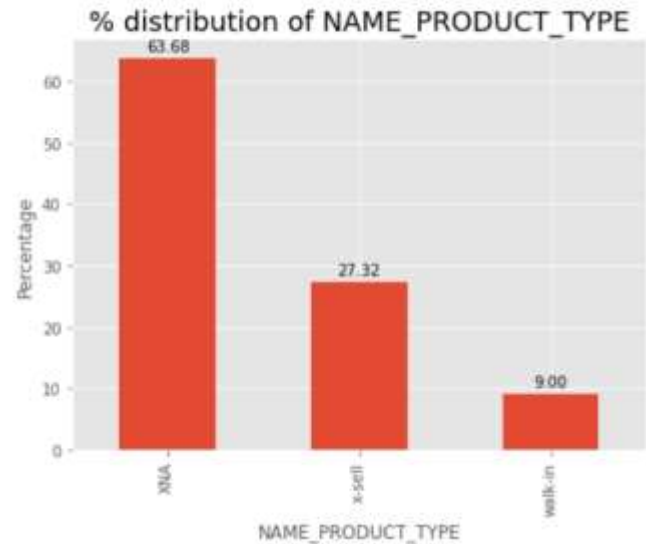
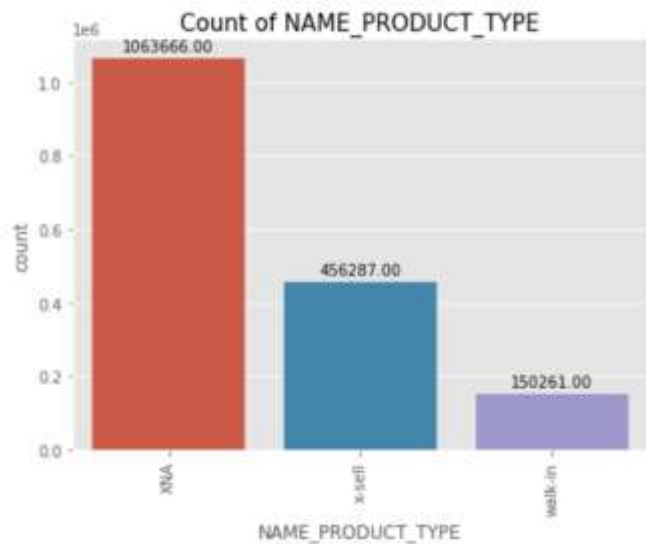
Categorical Univariate Analysis of NAME\_CONTRACT\_STATUS



Categorical Univariate Analysis of NAME\_PORTFOLIO

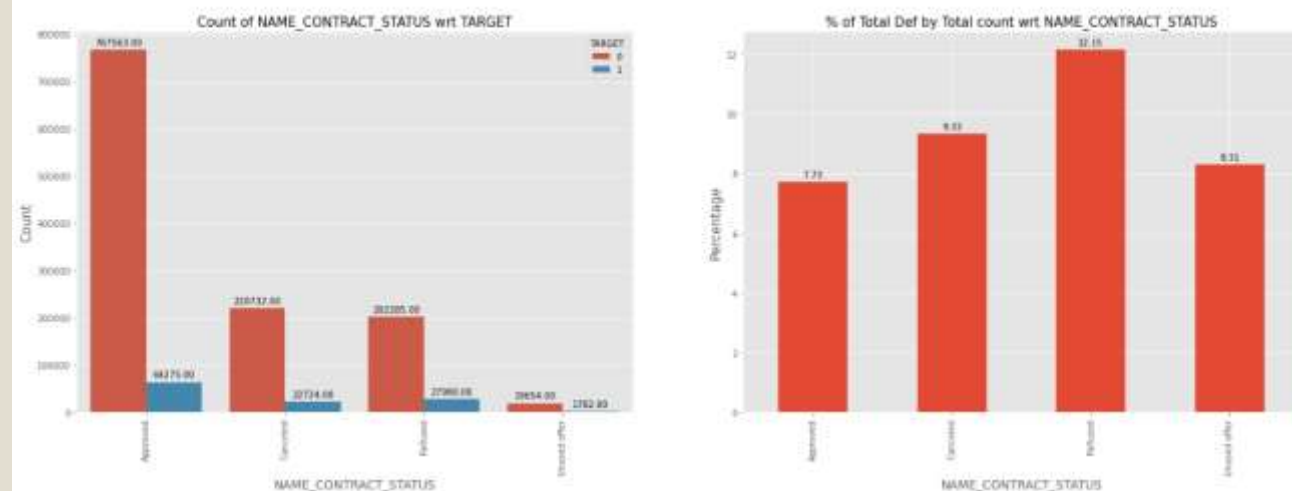


Categorical Univariate Analysis of NAME\_PRODUCT\_TYPE

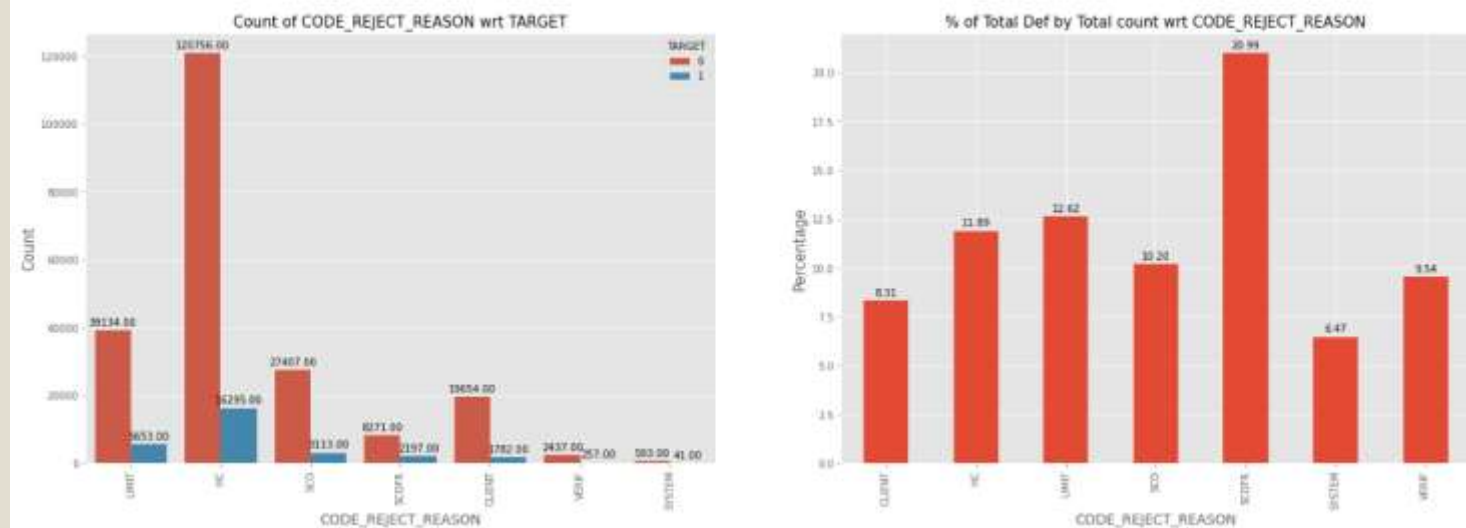


# **Bivariate Analysis for Merged Dataset**

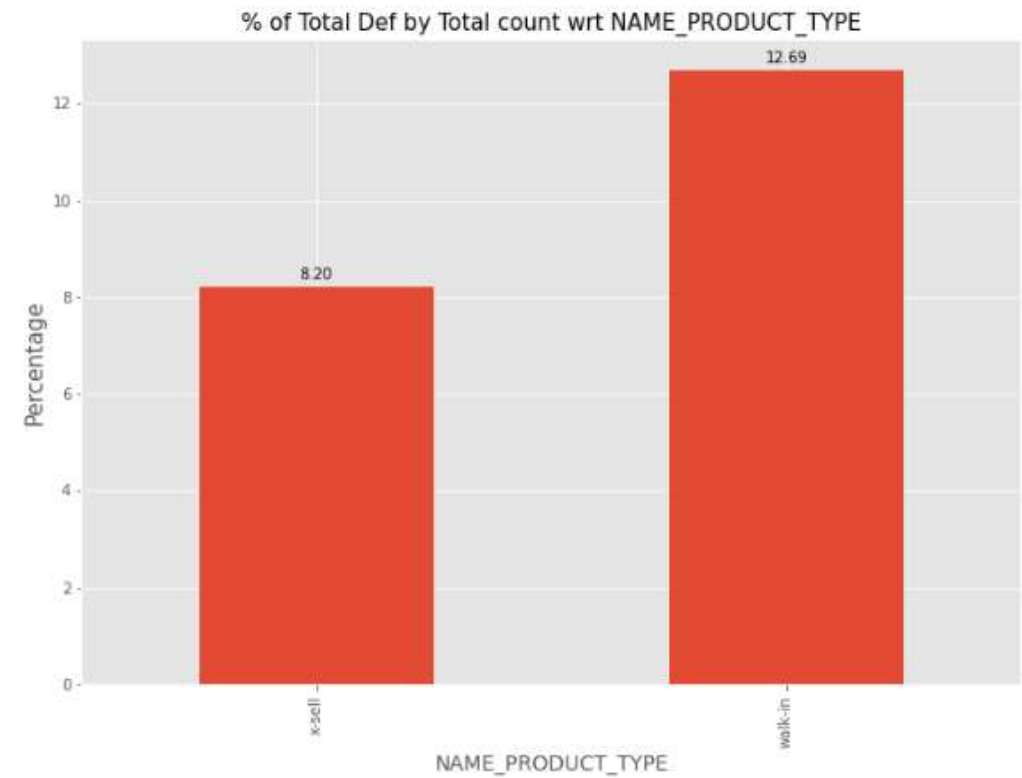
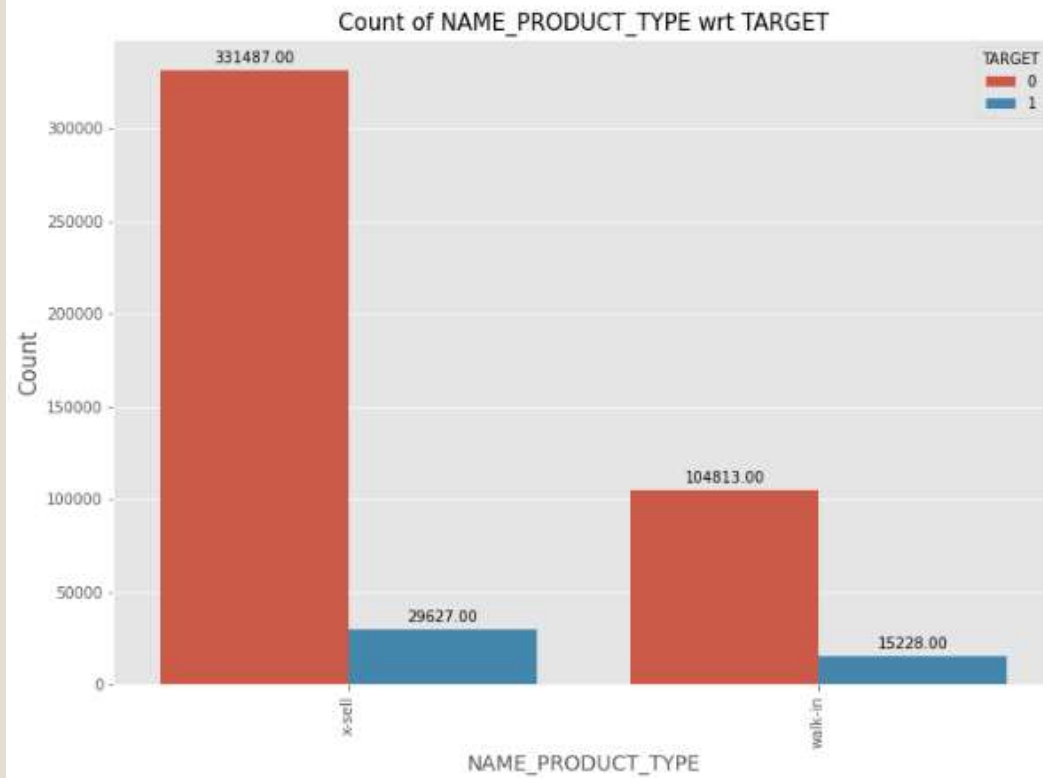
### Segmented Categorical Univariate Analysis of NAME\_CONTRACT\_STATUS



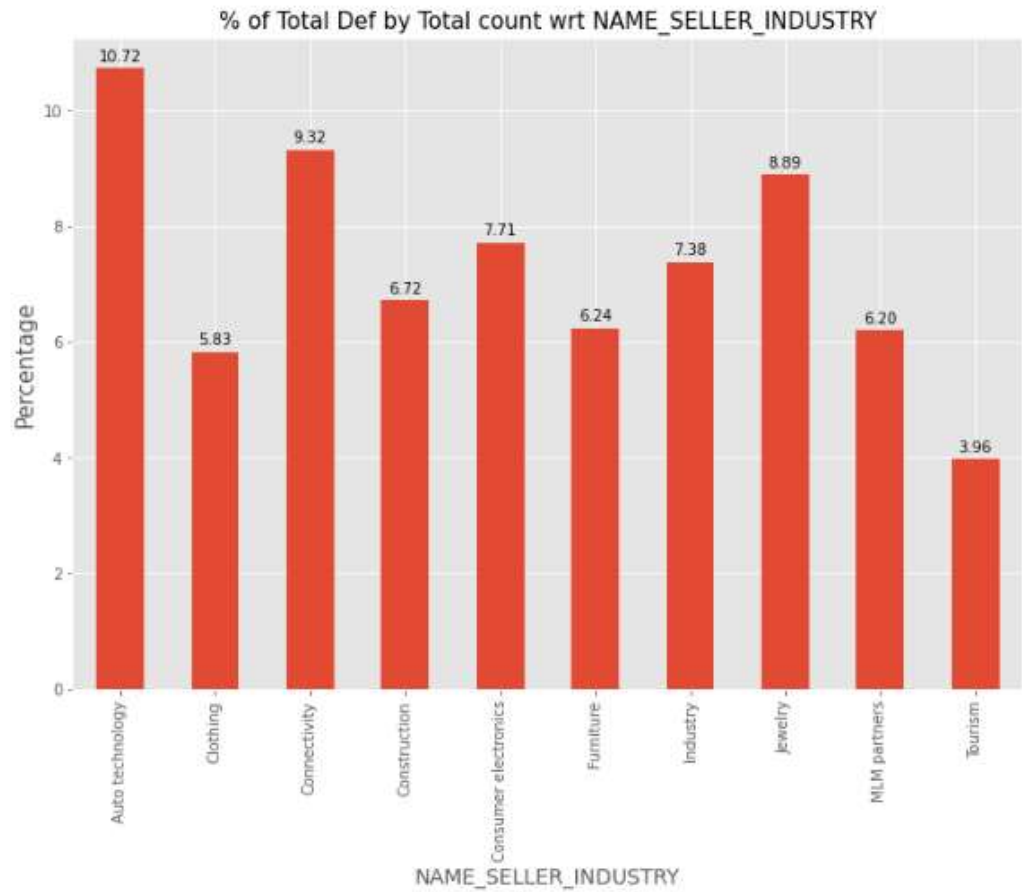
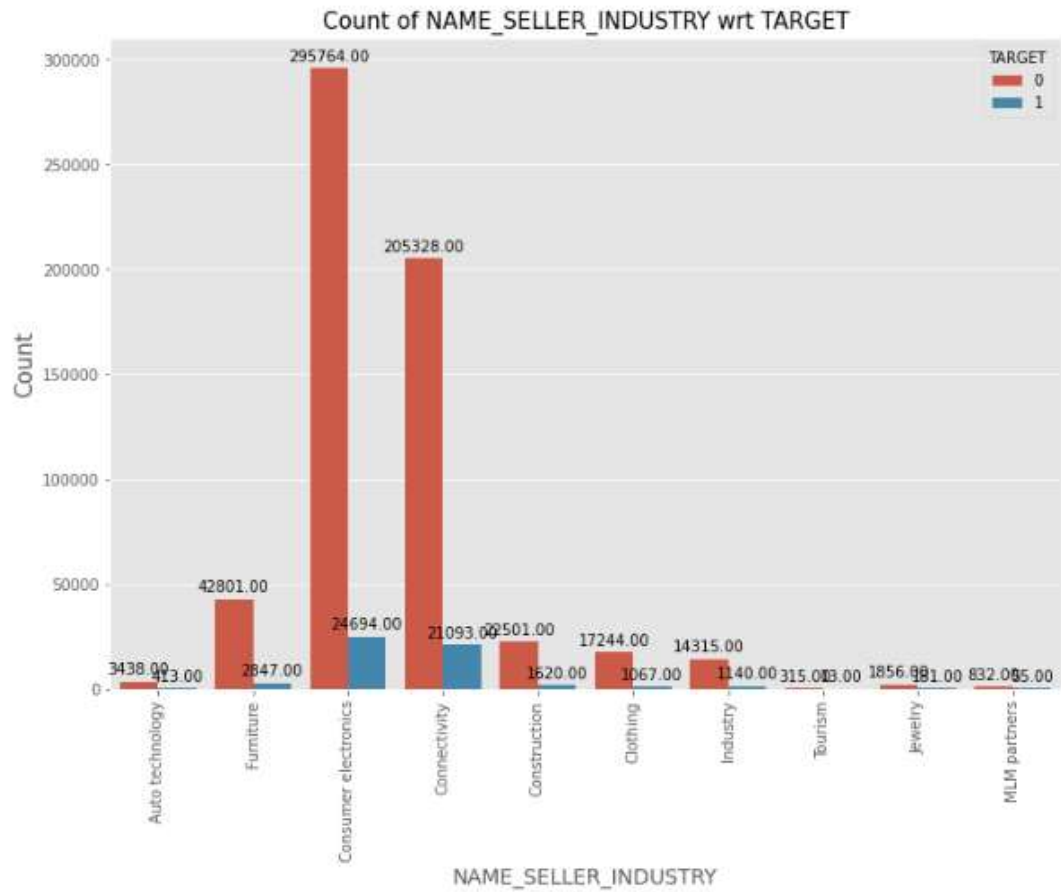
### Segmented Categorical Univariate Analysis of CODE\_REJECT\_REASON



## Segmented Categorical Univariate Analysis of NAME\_PRODUCT\_TYPE

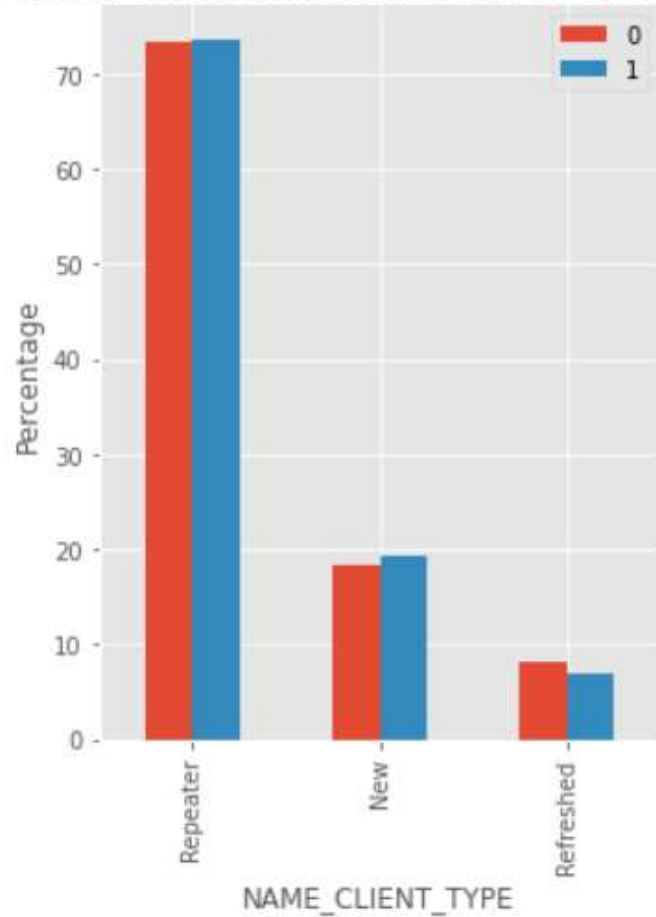


Segmented Categorical Univariate Analysis of NAME\_SELLER\_INDUSTRY

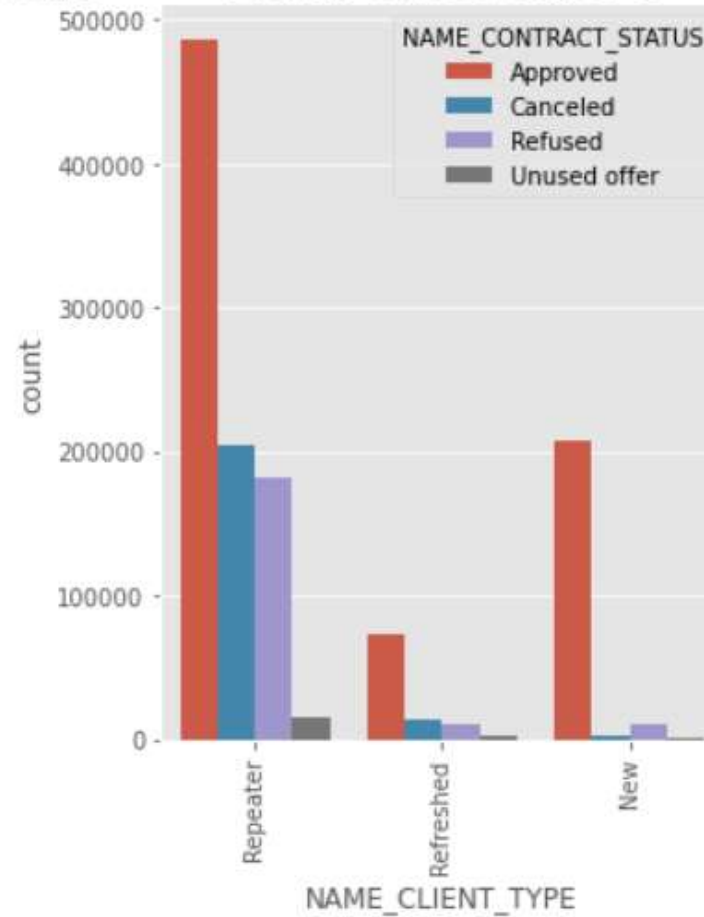


Plotting for NAME\_CLIENT\_TYPE having hue = NAME\_CONTRACT\_STATUS

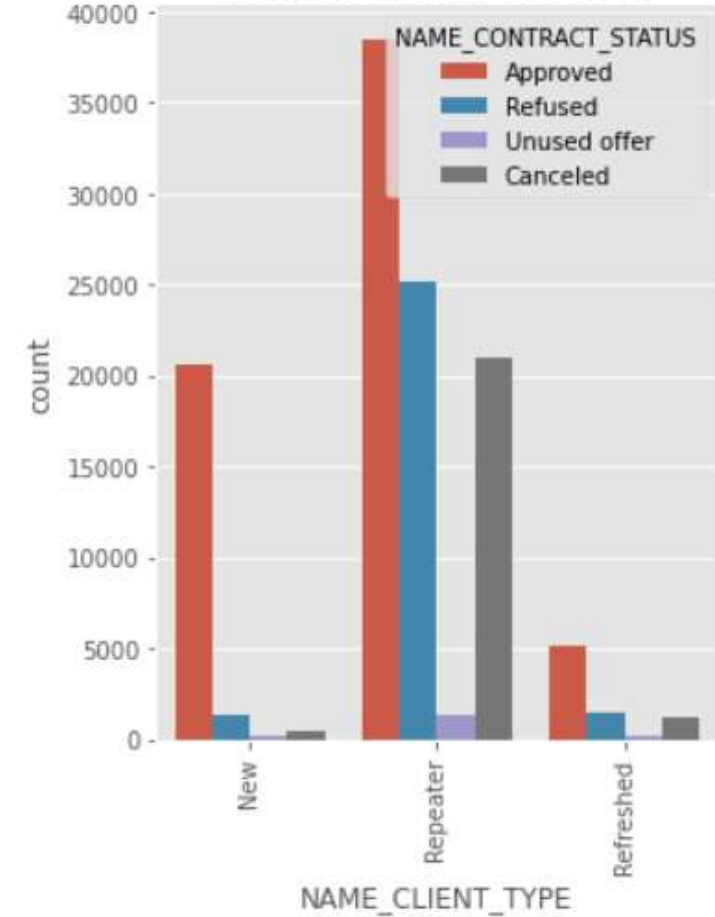
Plotting data for target in terms of Percentage



Plotting data for Target=0



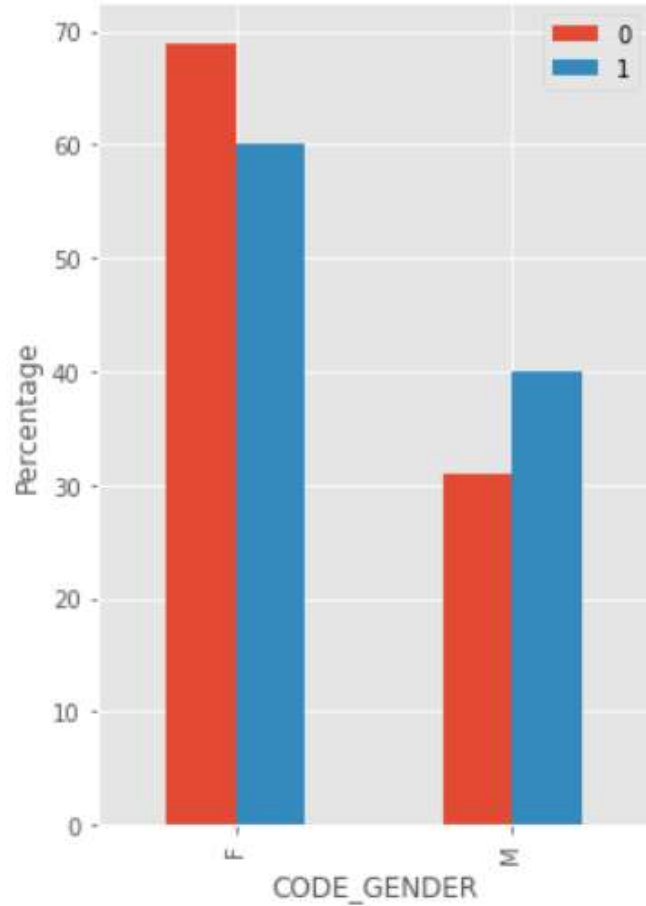
Plotting data for Target=1



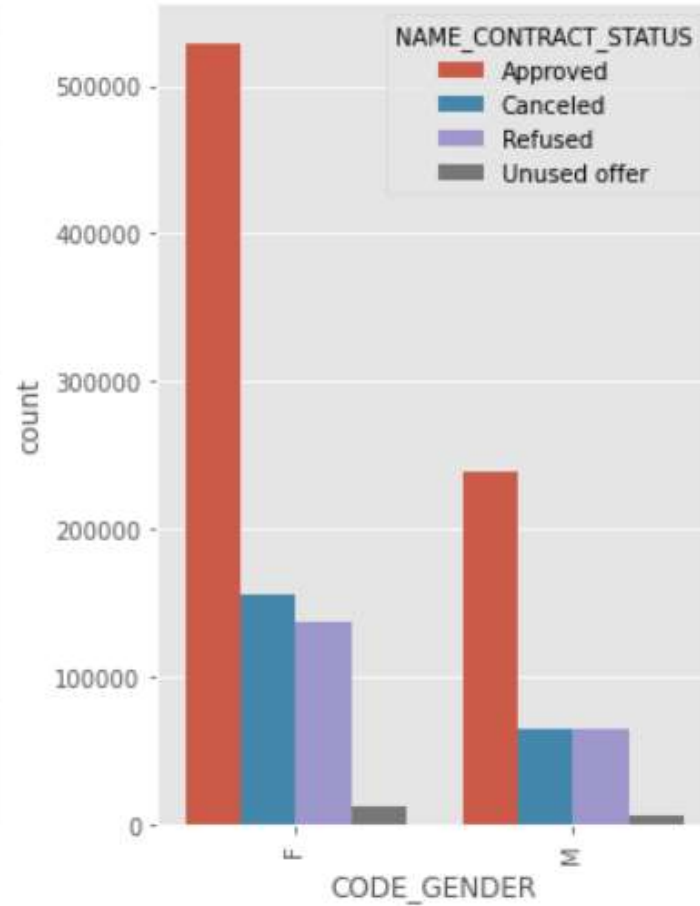


Plotting for CODE\_GENDER having hue = NAME\_CONTRACT\_STATUS

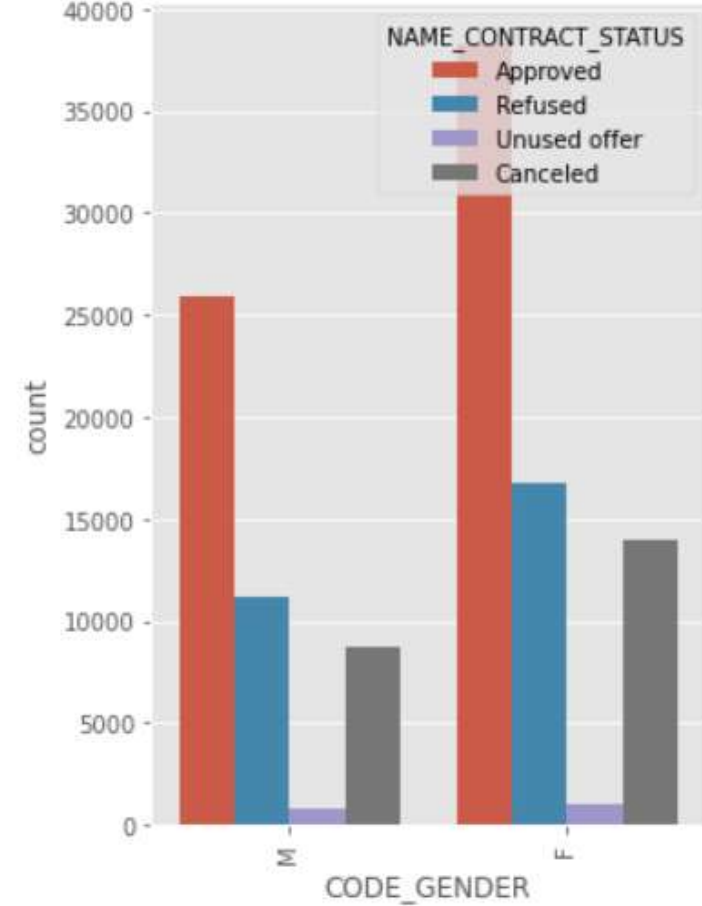
Plotting data for target in terms of Percentage



Plotting data for Target=0

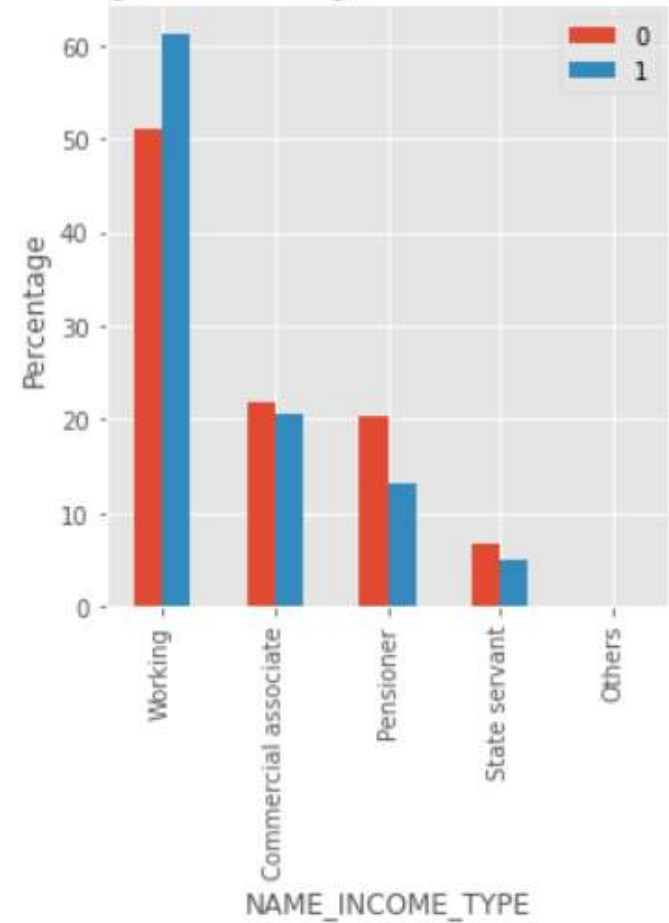


Plotting data for Target=1

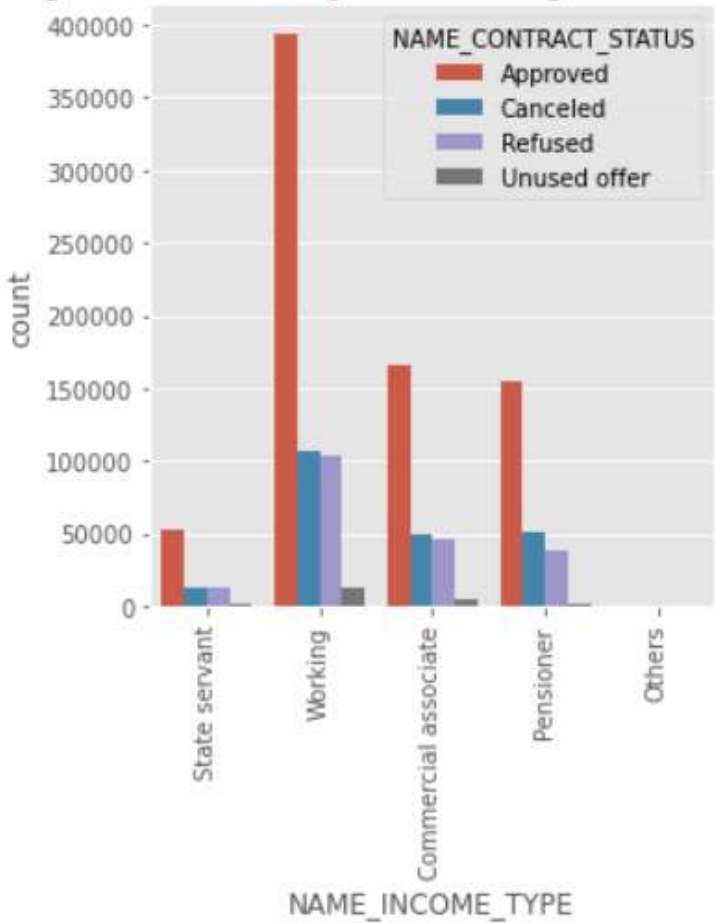


Plotting for NAME\_INCOME\_TYPE having hue = NAME\_CONTRACT\_STATUS

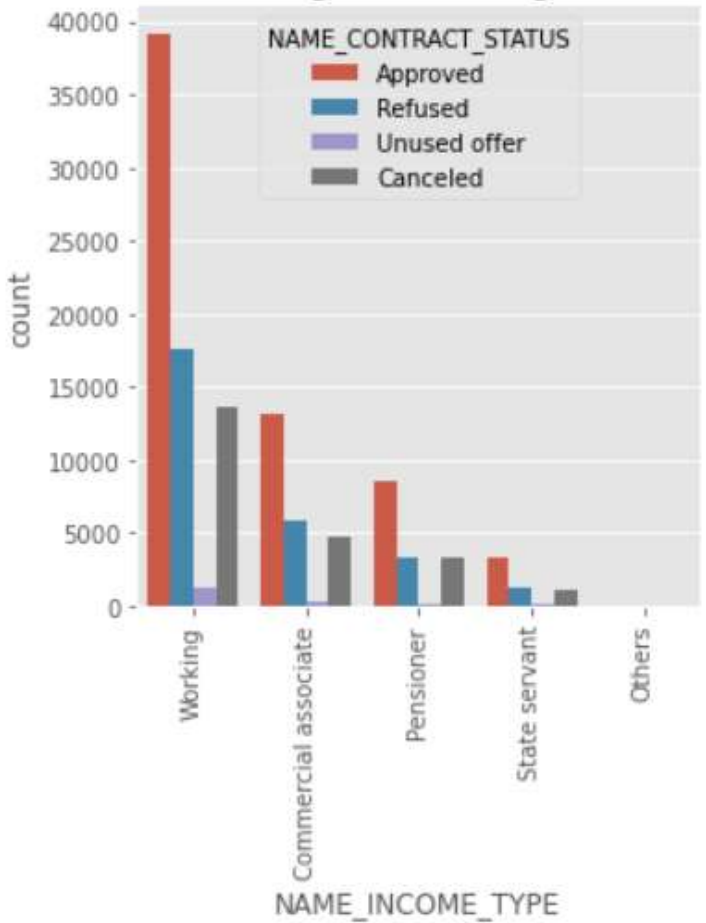
Plotting data for target in terms of Percentage



Plotting data for Target=0



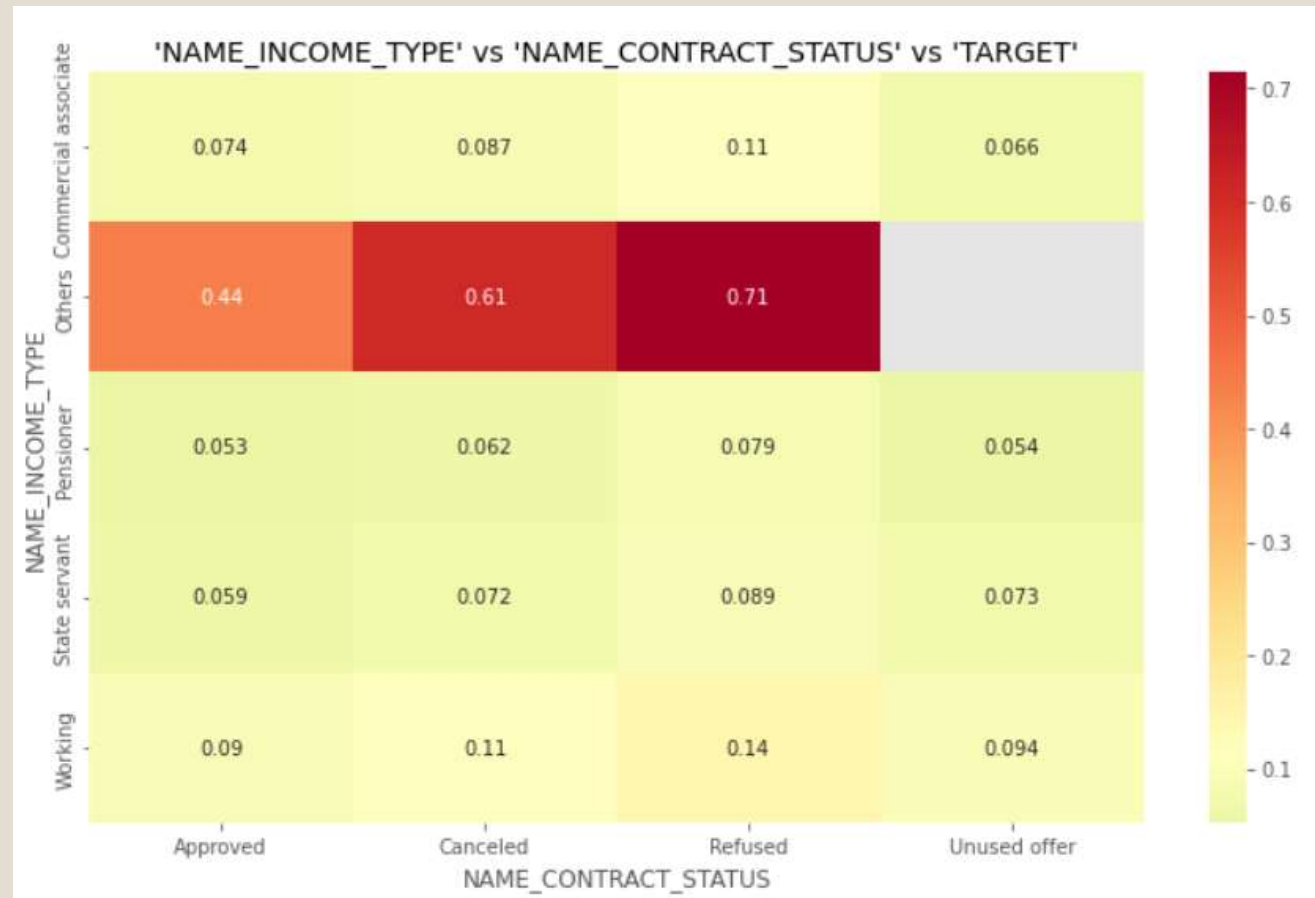
Plotting data for Target=1

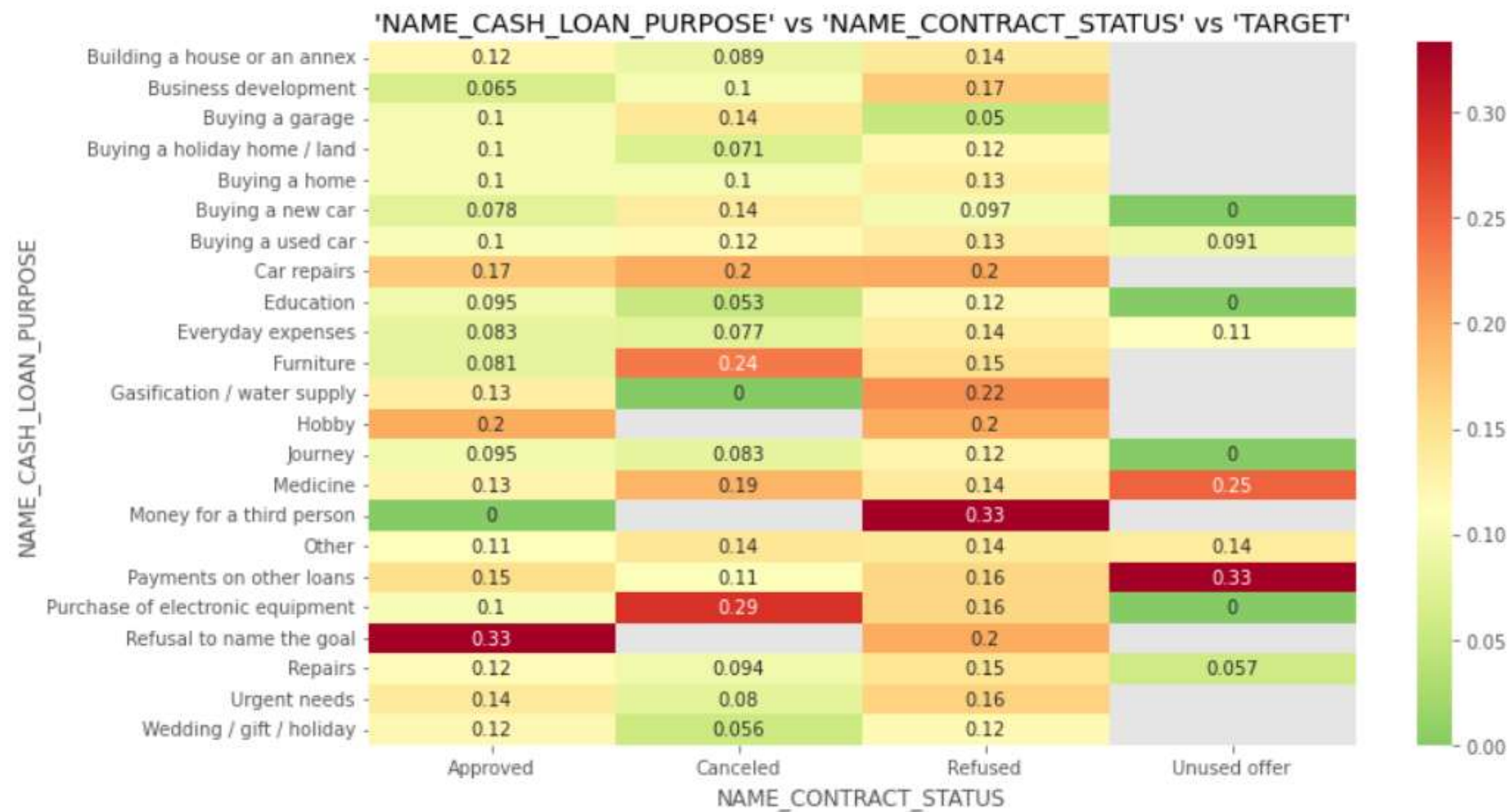


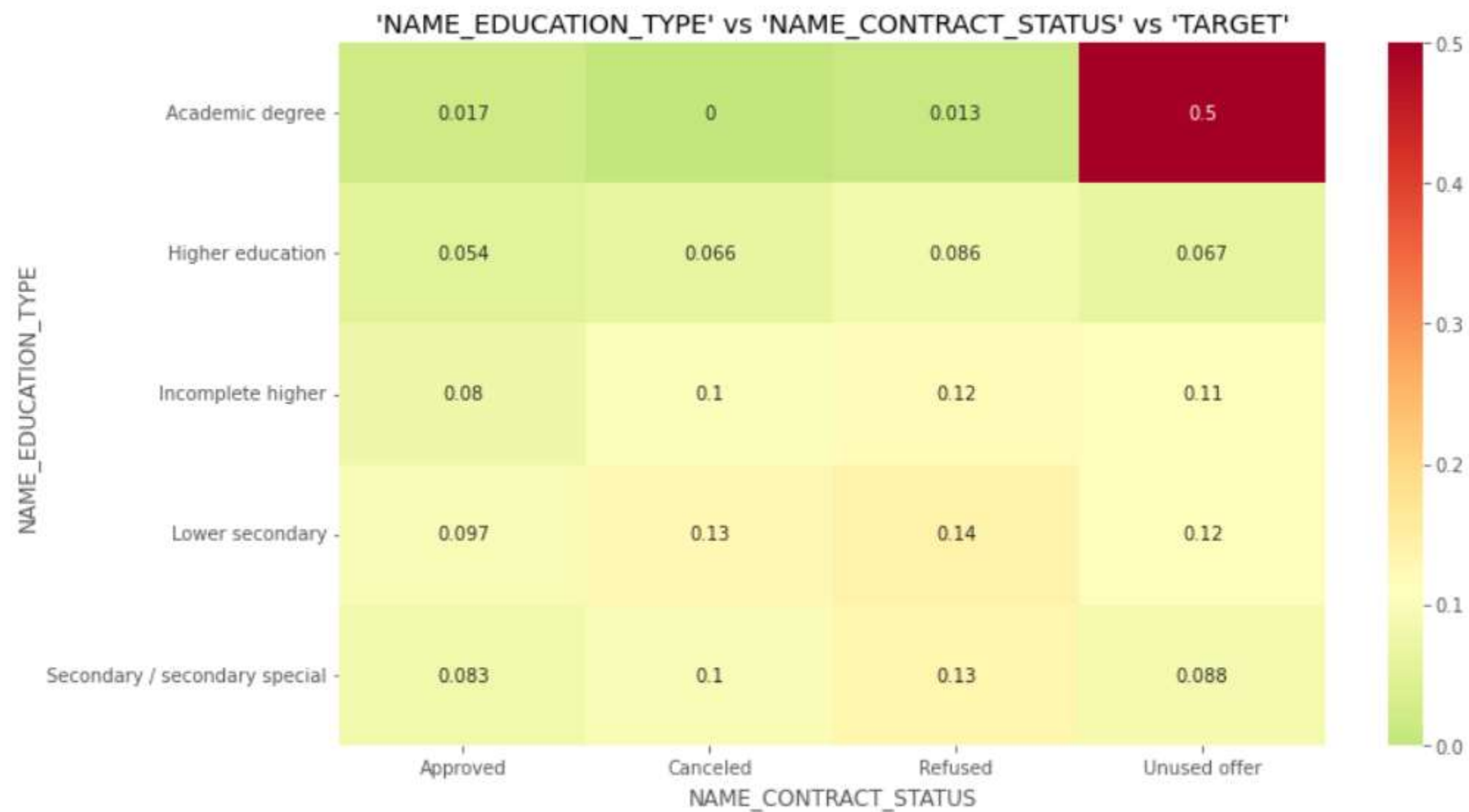
# **Multivariate Analysis**

## Plotting heatmap to see the relationship between few key categorical columns vs 'NAME\_CONTRACT\_STATUS' vs 'TARGET'

*Desired result is to have Green boxes over Approved and Red boxes over Refused. This indicates that people who are known high defaulters are being refused and low defaulters are being approved for loans*







# Inference:

- From the heatmaps, we can identify the category of people in spite of having very less defaulting rates being rejected. These are few of missed business opportunity to the bank and we can also see in some cases people having higher default rates being approved. Therefore, to improve the business/profits banks should try to provide loans to people having lesser defaulting rates.

**Thank You.**