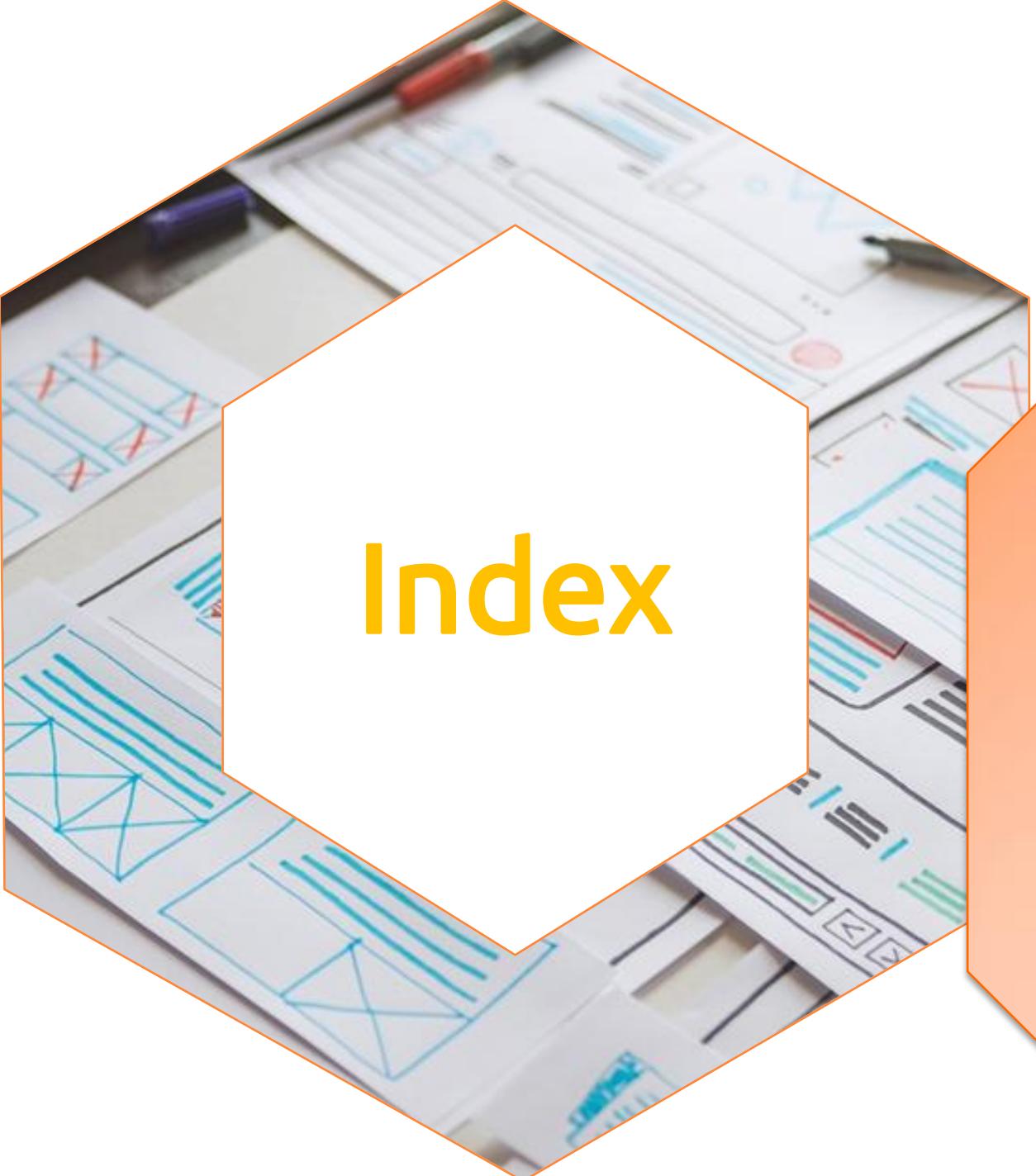


# CREDIT EDA ANALYSIS

BY GAYATRI BHINGE





# Index

1. Introduction
2. Problem Statement
3. Outliers
4. Imbalance in Data
5. Univariate Analysis
6. Bivariate Analysis
7. Conclusion

# Introduction



This case study-assignment is based on Exploratory Data Analysis. The aim of this assignment is to apply EDA on real time business case. In this case study we will learn how to develop understanding of risk analytics in banking and financial services and as this dataset is applications of loan by client, so by analyzing this data we understand how data is used to minimize the risk of losing money while approving the loan to applicants.

# Problem Statement

To take decision of loan application is very complicated task for any finance company, and it is big risk to approve loan without analysing applicant's profile. It is hard to find insufficient or non-existent credit history of any applicant and due to this applicant may be take advantage of becoming defaulter. As a data analyst for finance company we have to apply EDA on given data to analyses the patterns present in that datasets. When the company receives a loan application, we must ensure that applicant capable to repaying loan amount without fault, and loan is approved after doing analysis of applicant's profile.

There are 2 types of risks depend on decision of loan approval are as below:

- 1] Rejecting the loan application to capable applicant leads to business loss.
- 2] Accepting application of defaulter applicant, and if not able to replay loan leads to financial loss. In given dataset there are 2 scenarios, 'Target': 1 and 0.

As task performing EDA on this problem statement we have to understand the driving factors (or driver variables) behind loan default, i.e. the variables which are strong indicators of default. The company can utilise this outcome before approving loan to any applicant.

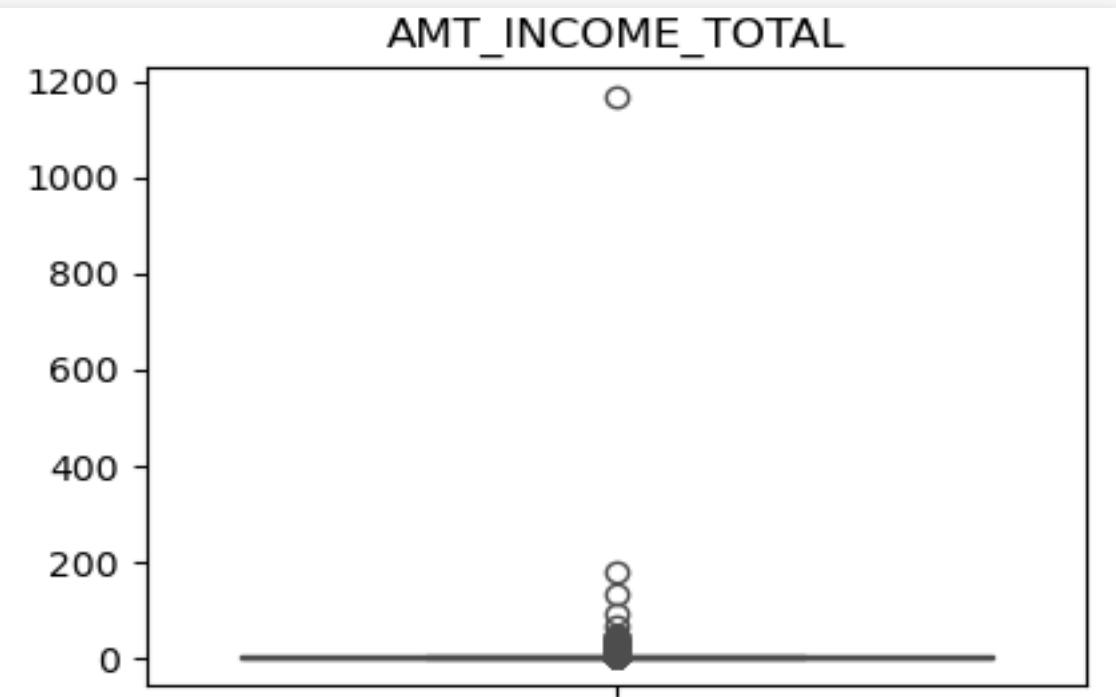
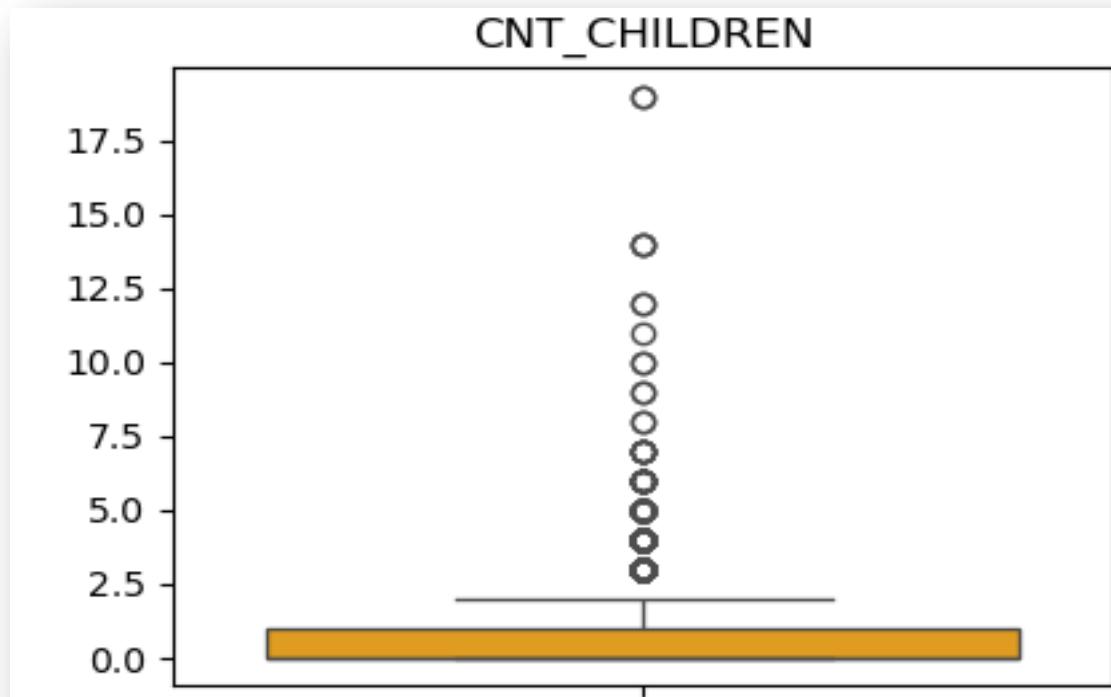
# Data Understanding

This dataset has 3 files as explained below:

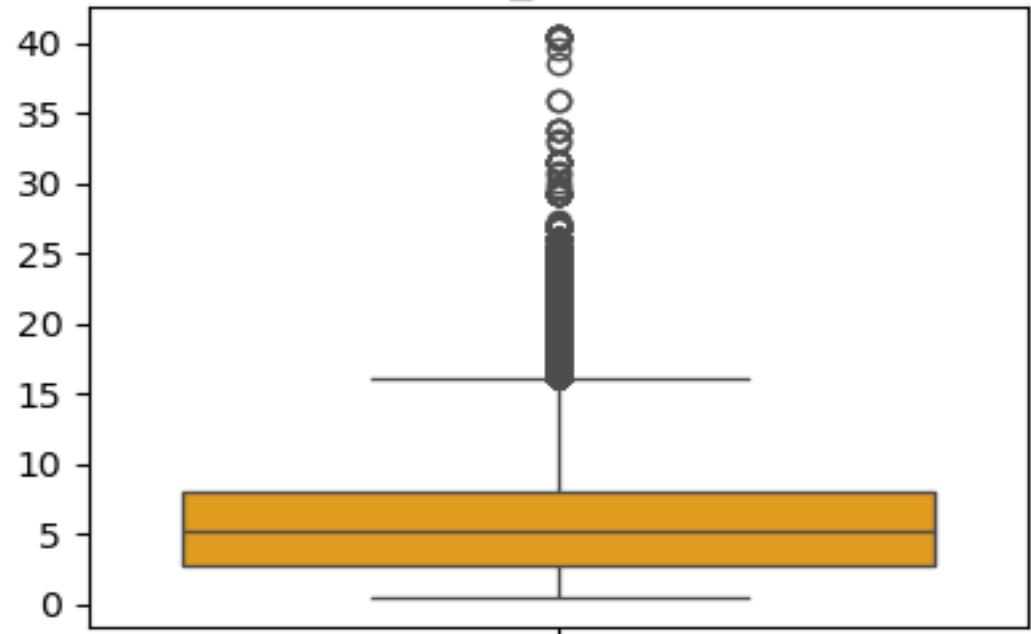
- ❖ 'application\_data.csv' contains all the information of the client at the time of application. The data is about whether a client has payment difficulties.
- ❖ 'previous\_application.csv' contains information about the client's previous loan data. It contains the data on whether the previous application had been Approved, Cancelled, Refused or Unused offer.
- ❖ 'columns\_description.csv' is data dictionary which describes the meaning of the variables.

# Outliers

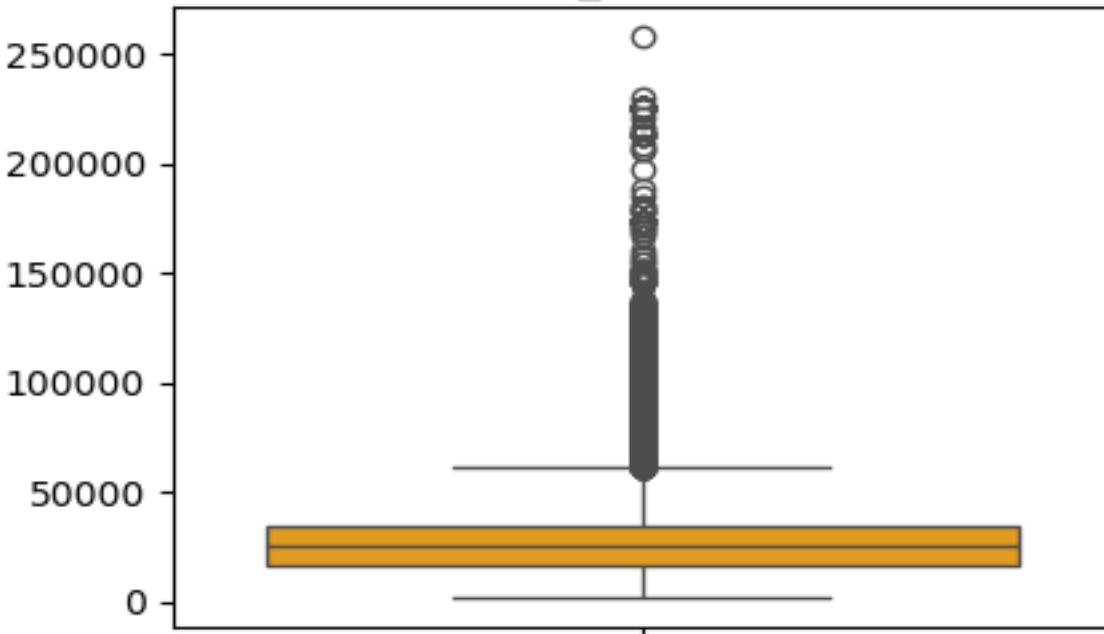
Application Data -



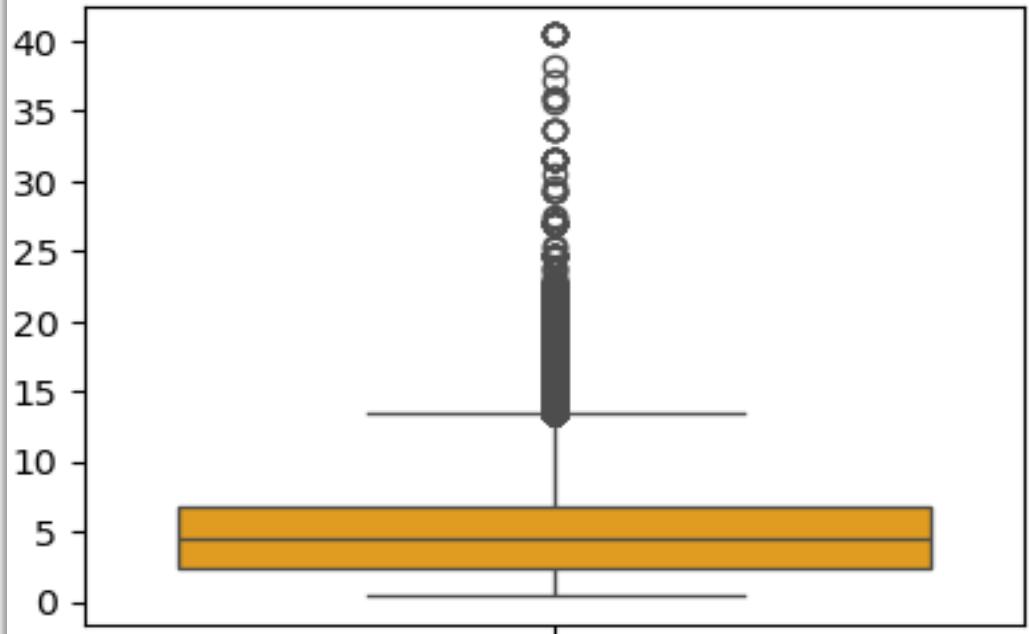
AMT\_CREDIT



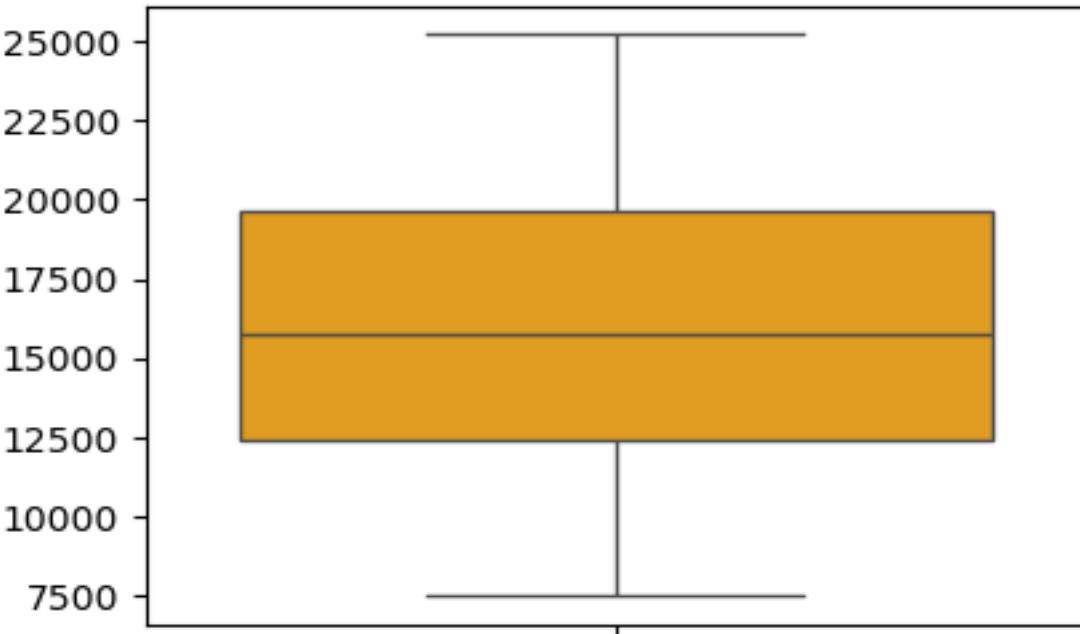
AMT\_ANNUITY



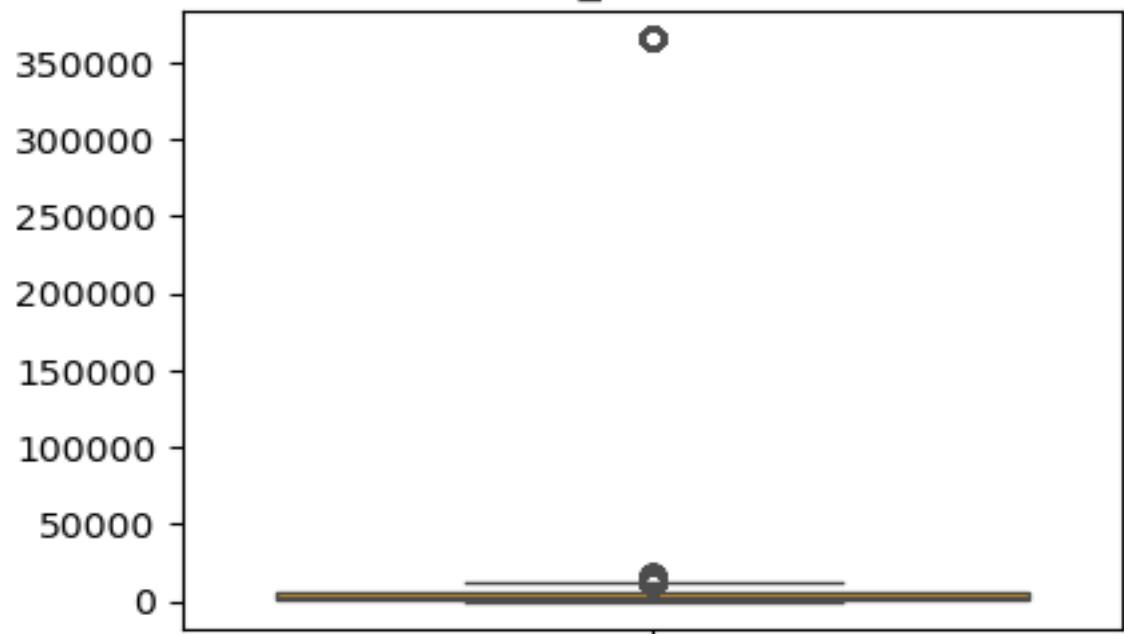
AMT\_GOODS\_PRICE



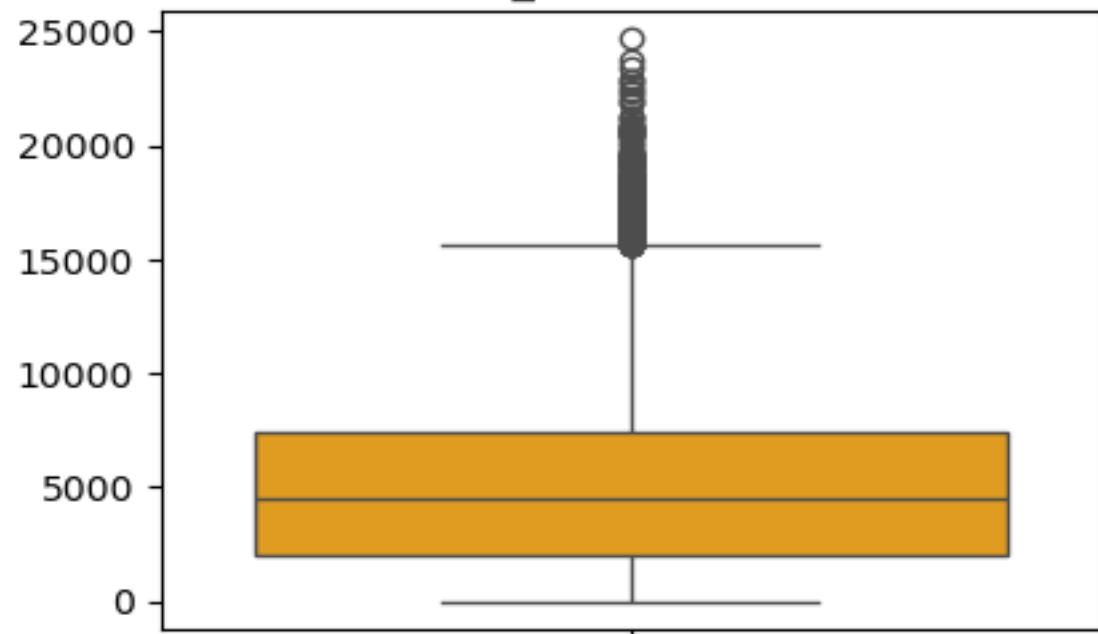
DAYS\_BIRTH



DAYs\_EMPLOYED



DAYs\_REGISTRATION

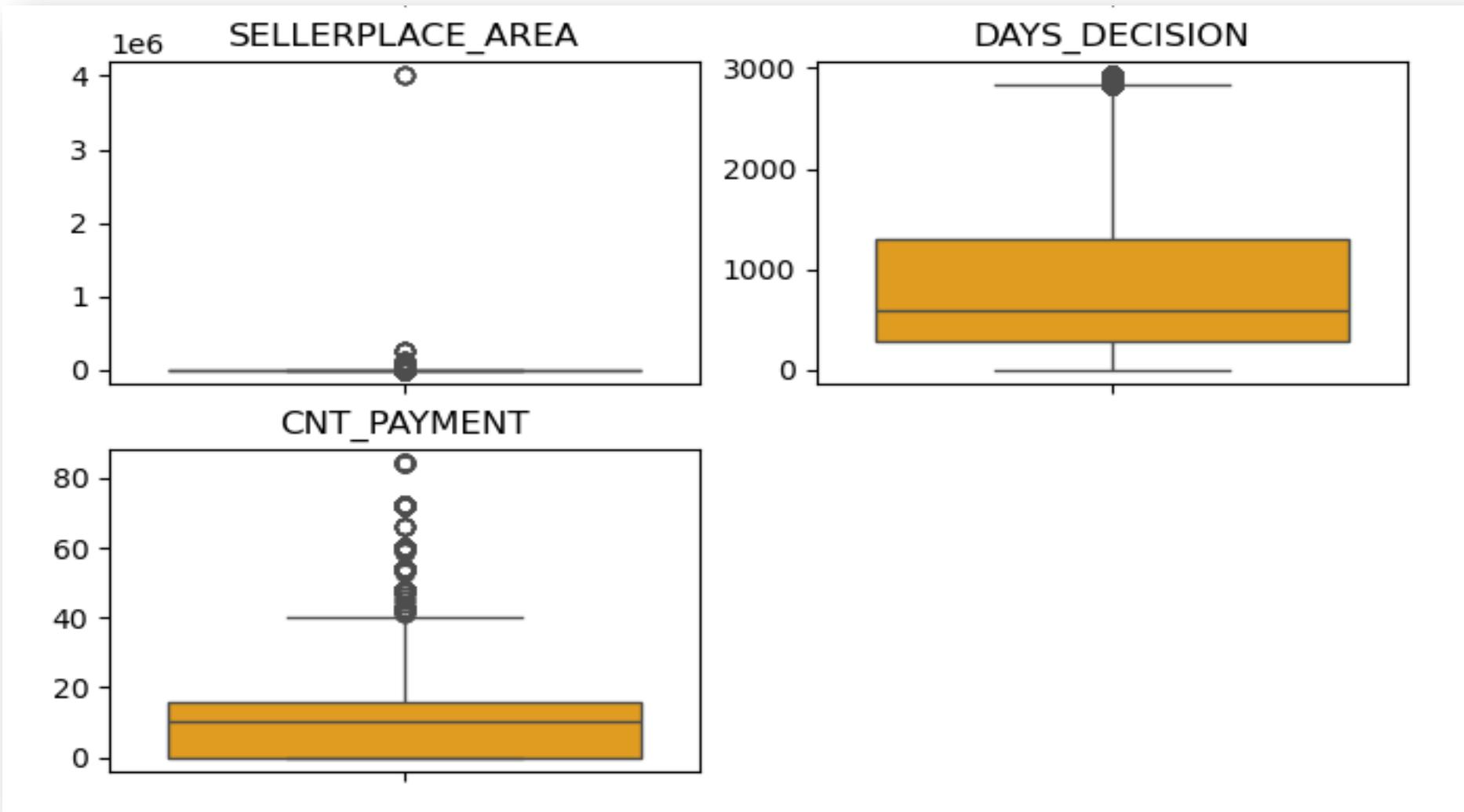


## Inferences

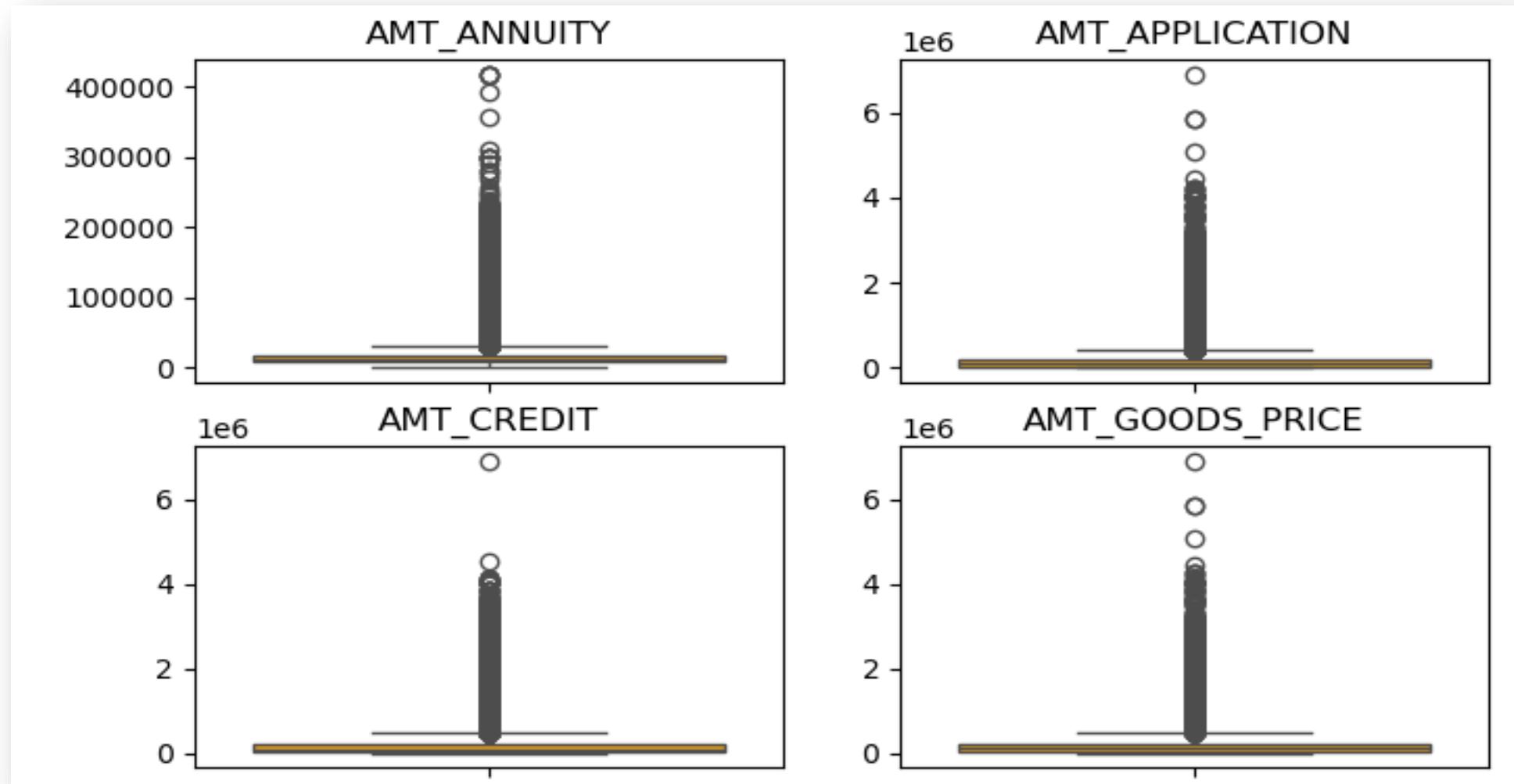
Above boxplots analysis –

- AMT\_ANNUITY, AMT\_CREDIT, AMT\_GOODS\_PRICE, CNT\_CHILDREN have some number of outliers.
- Few of the loan applicants have high income when compared to the others as data in AMT\_INCOME\_TOTAL has huge number of outliers.
- DAYS\_BIRTH has no outliers which means the data available is reliable.
- DAYS\_EMPLOYED has outlier values around 350000(days) which is around 958 years which is impossible and hence this has to be incorrect entry.
- Outcome for CNT\_CHILDREN looks like not correct as number of children's showing outliers.

## Previous Applications Data



## Previous Applications Data





## Inferences

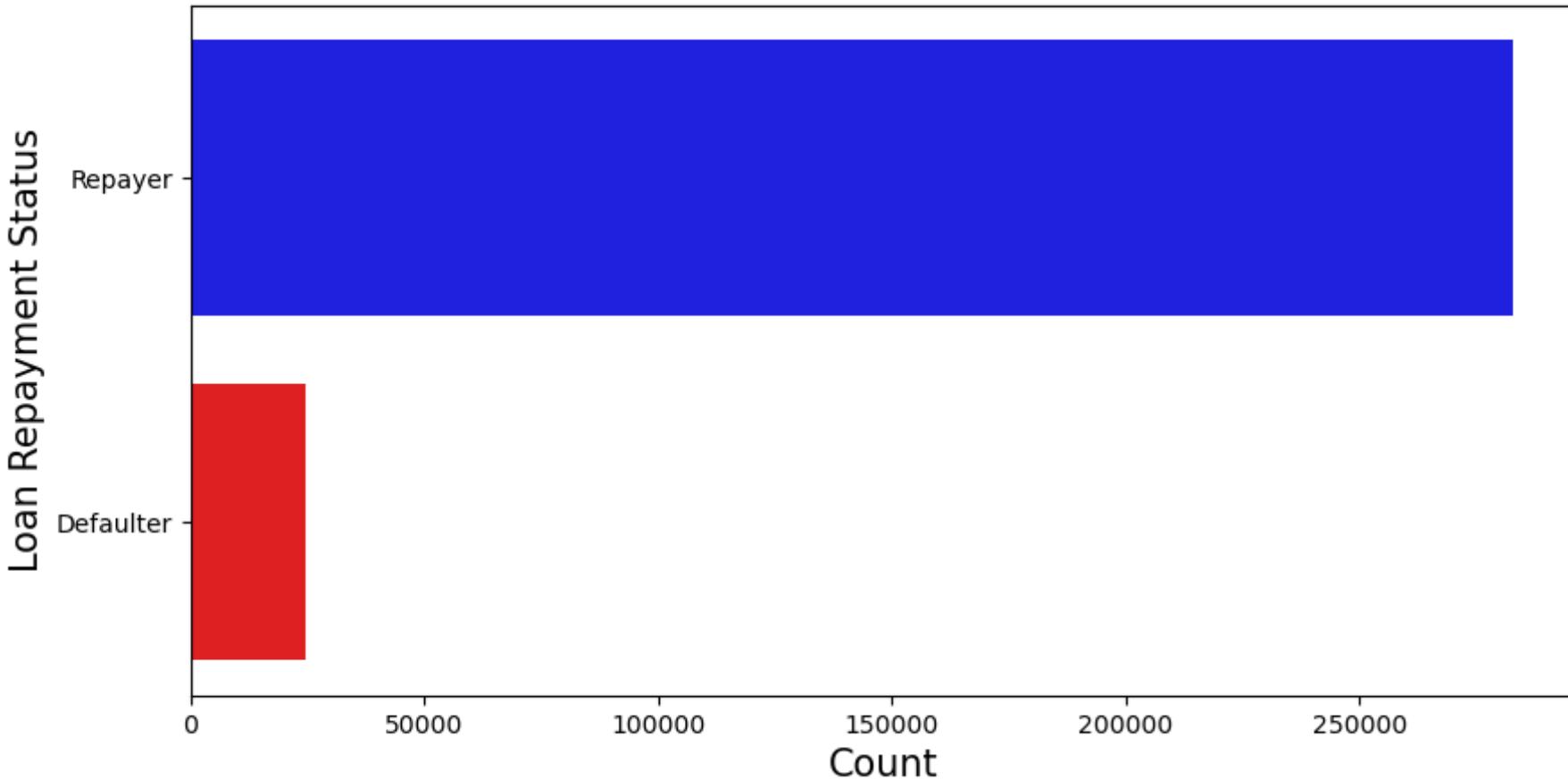
Above boxplots analysis –

- DAYS\_DECISION has little number of outliers indicating that these previous applications decisions were taken long back.
- CNT\_PAYMENT has few outlier values.
- AMT\_ANNUITY, AMT\_APPLICATION, AMT\_CREDIT, AMT\_GOODS\_PRICE, SELLERPLACE\_AREA have huge number of outliers.

# Imbalance in Data

## Application Data

Imbalance Plotting (Repayer Vs Defaulter)



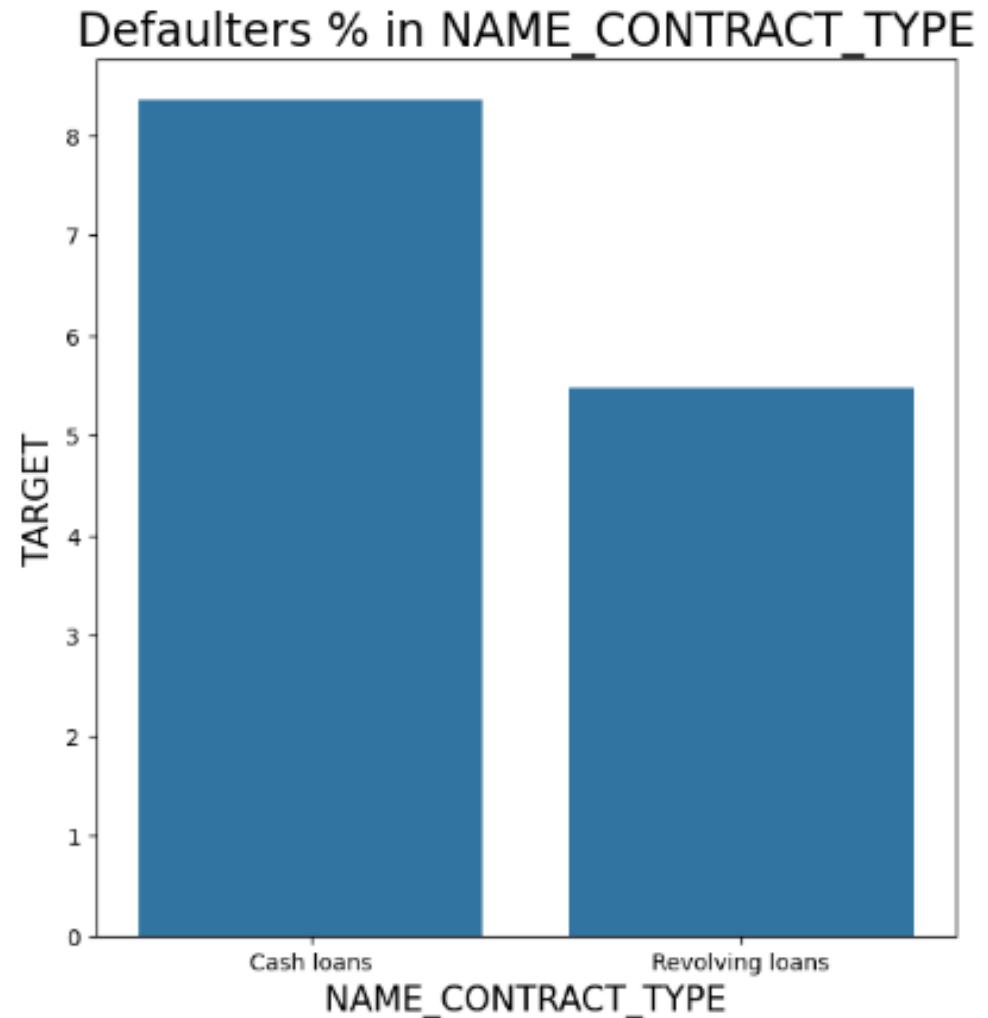
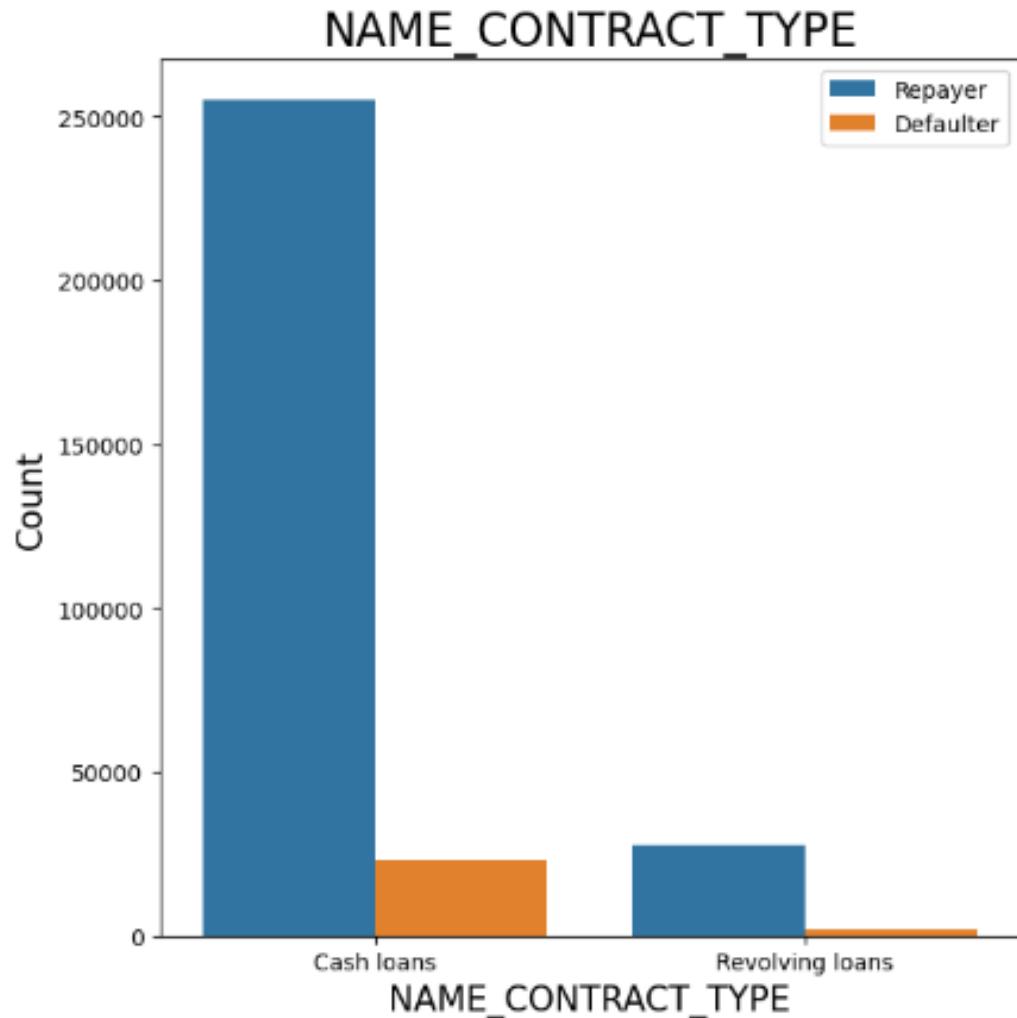
## Inferences

As we see in this graph, there is significant difference between 'Repayer' and 'Defaulter'. We will check ratio as below –

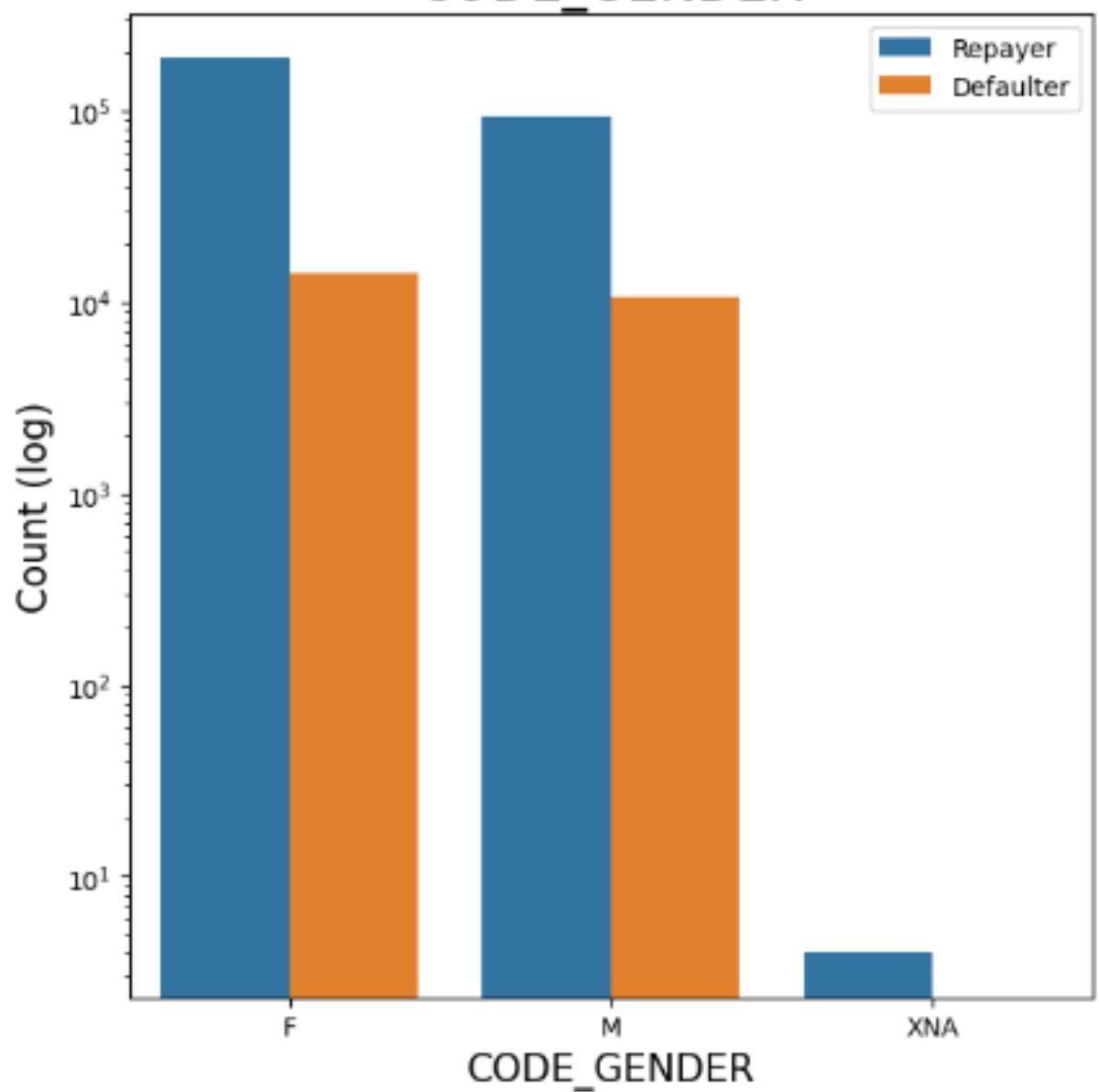
- Repayer Percentage is 91.93%
- Defaulter Percentage is 8.07%
- Imbalance Ratio with respect to Repayer and Defaulter is given: 11.39/1 (approx)

# Univariate Analysis

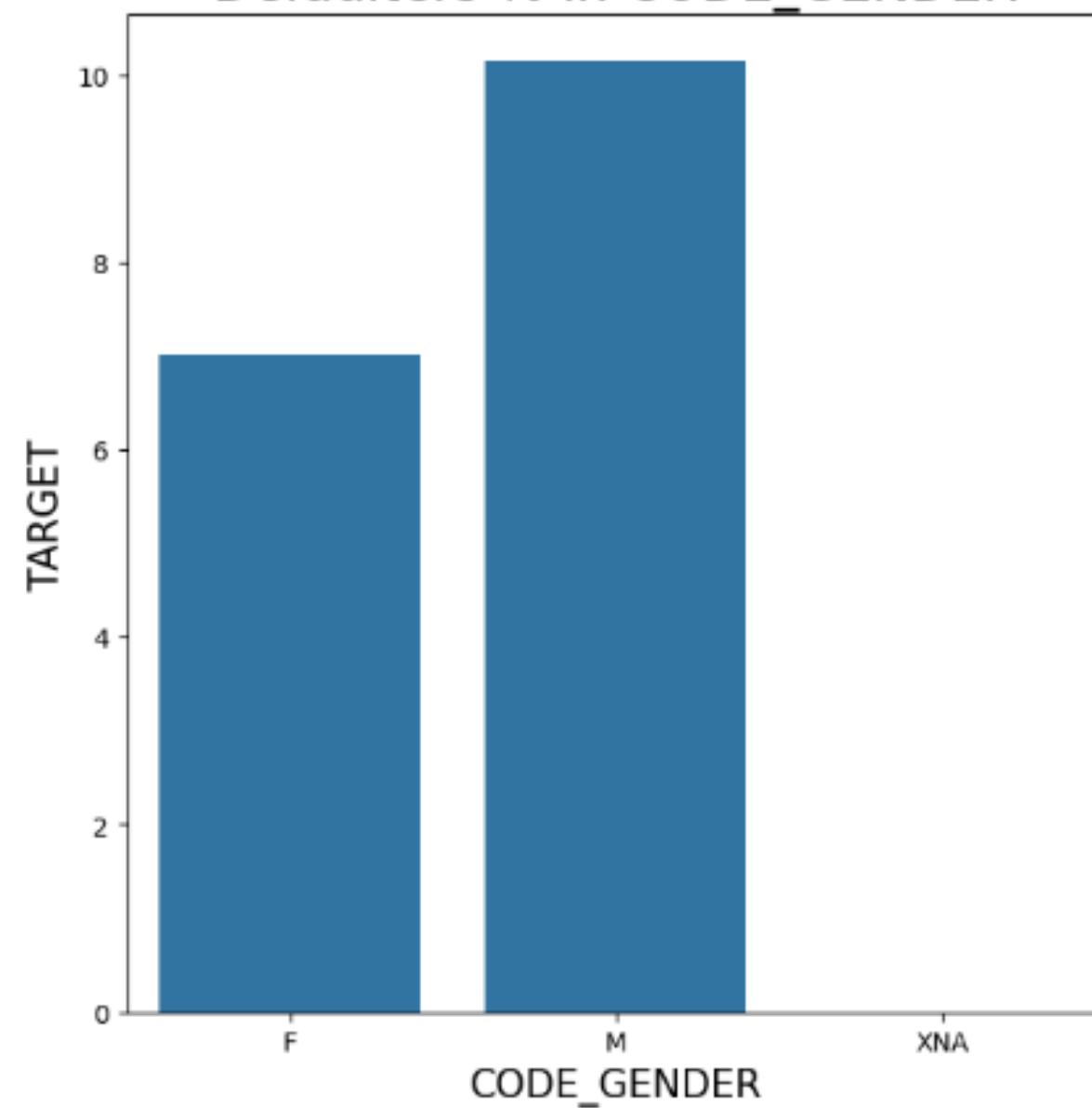
## Application Data



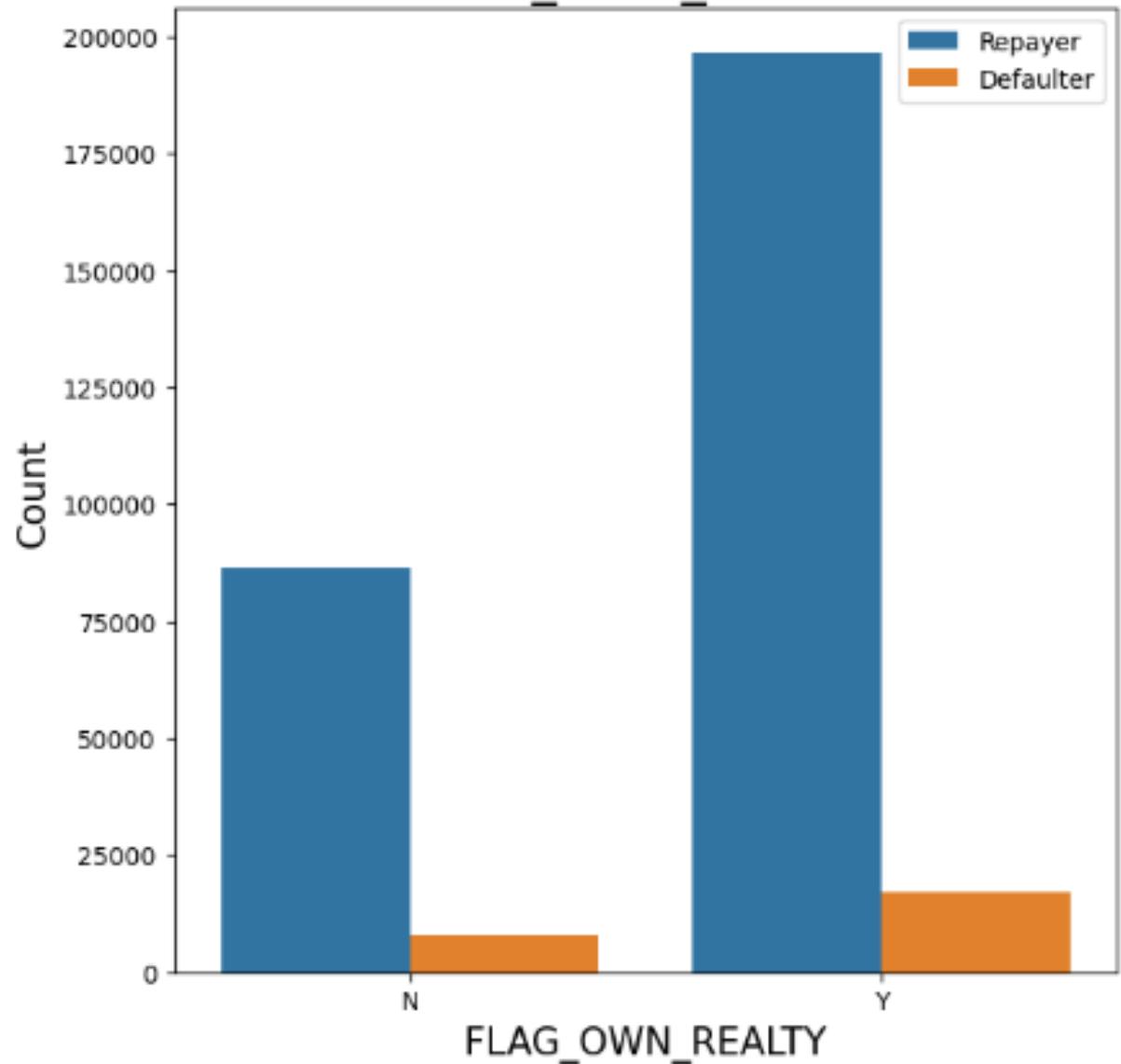
## CODE\_GENDER



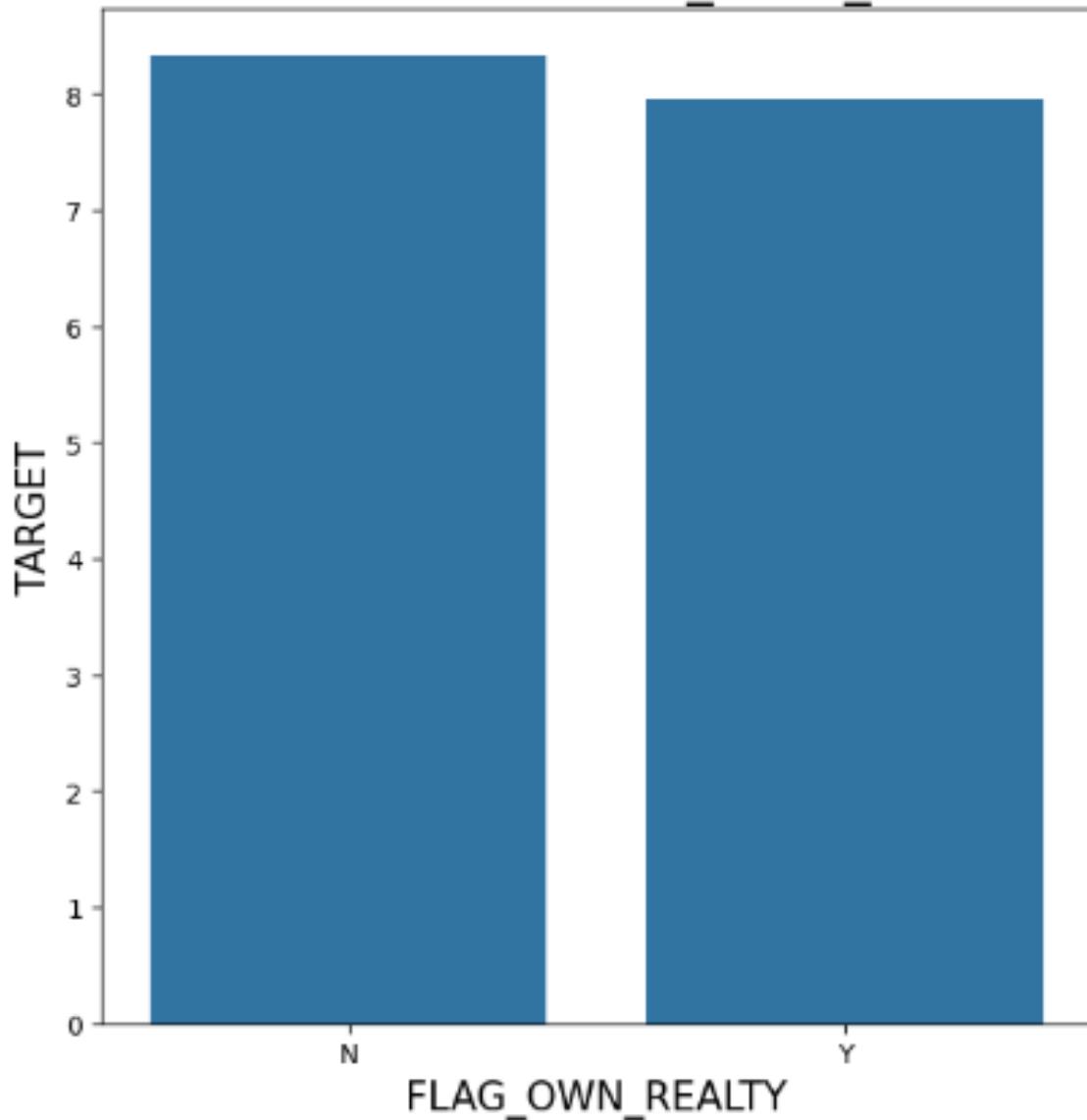
## Defaulters % in CODE\_GENDER



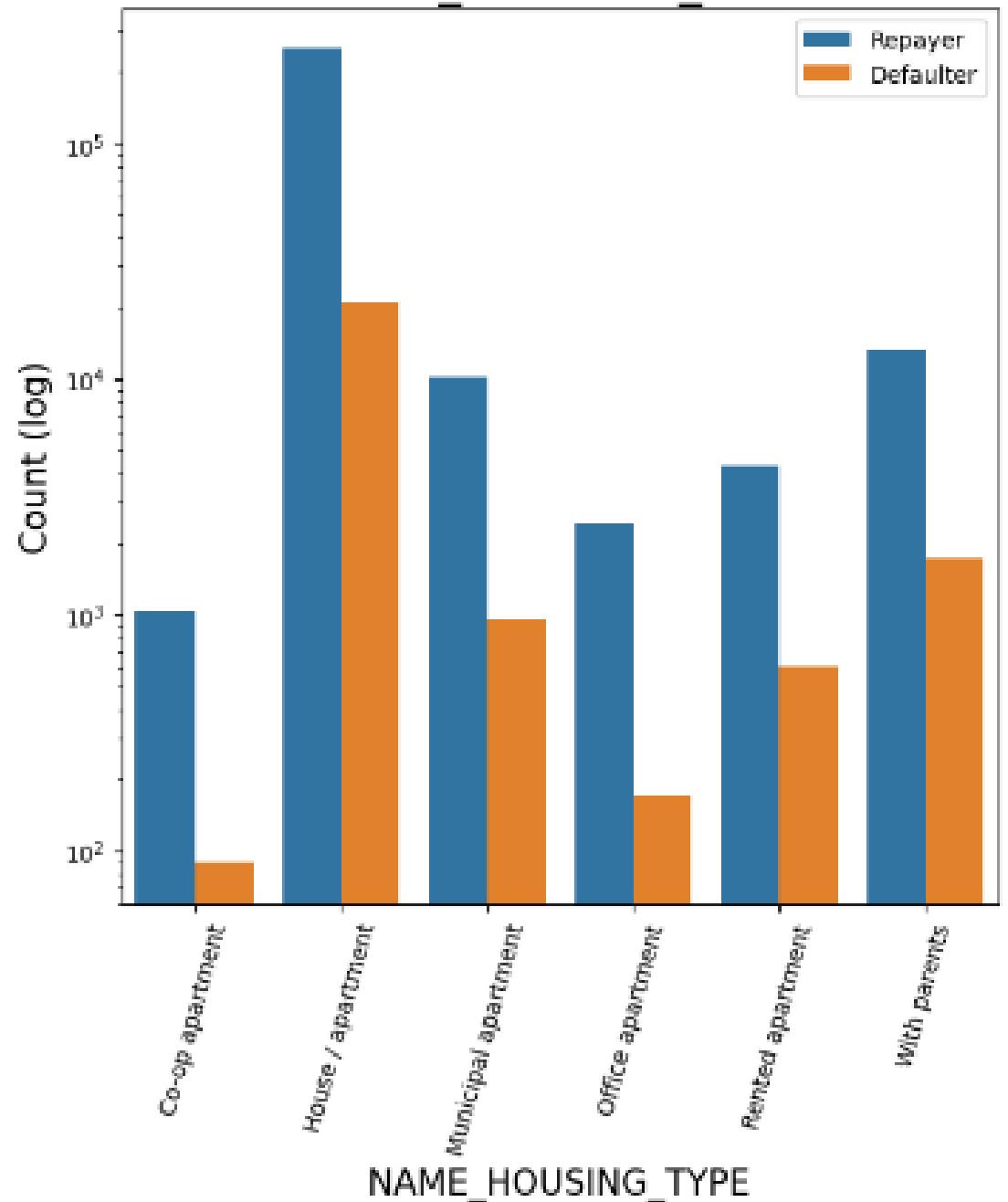
### FLAG\_OWN\_REALTY



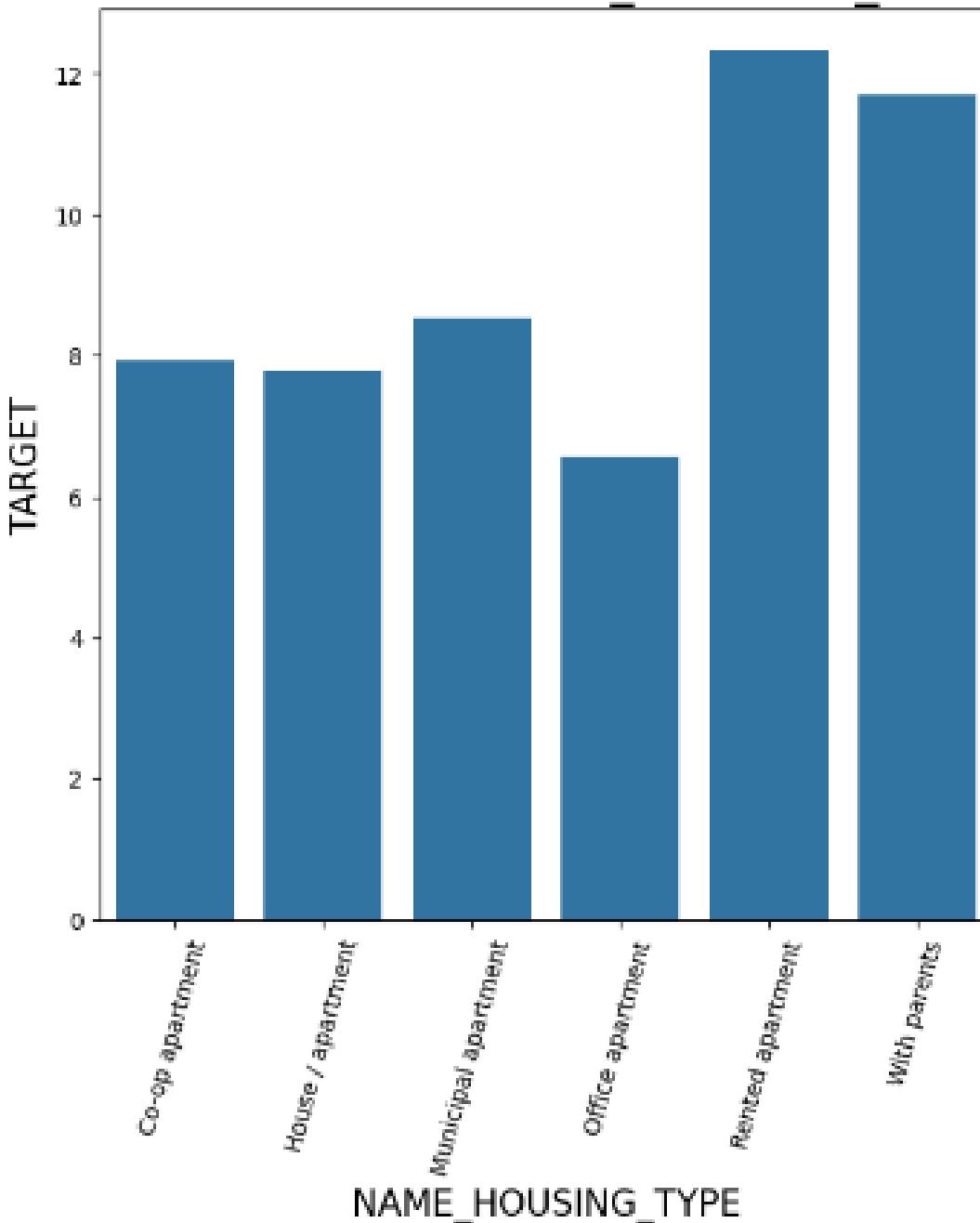
### Defaulters % in FLAG\_OWN\_REALTY



## NAME\_HOUSING\_TYPE

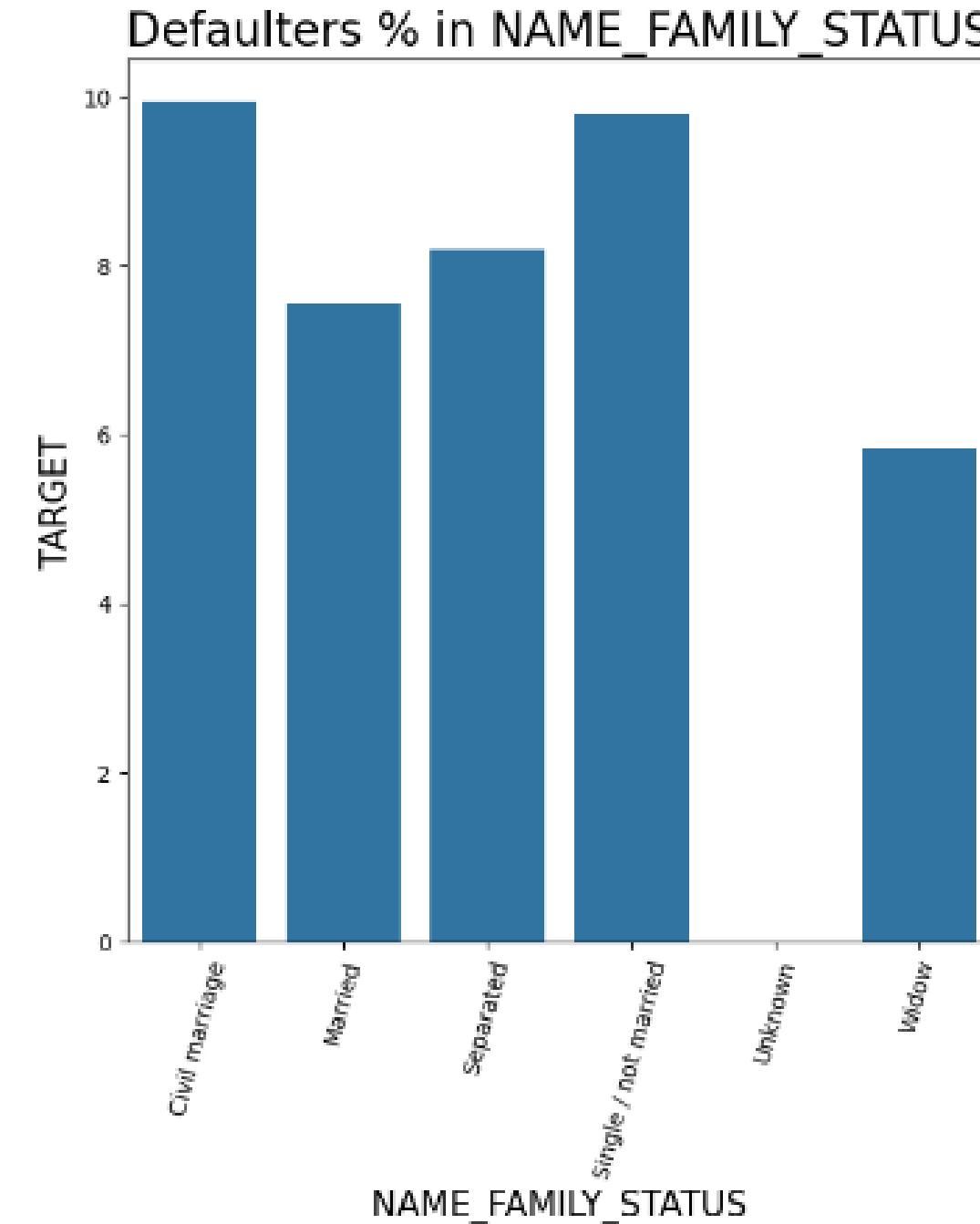
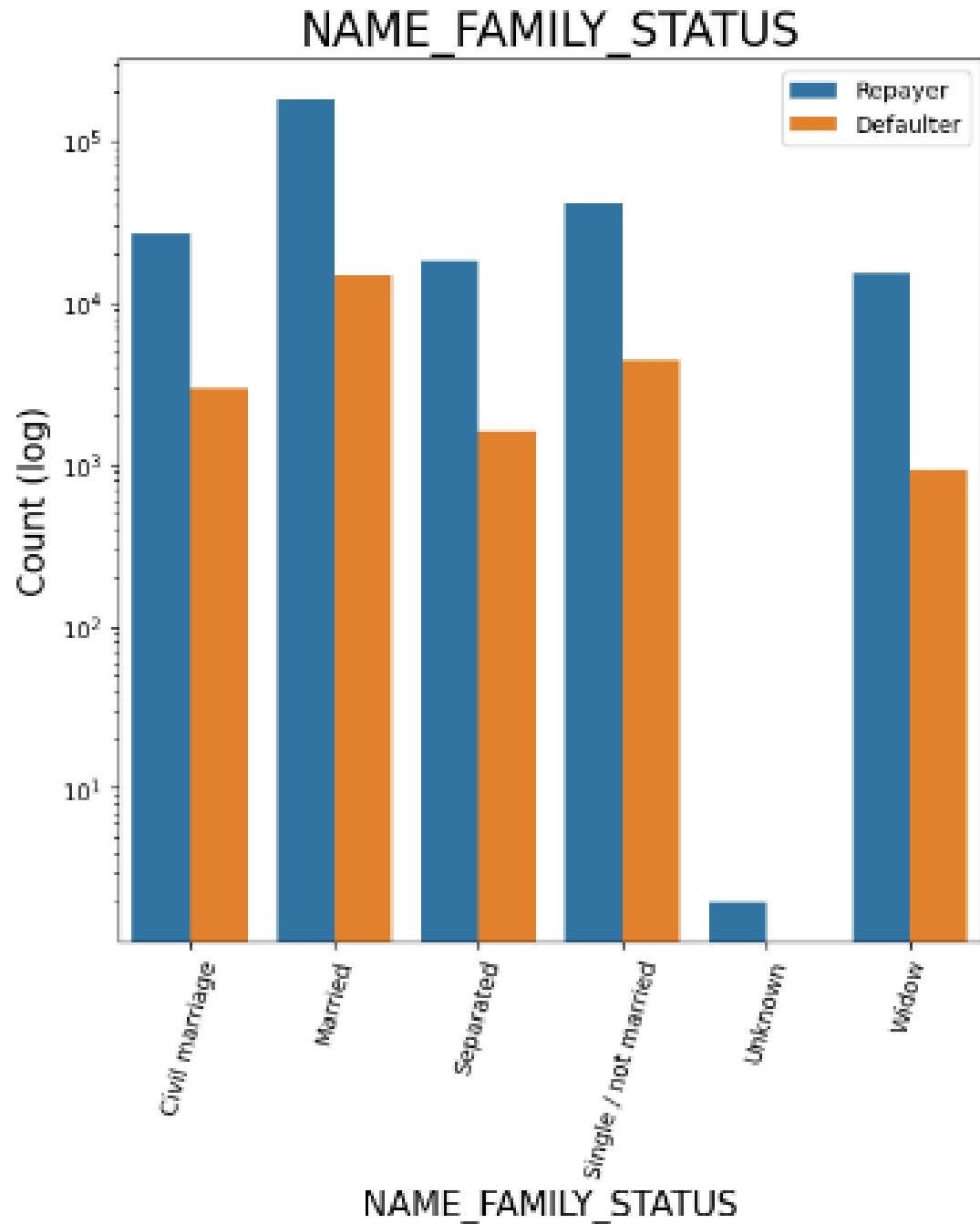


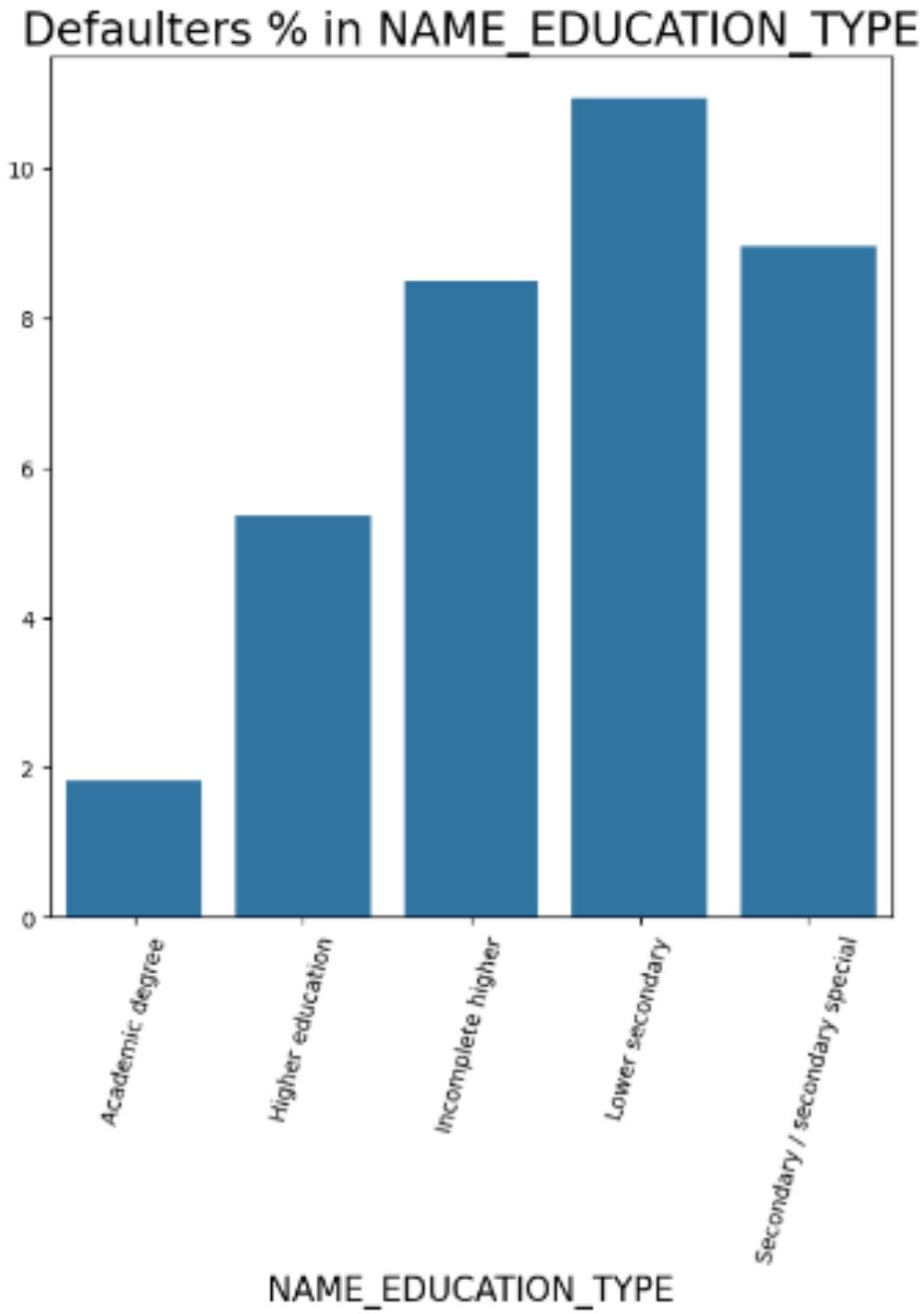
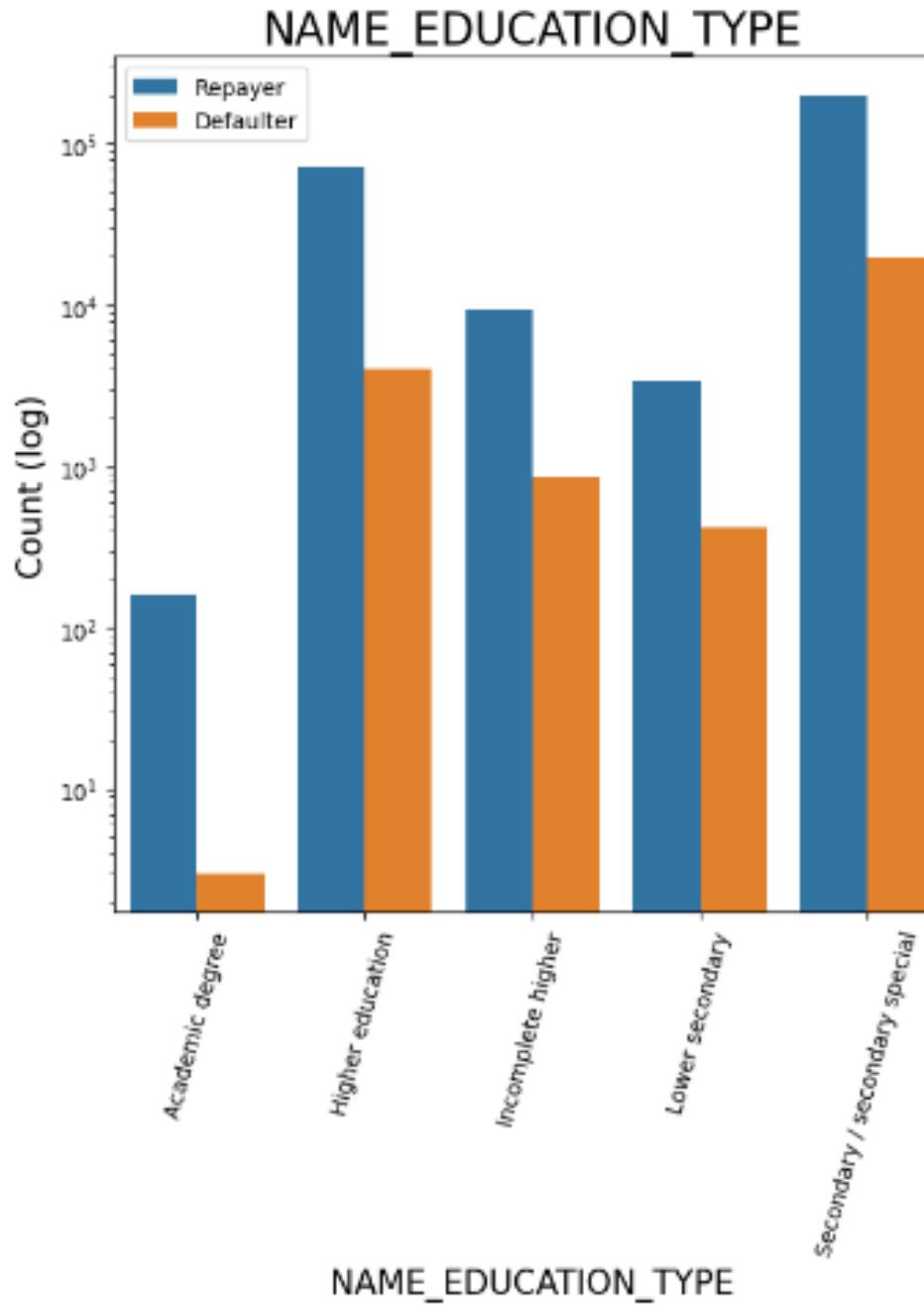
## Defaulters % in NAME\_HOUSING\_TYPE

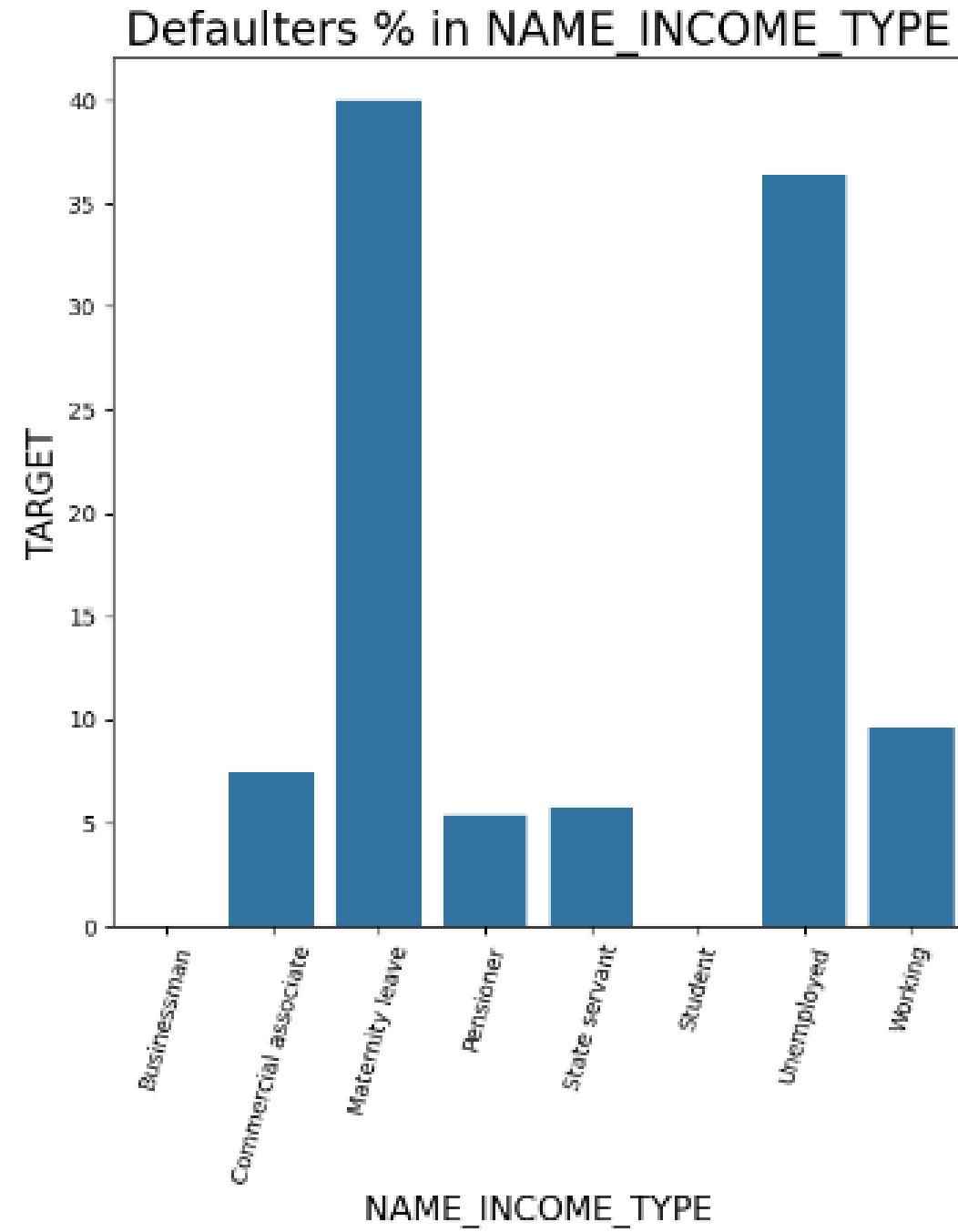
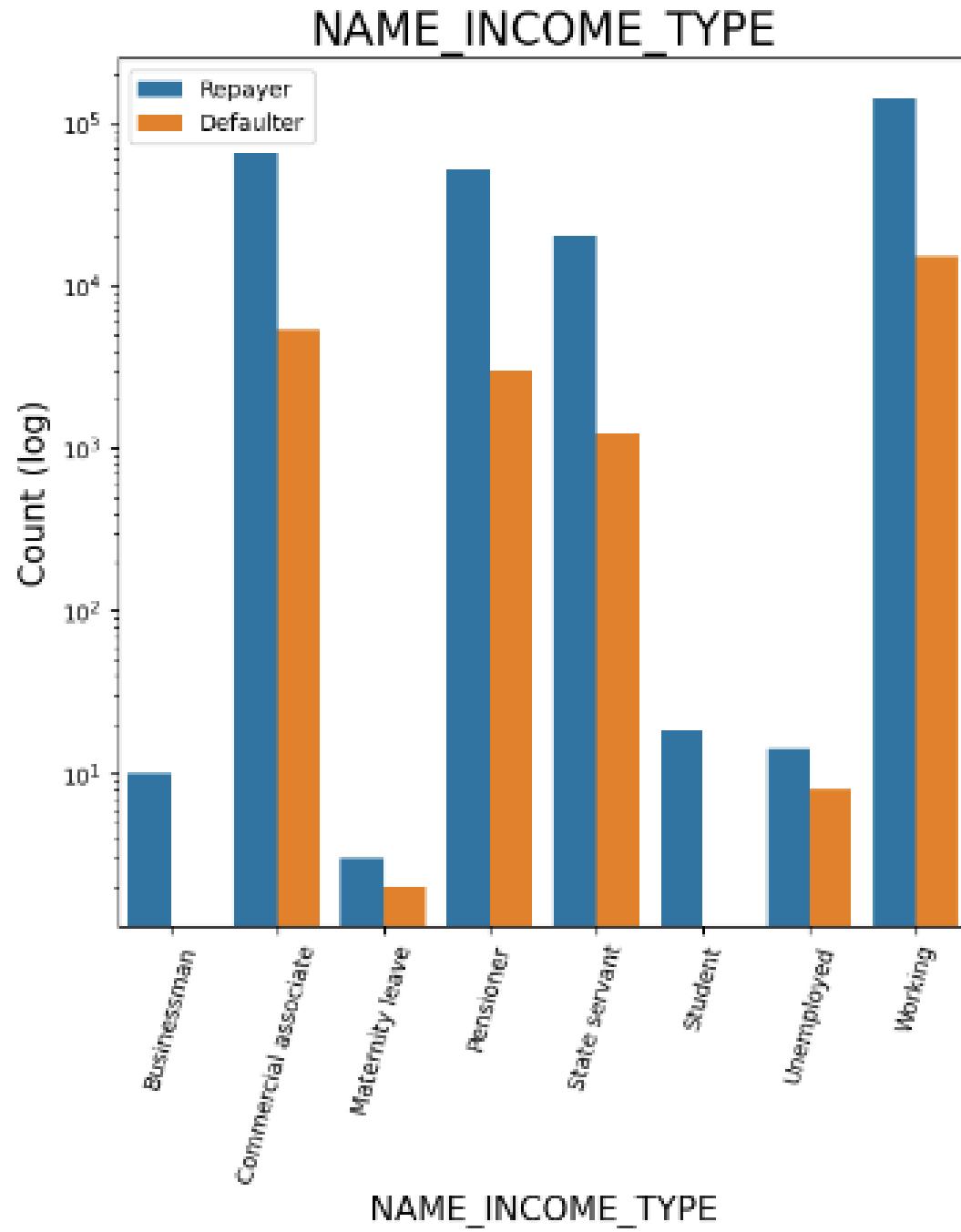


## Inferences

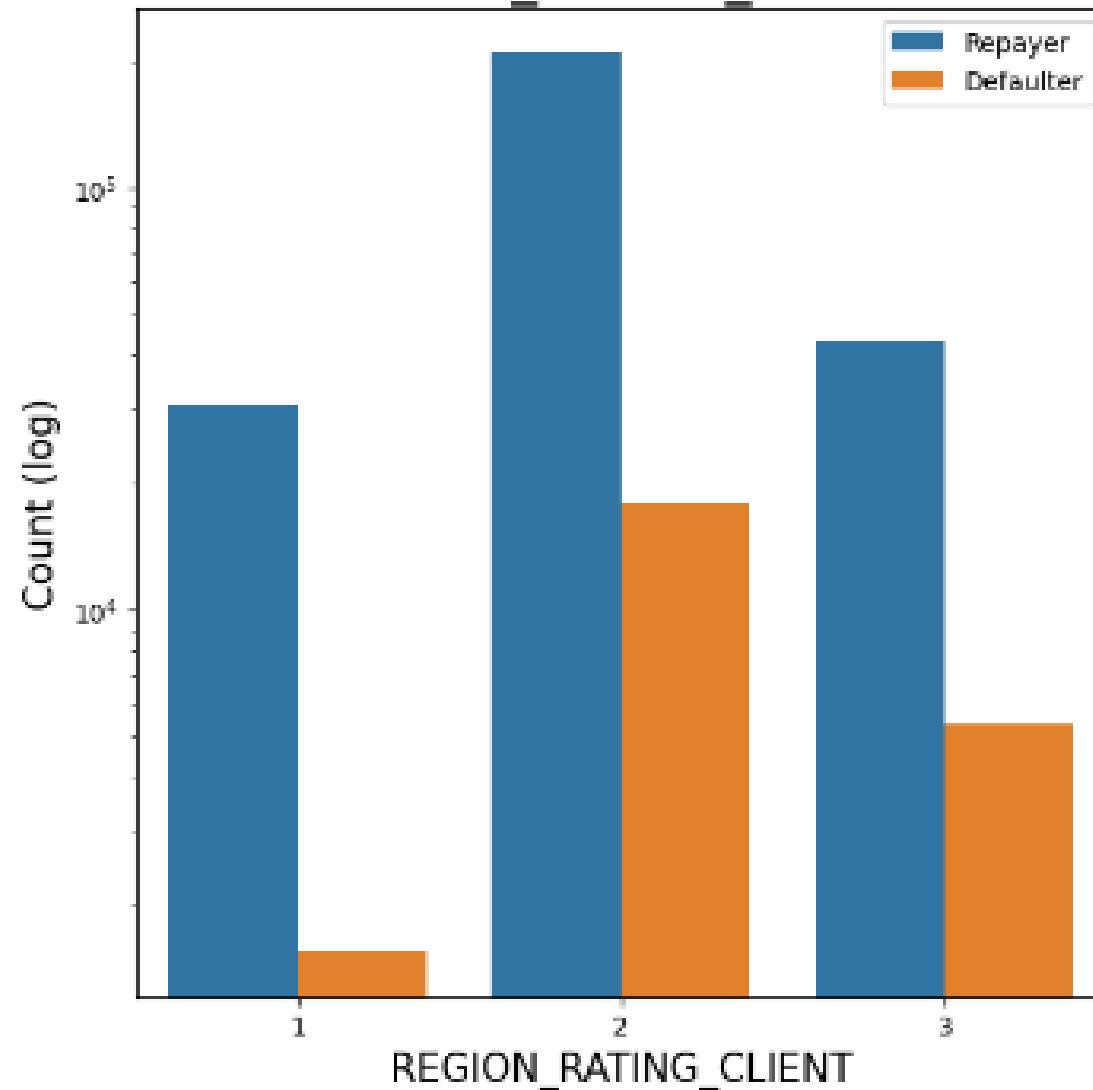
- Revolving loans, means loan for buying any goods, are just a small fraction (10%) from the total number of loans.
- Defaulters are in 8-9% Cash loan applicants and 5-6% Revolving loan applicant.
- The number of female loan applicants is almost double the number of male applicants.
- Based on the percentage of defaulted credits, chance of males' defaulters 10% is higher than female 7%.
- There are property owner are more applicants than don't own property.
- The defaulting rate is same for both, own and don't own property (8%), which means correlation between owning a property and defaulting the loan.
- Majority of applicants are live in House/apartment and applicants living in office apartments have lowest default rate
- The applicants who are living with parents (11.5%) and living in rented apartments(>12%) have higher probability of defaulting



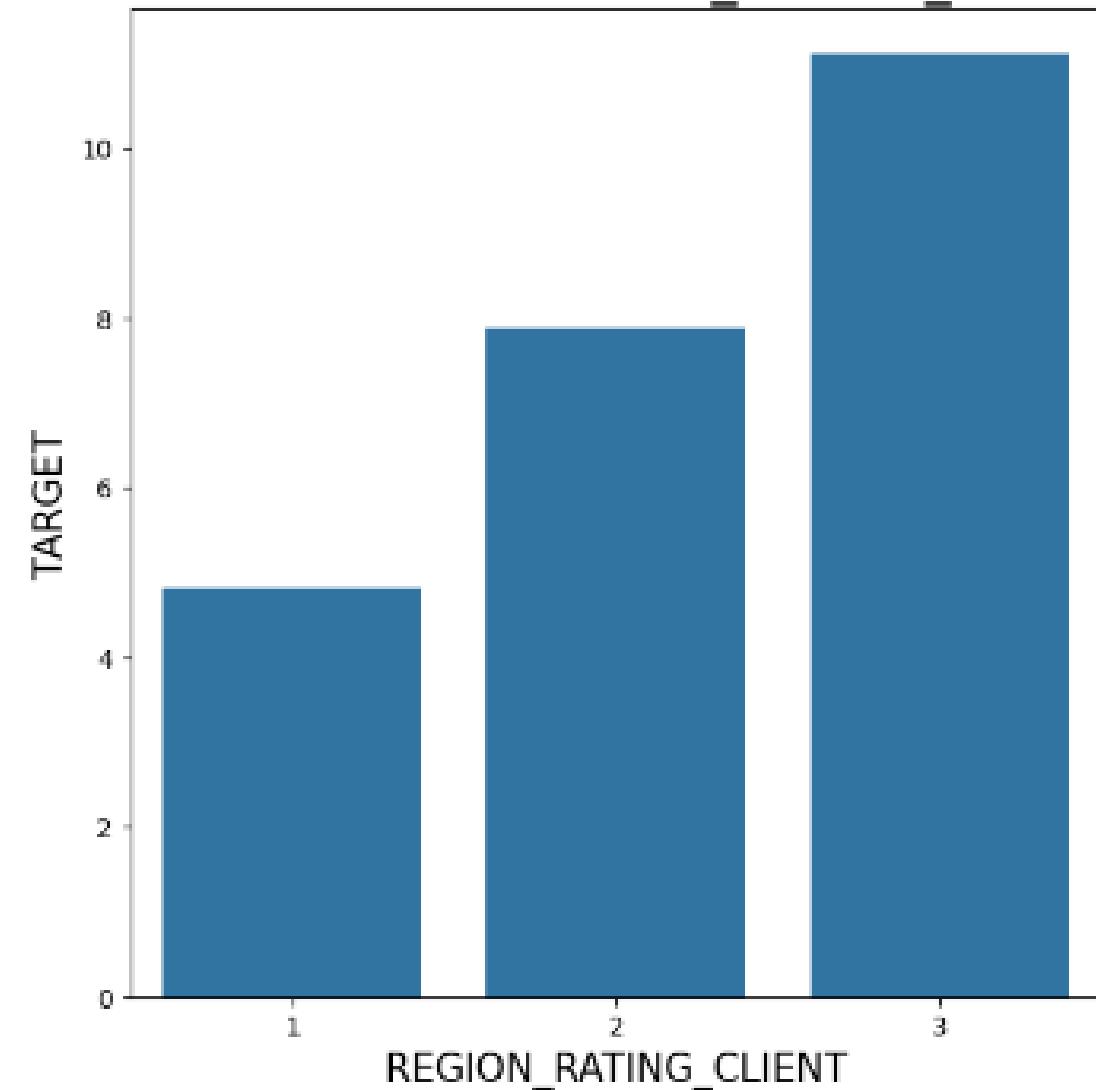




## REGION\_RATING\_CLIENT

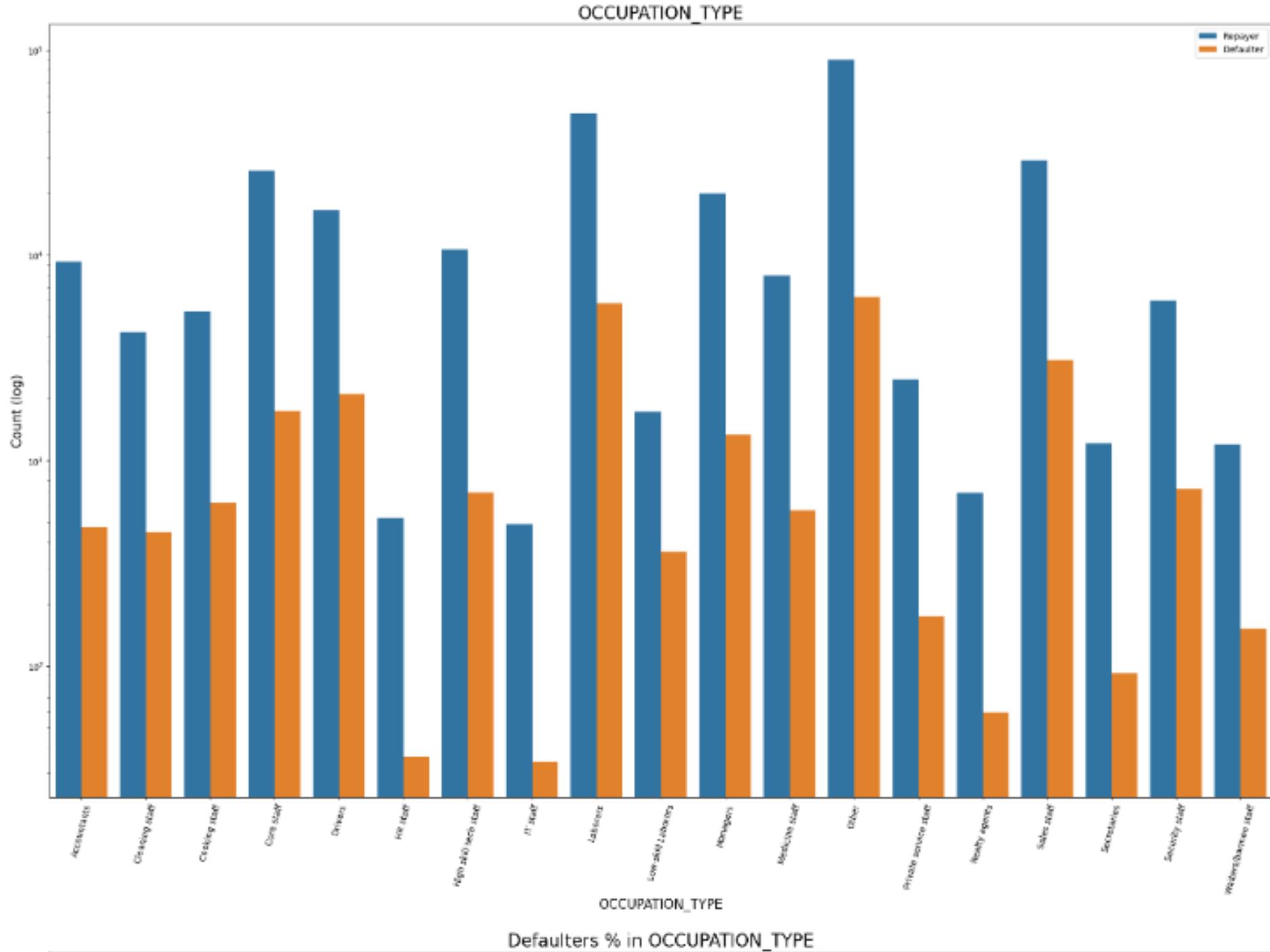


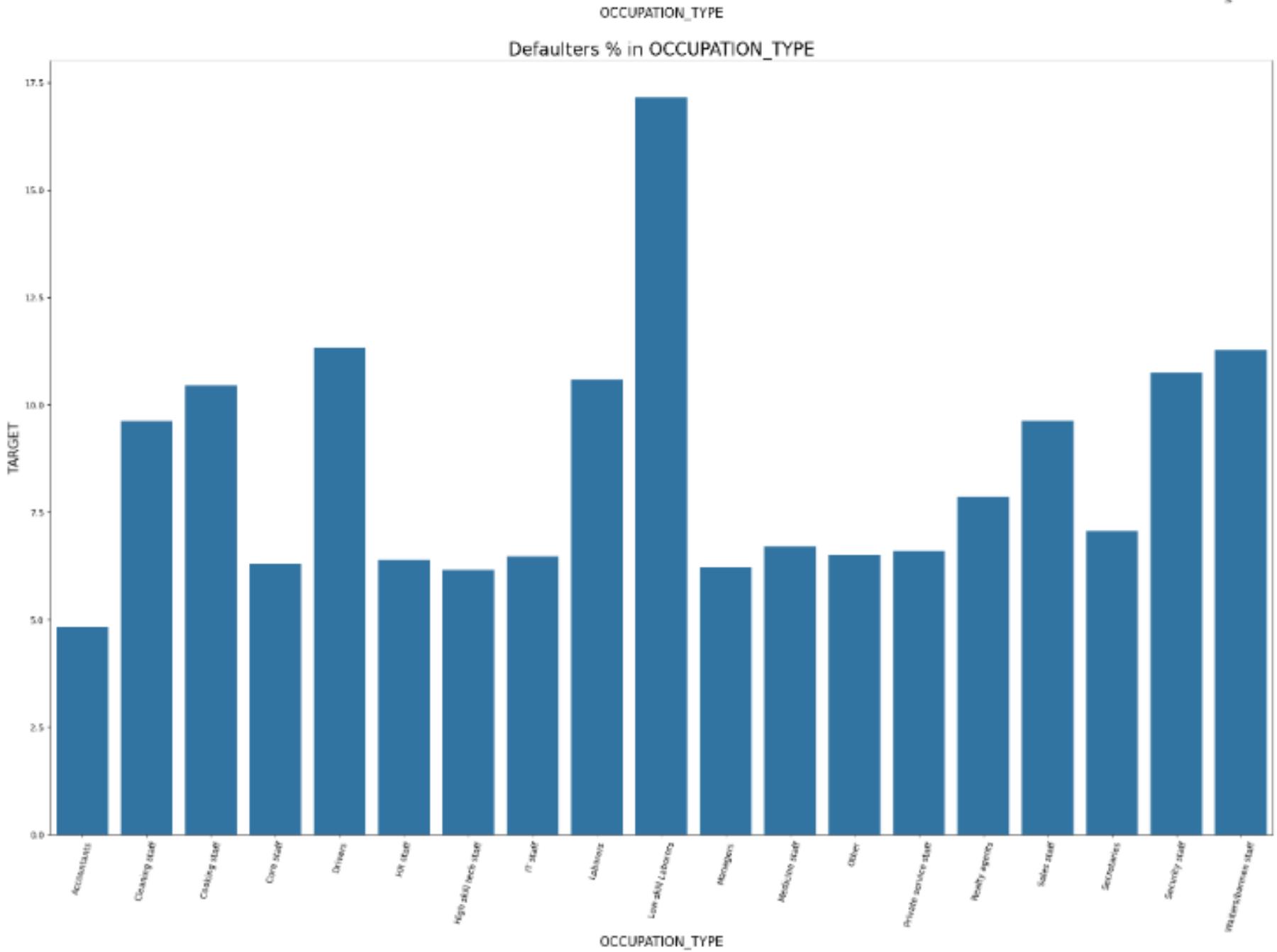
## Defaulters % in REGION\_RATING\_CLIENT

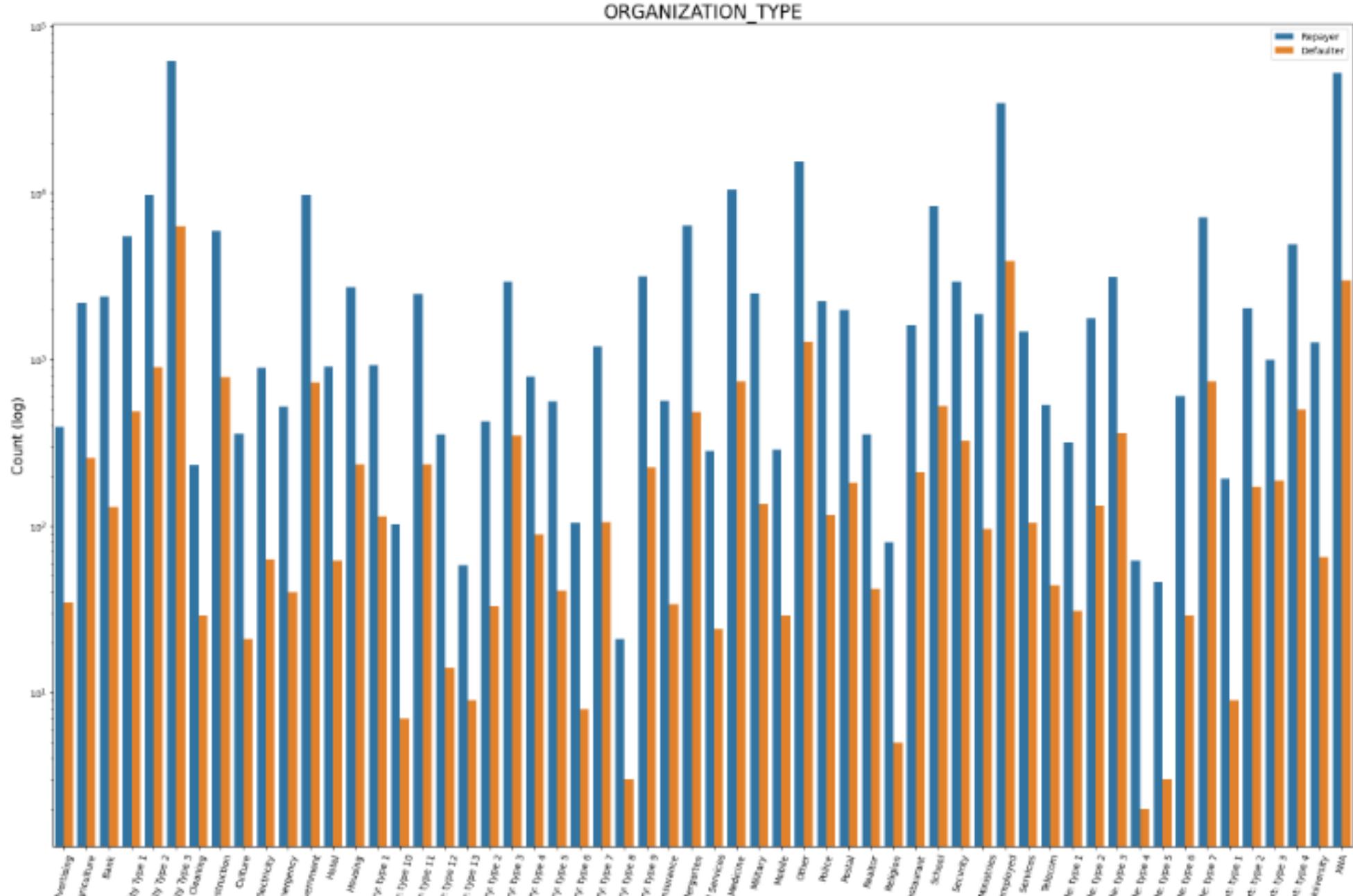


## Inferences

- Married have applied more then followed by Single/not married and civil marriage, applicants.
- In Percentage of defaulters, Civil marriage has the highest percent around (10%) and widow has the lowest around 6% (exception being Unknown).
- Majority of applicants have Secondary/secondary special education, followed by applicants with Higher education.
- Lower secondary category have highest rate of defaulting around 11%. People with Academic degree are least likely to default and very few applicants have an academic degree.
- Most of the applicants are living in Region with Rating 2 place. Region Rating 3 has the highest default rate (11%)
- Applicant living in Region\_Rating 1 are safer for approving loanshas the lowest probability of defaulting.



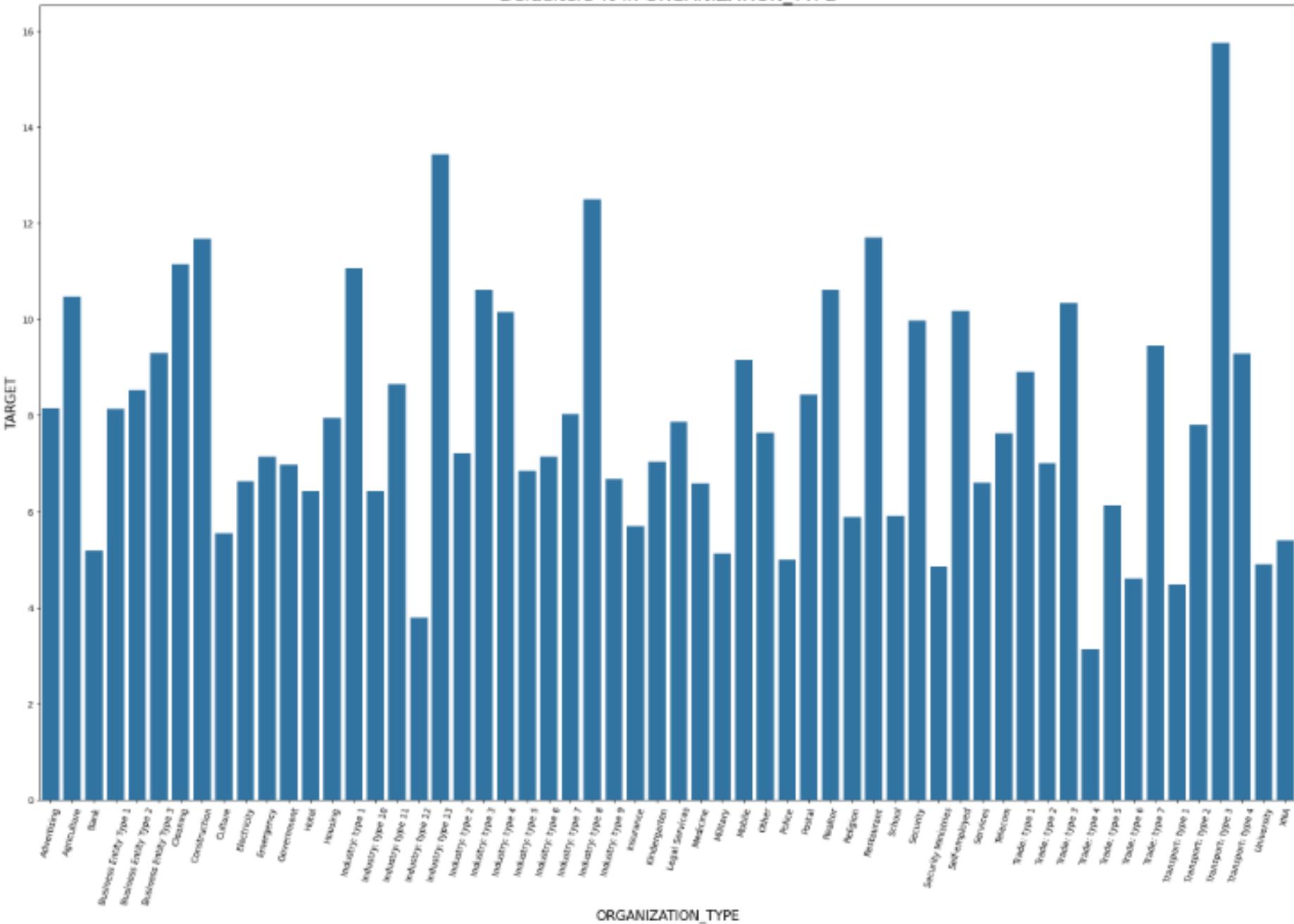




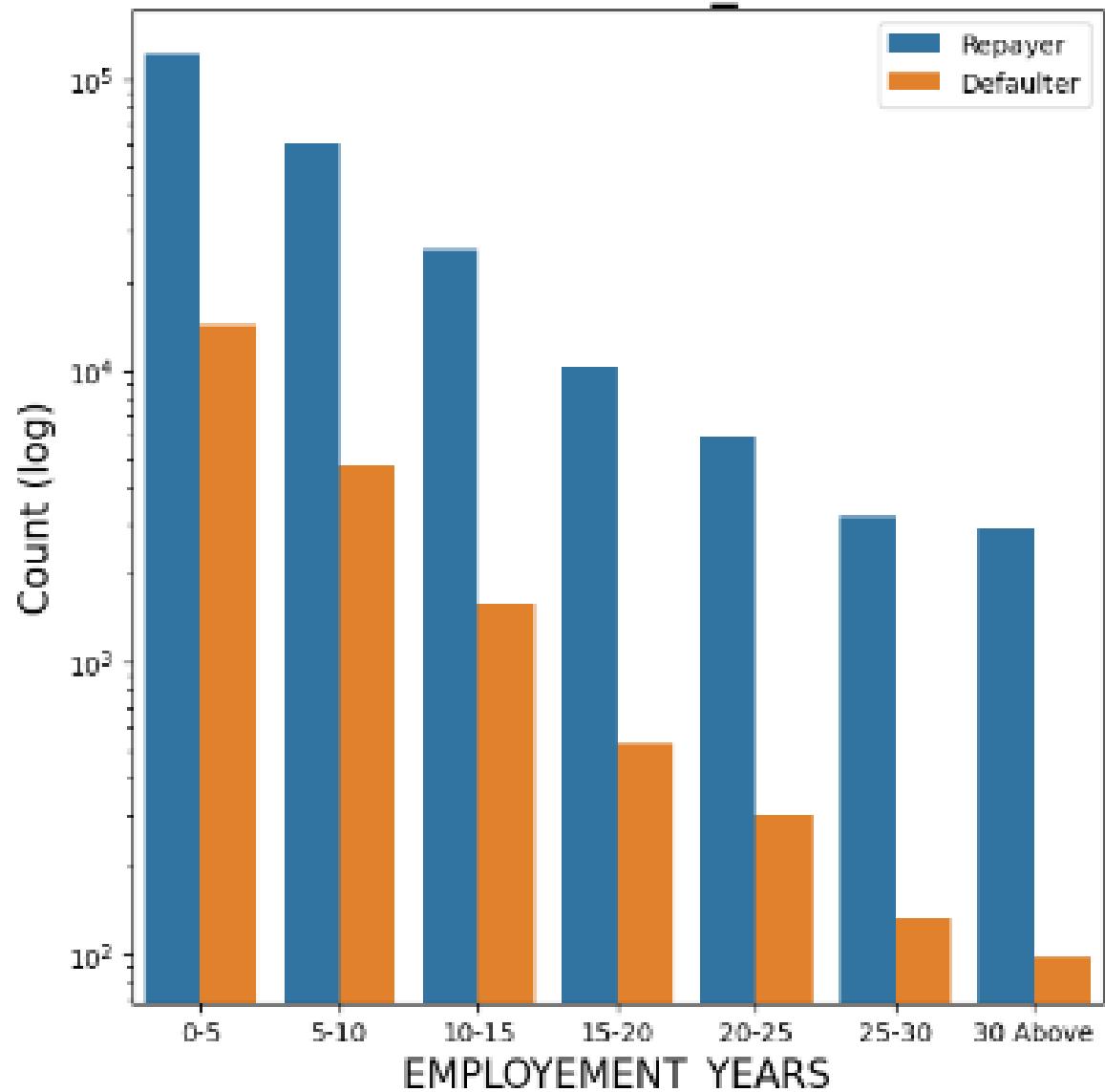
Defaulters % in ORGANIZATION\_TYPE



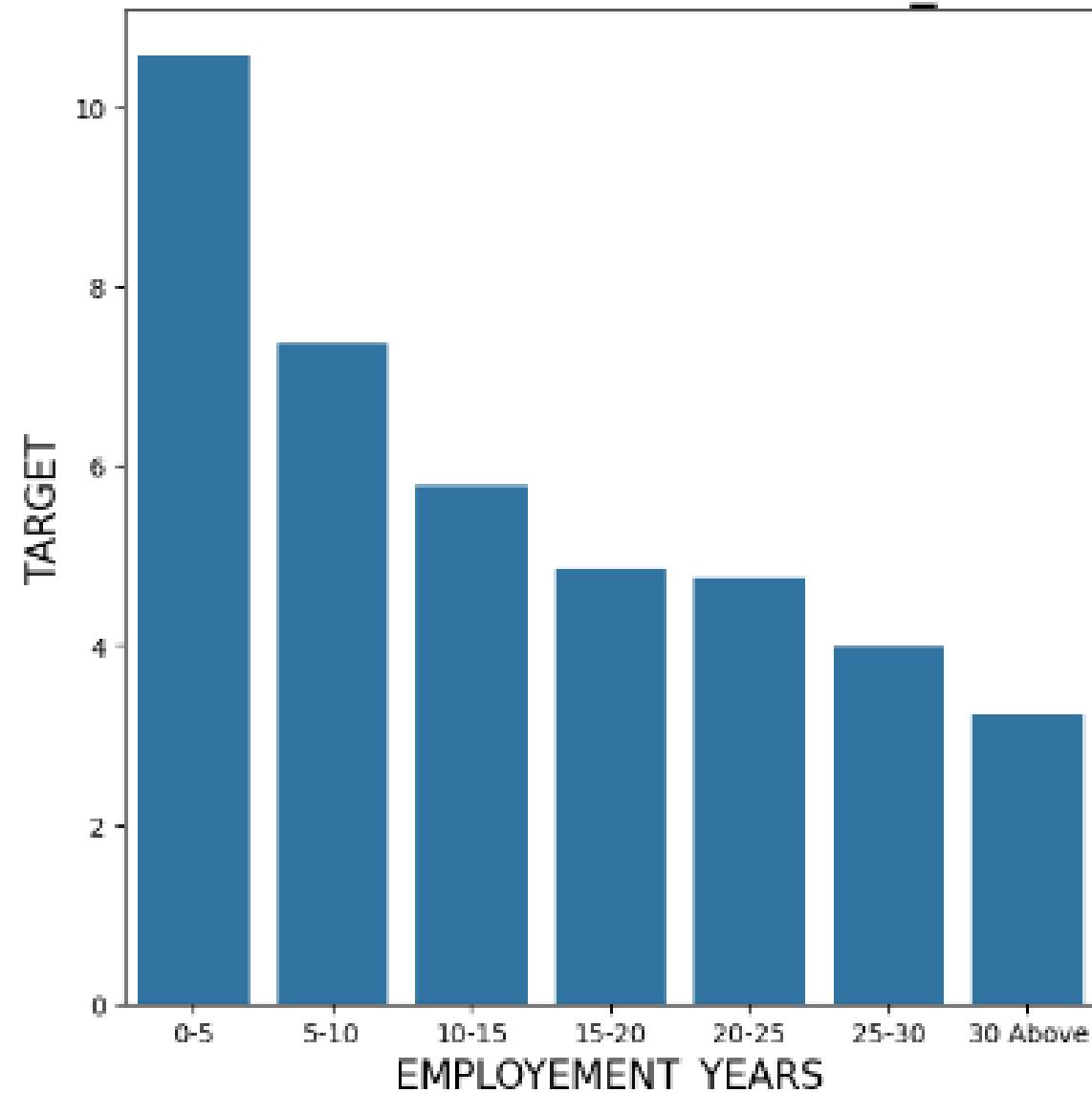
## Defaulters % in ORGANIZATION\_TYPE

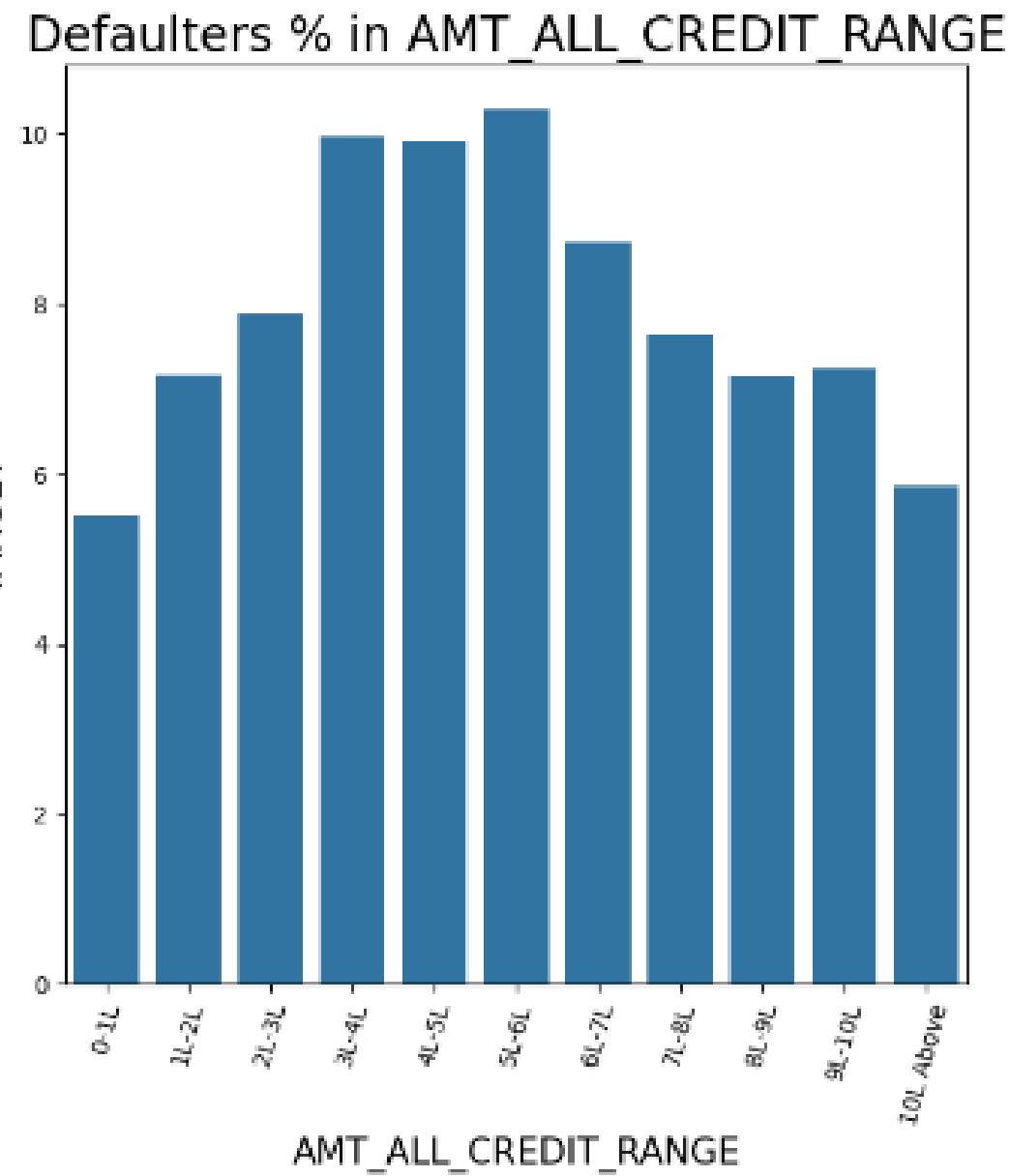
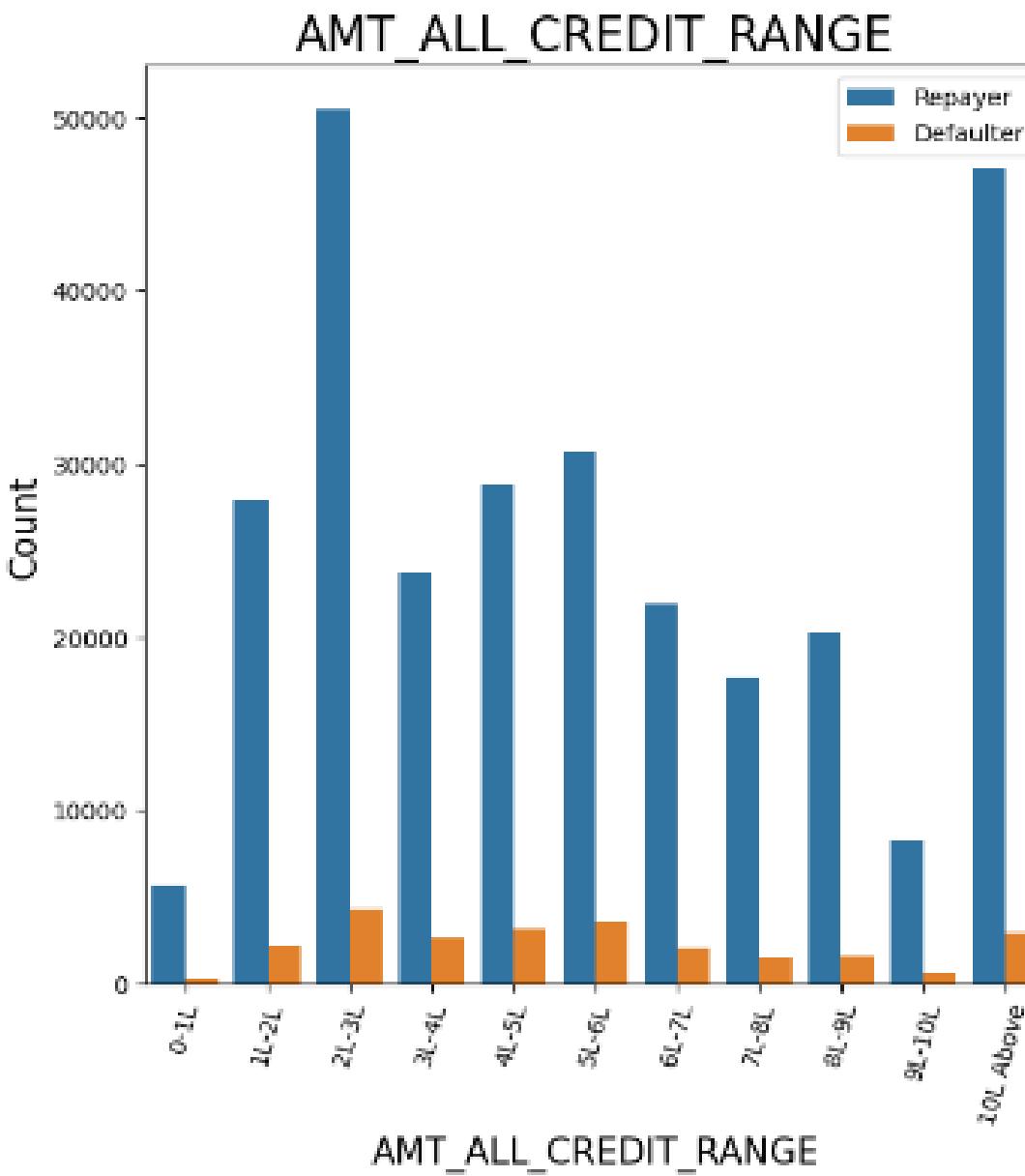


## EMPLOYEMENT\_YEARS



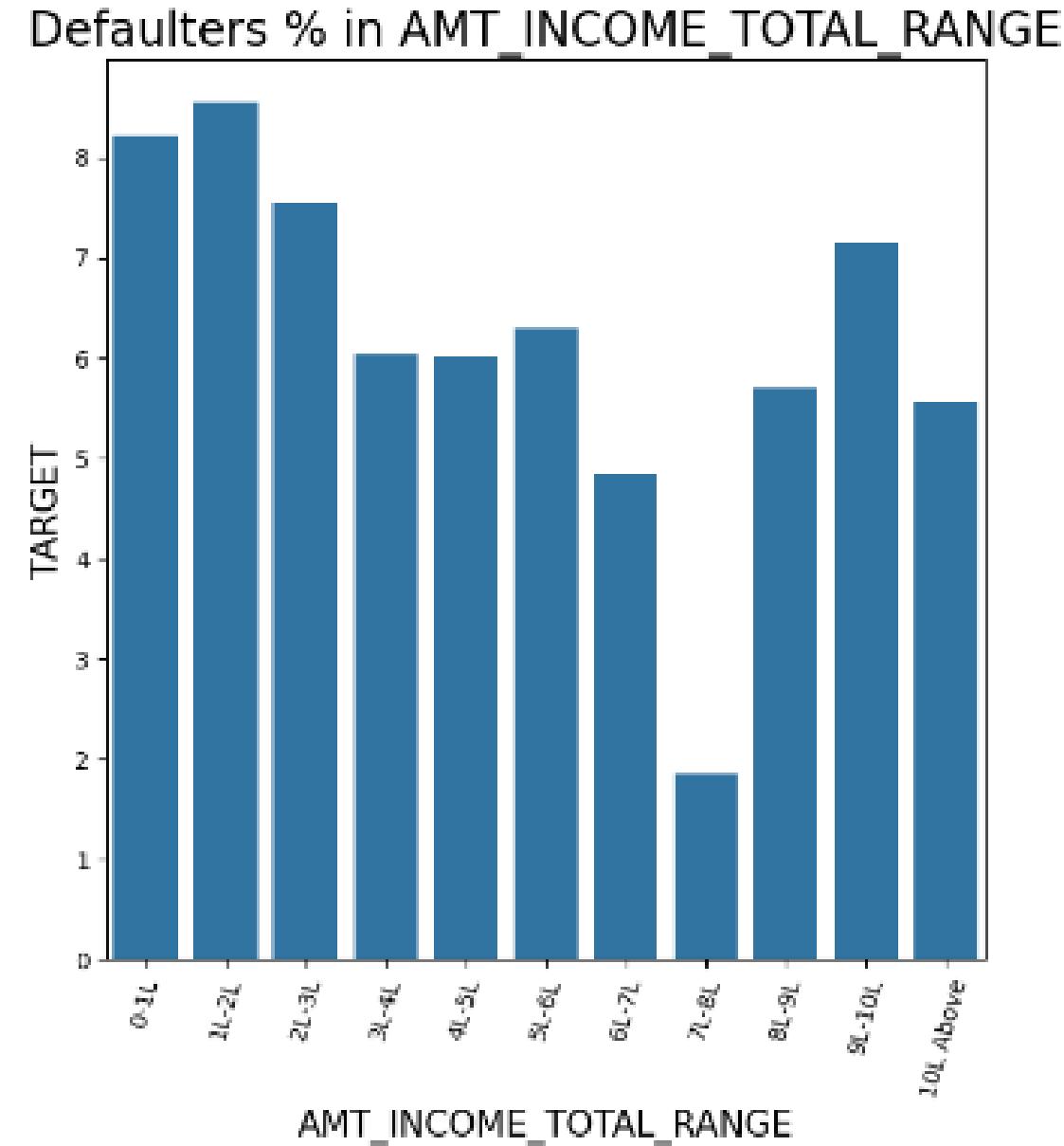
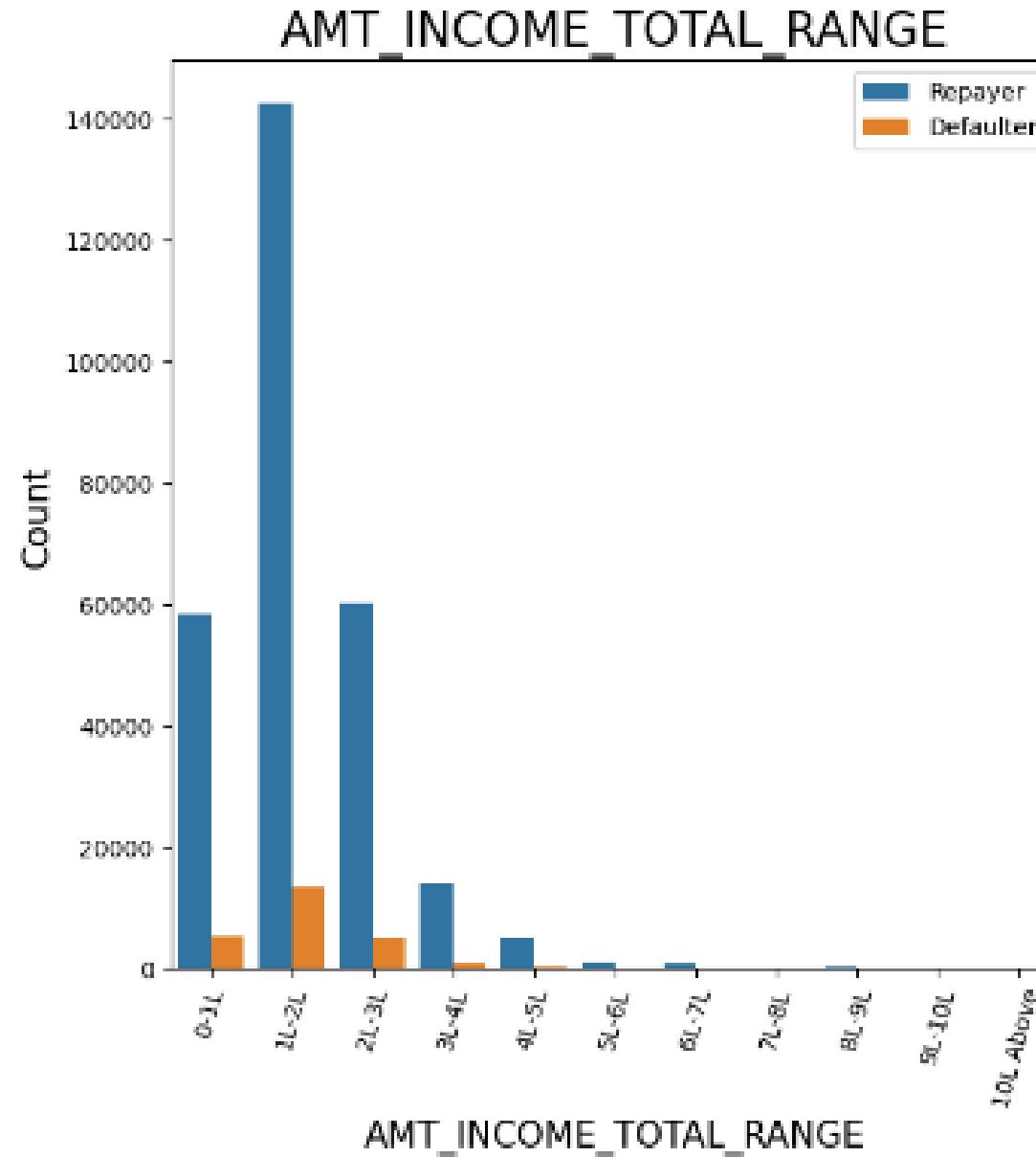
## Defaulters % in EMPLOYEMENT\_YEARS



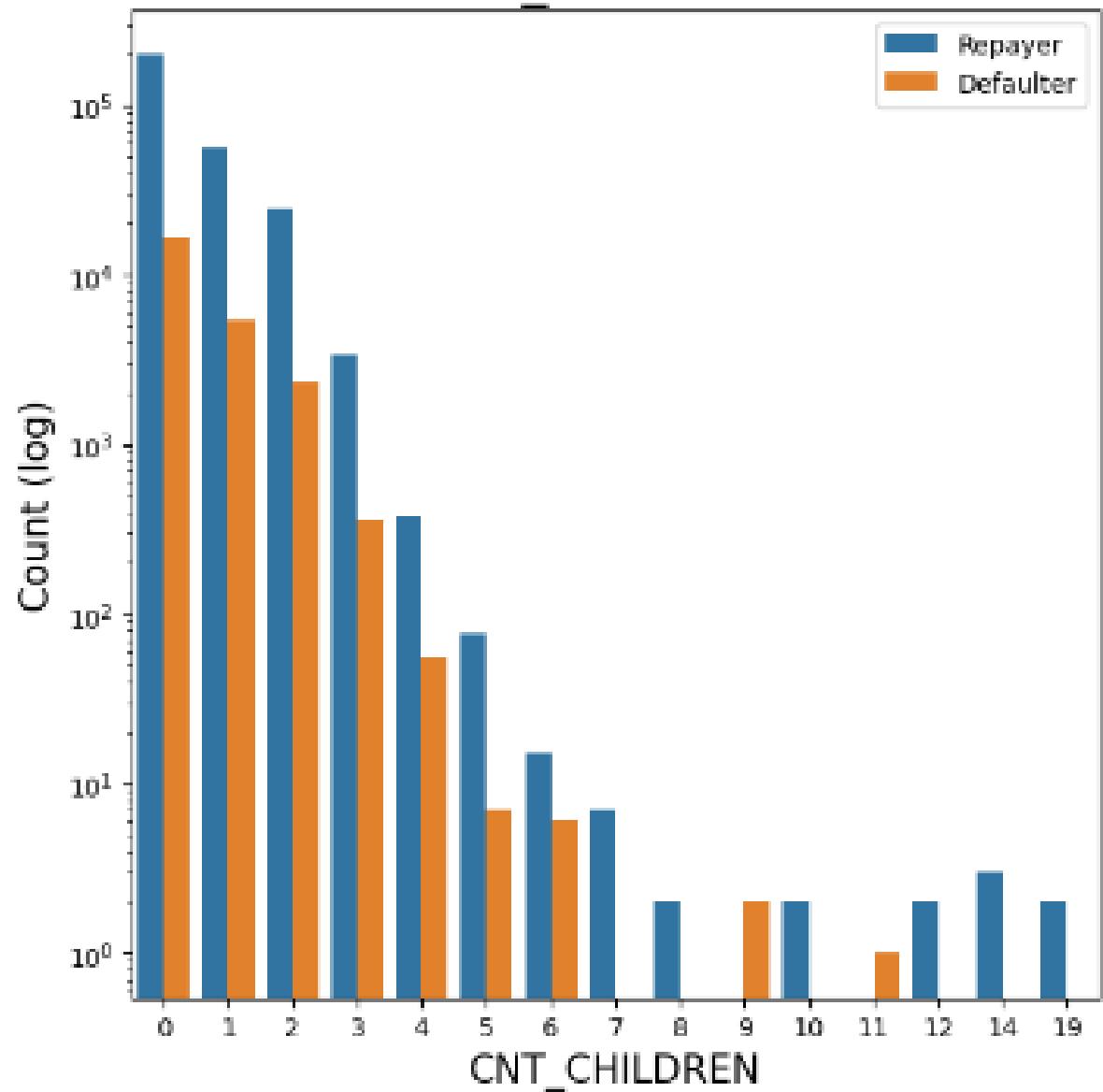


## Inferences

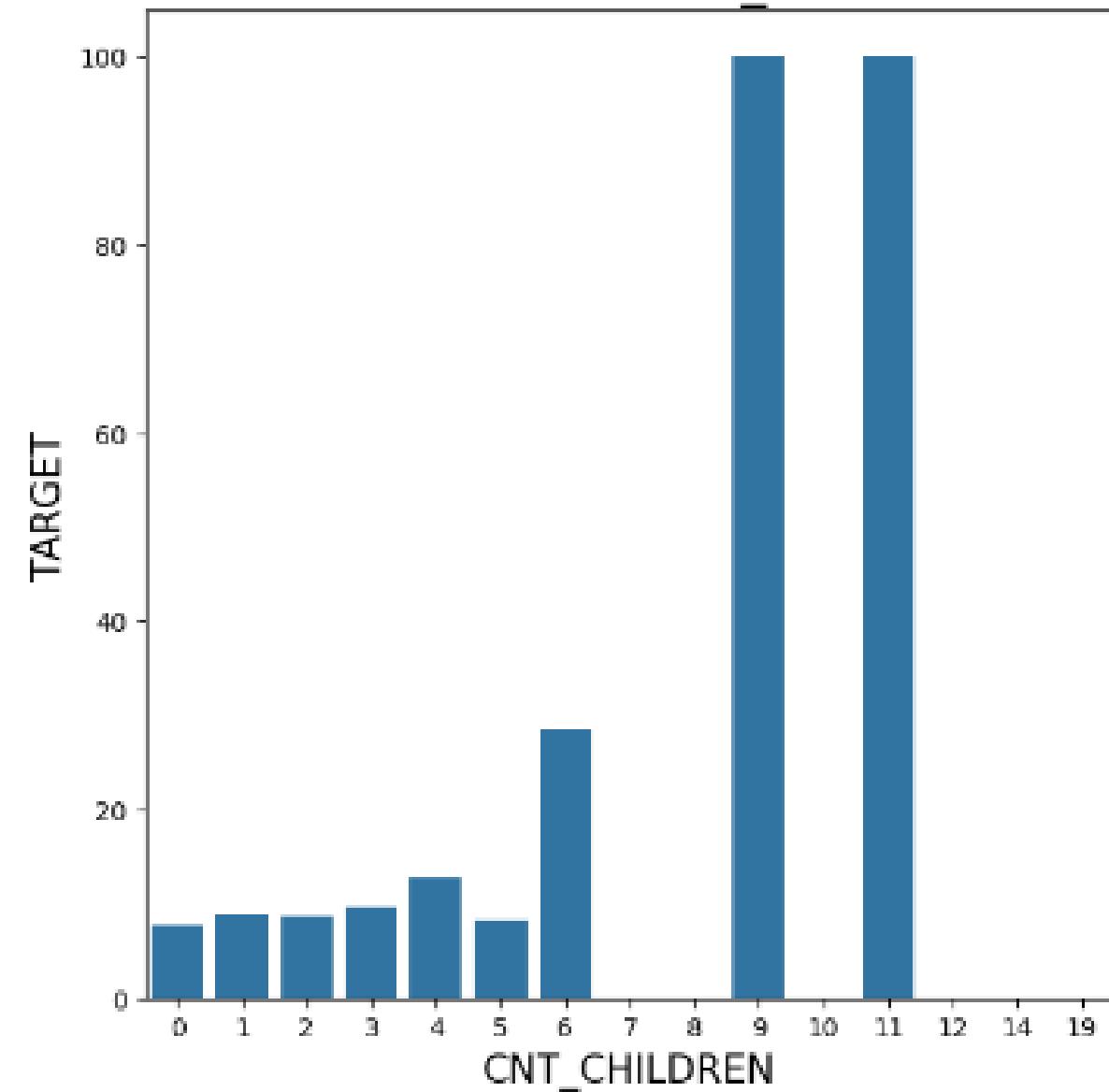
- Category with highest percent of defaulters are Low-skill Laborers (above 17%), followed by Drivers and Waiters/barmen staff, Security staff, Laborers and Cooking staff.
- IT staff are less likely to apply for Loan, and most of the loans are taken by Laborers, followed by Sales staff.
- Self employed people have relative high defaulting rate, to be safer side loan disbursement should be avoided or provide loan with higher interest rate to mitigate the risk of defaulting.
- Organizations with highest percent of defaulters are Transport: type 3 (16%), Industry: type 13 (13.5%), Industry: type 8 (12.5%) and Restaurant (less than 12%).
- With people having 40+ year experience have less than 1% default rate. With increase of employment year, defaulting rate is gradually decreasing.
- Majority of the applicants having working experience between 0-5 years are defaulters. The defaulting rating of this group is also the highest which is around 10%.

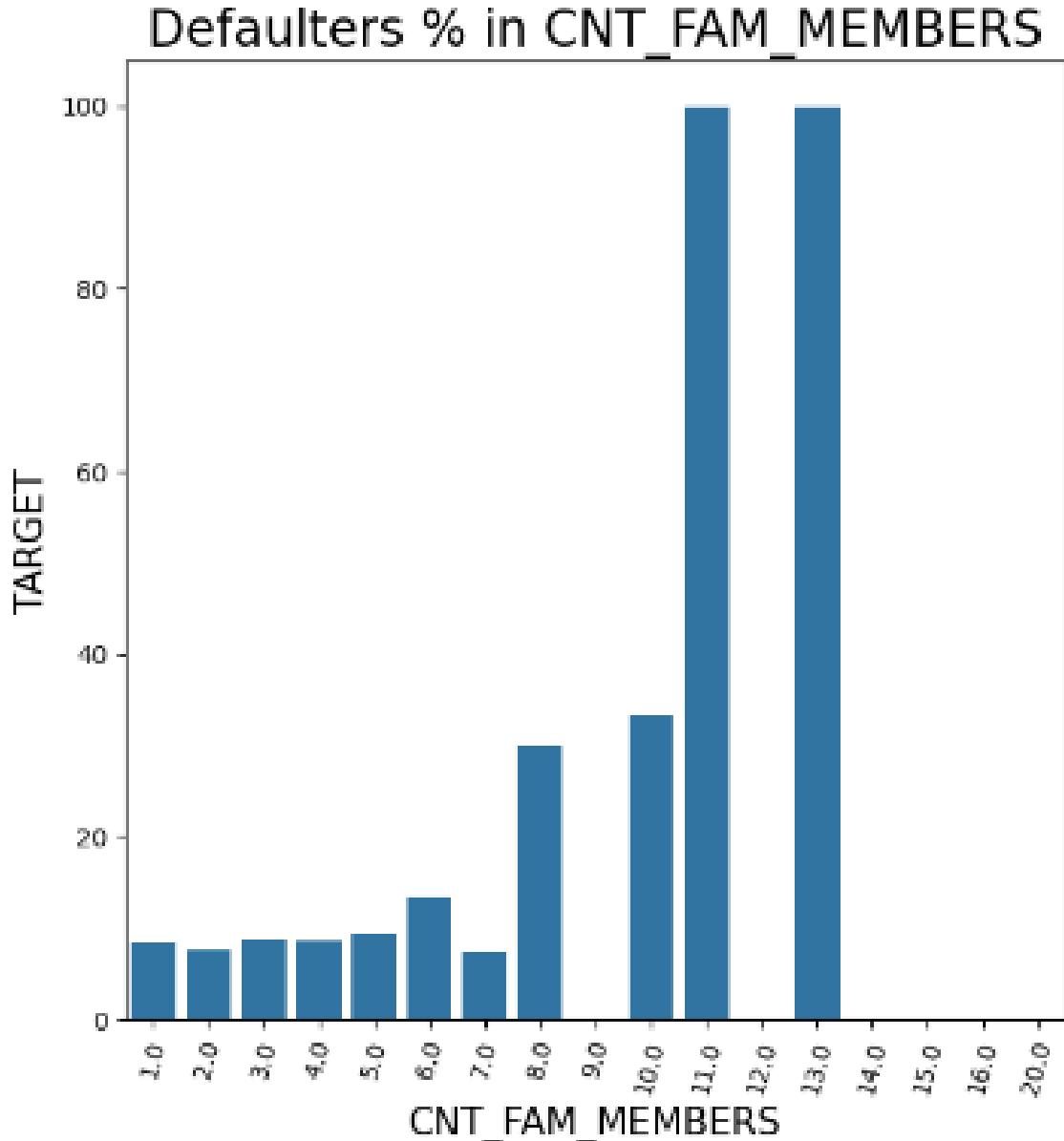
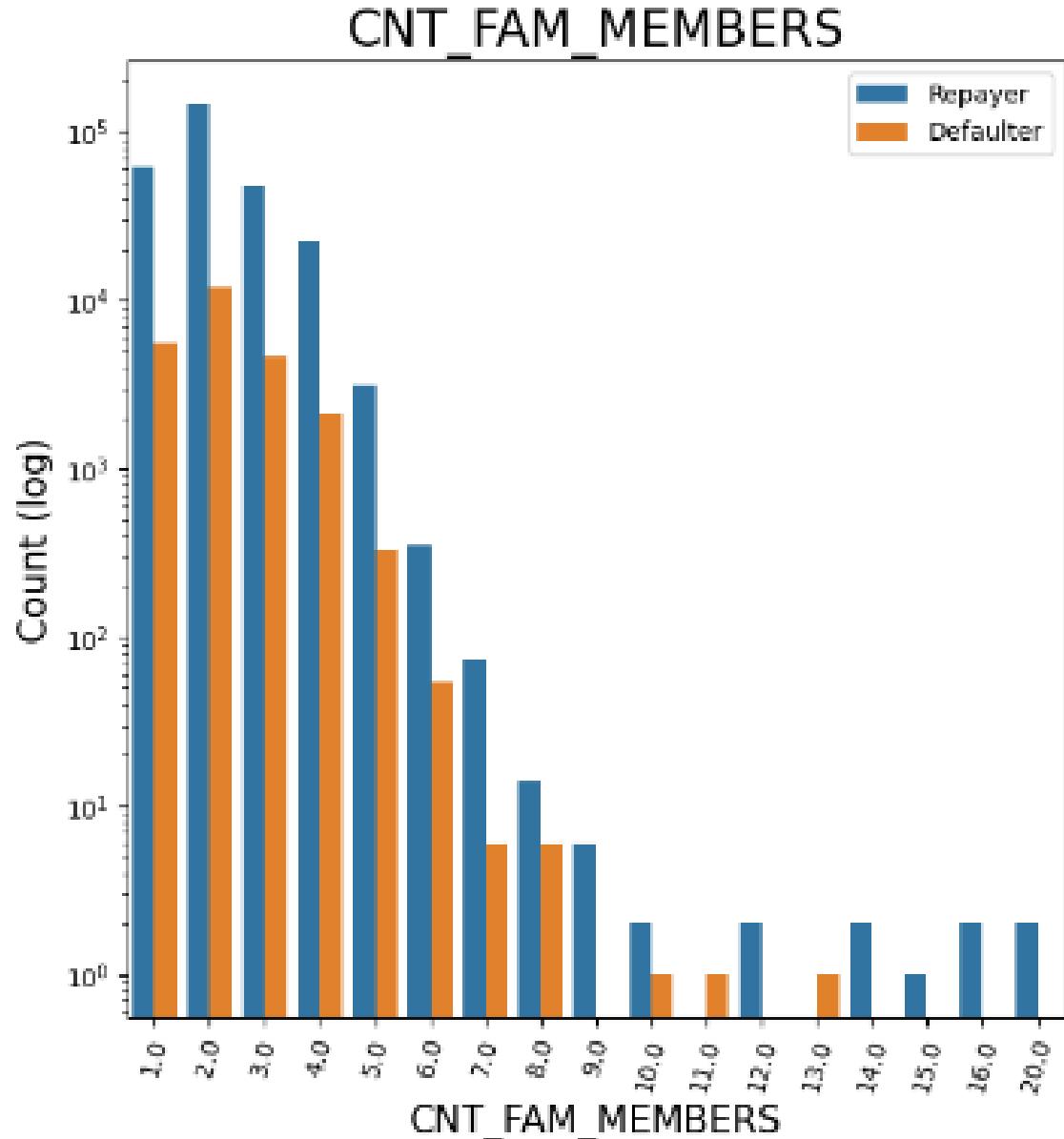


## CNT\_CHILDREN



## Defaulters % in CNT\_CHILDREN

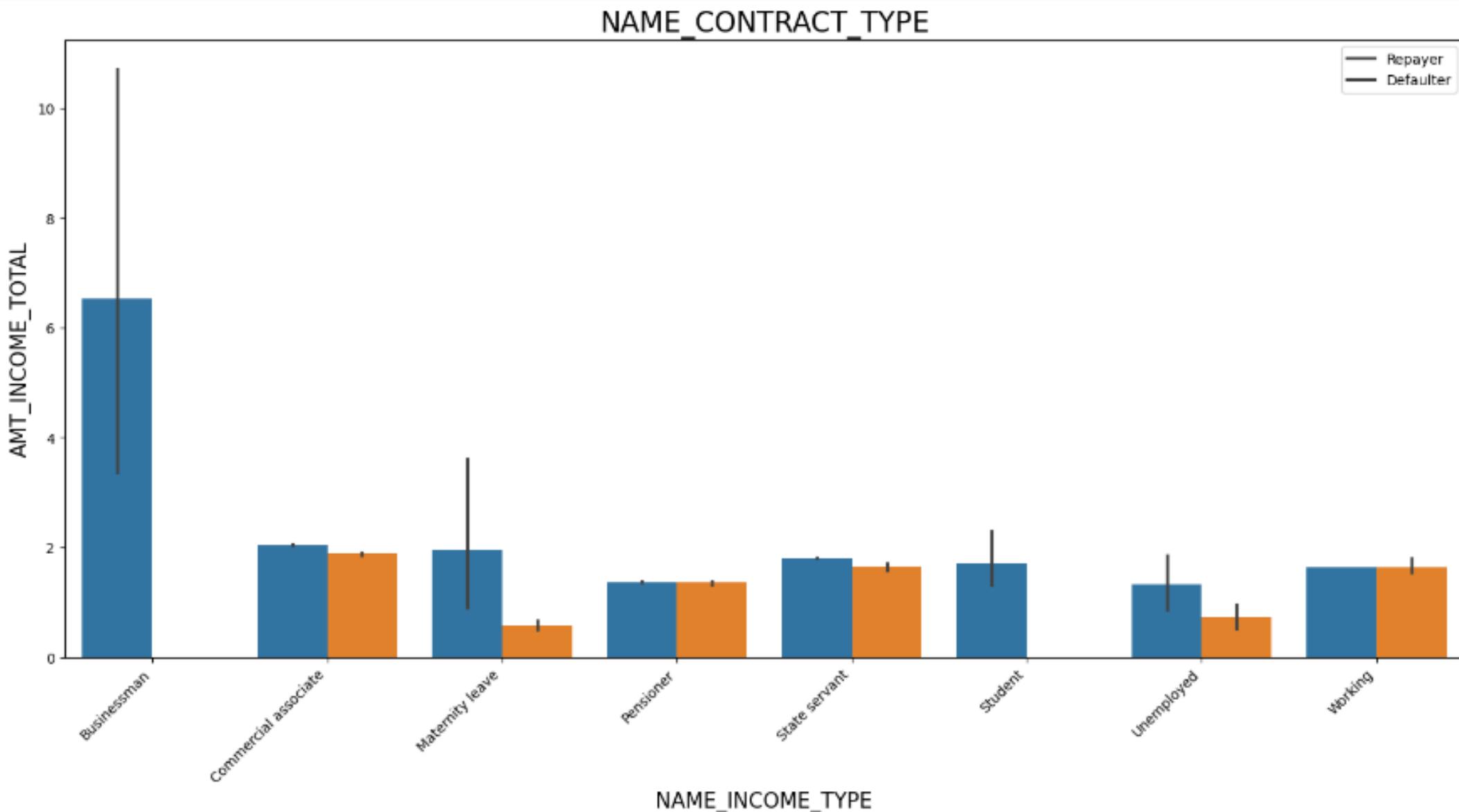




## Inferences

- Applicants who get loan for 3L-6L have most number of defaulters than other loan range.
- Applicants with high number have loan in range of 2L-3L.
- Majority of the applications have Income total less than 3L.
- Application with Income less than 3L, has high probability of defaulting, and applicant with Income 7L-8L are less likely to default.
- Applicants who have more than 4 children has a very high default rate with child count 9 and 11 showing 100% default rate
- Most of the applicants do not have children, and very few clients have more than 3 children.
- Family member follows the same trend as children where having more family members increases the risk of defaulting.

# Categorical Bivariate / Multivariate Analysis

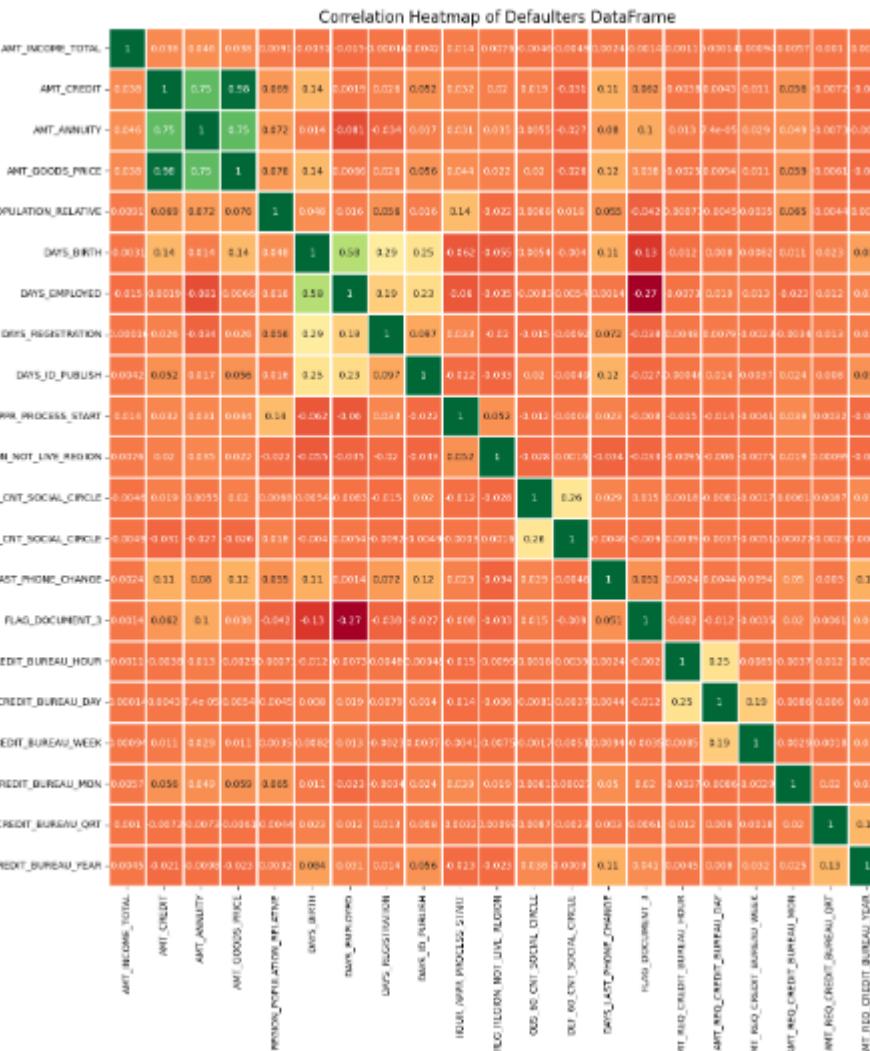
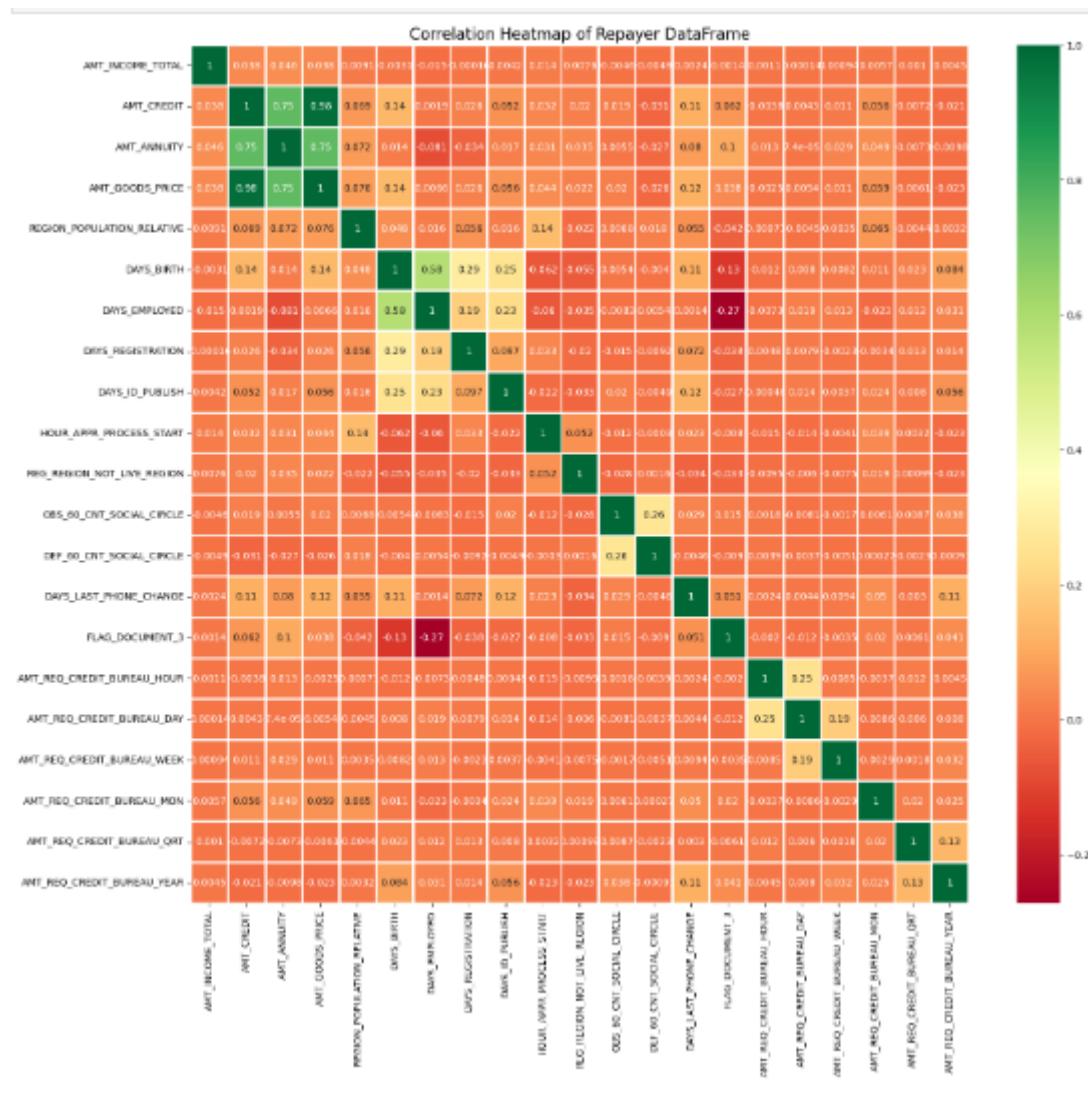


## Inferences

- Businessman income is the highest and the estimated range with default 95% confidence level and it seems to indicate that the income of a Businessman could be in the range of slightly close to 4L and slightly more than 10L.

# Numeric Multivariate Analysis

This is based on Target values 0 and 1 for correlation and other analysis.

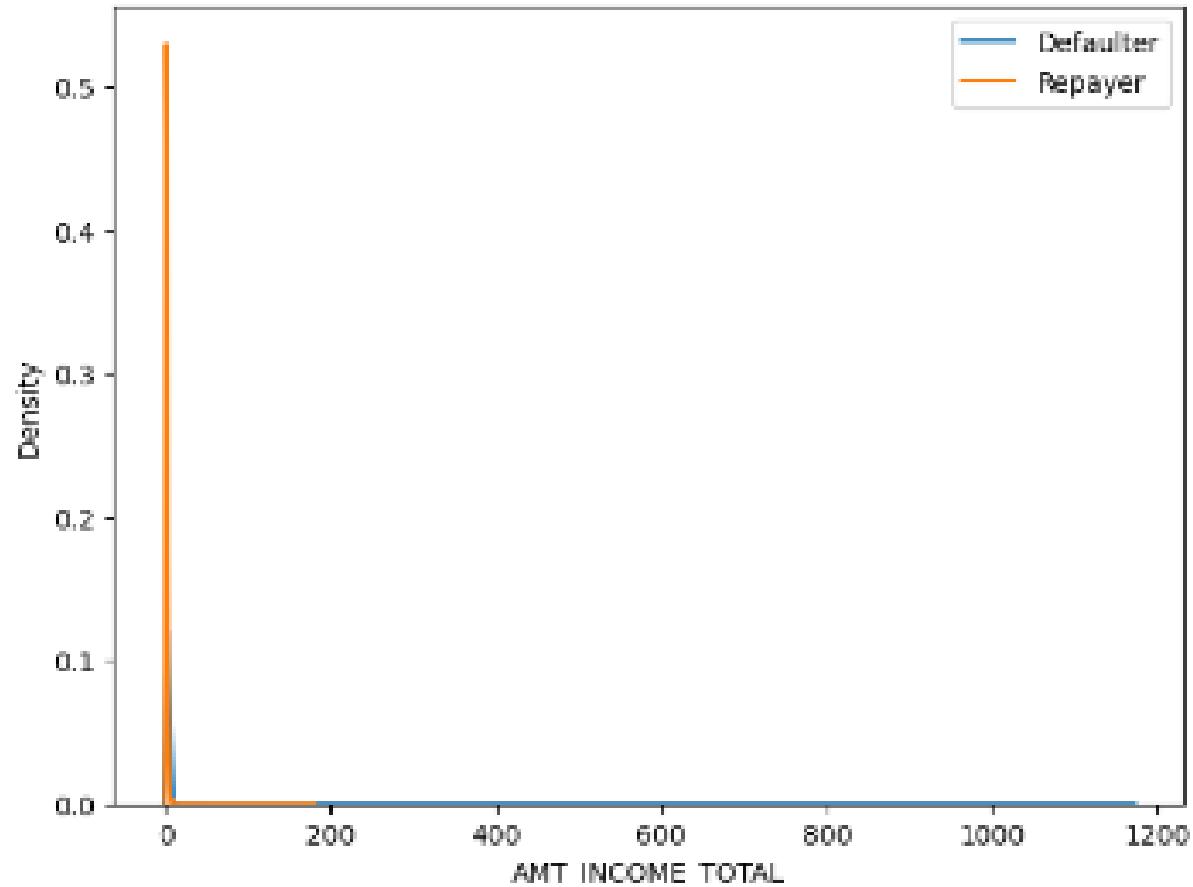


## Inferences

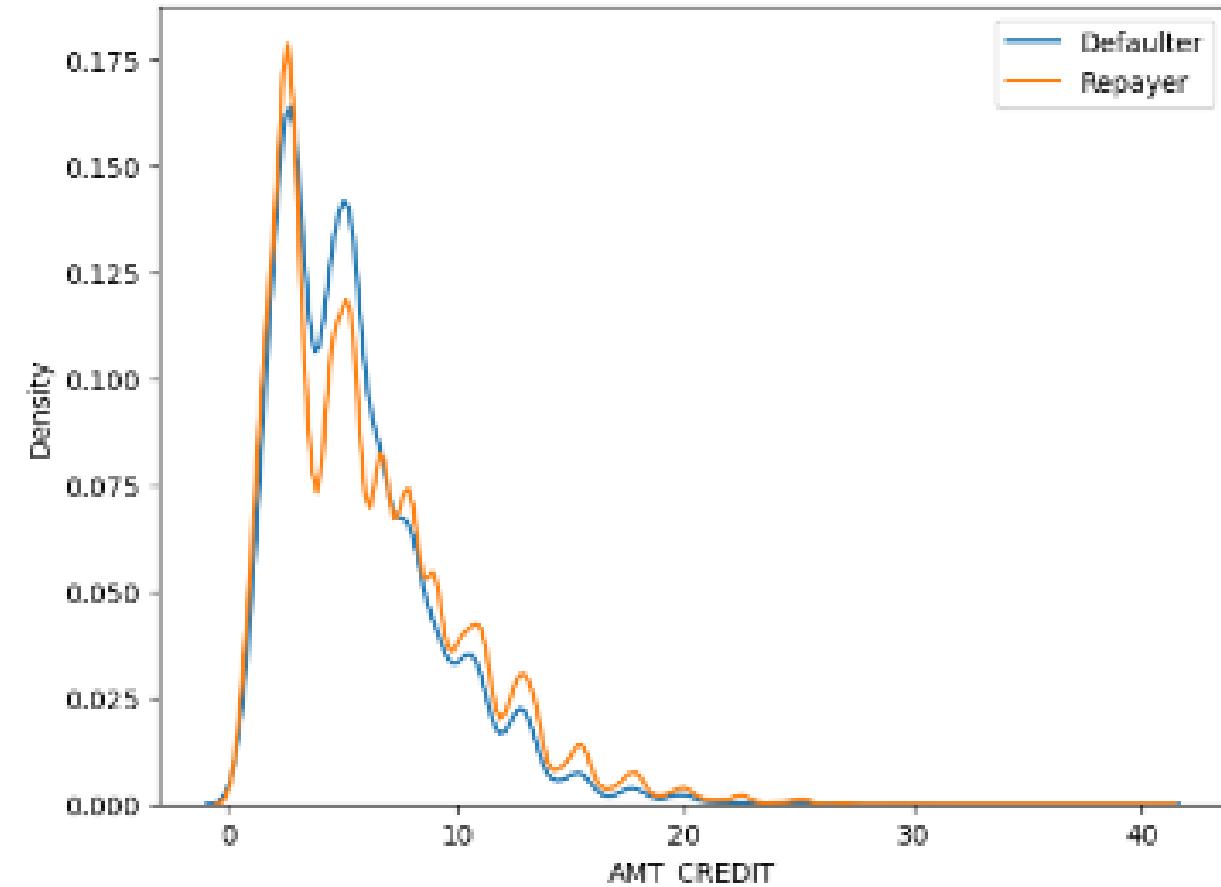
- Credit amount is highly correlated with:
  - Goods Price Amount
  - Loan Annuity
  - Total Income
- We can also see that Repayers have high correlation in number of days employed.
- Credit amount is highly correlated with good price amount which is same as Repayers.
- There is a slight increase in defaulted to observed count in social circle among Defaulters is 0.264) when compared to Repayers is 0.254
- Days\_birth and number of children correlation has reduced to 0.259 in Defaulters when compared to 0.337 in Repayers.
- Loan annuity correlation with credit amount has slightly reduced in Defaulters is 0.75 when compared to Repayers is 0.77
- We can also see that Repayers have high correlation in number of days employed is 0.62 when compared to Defaulters is 0.58. 8 .There is a severe drop in the correlation between total income of the client and the credit amount (0.038) amongst Defaulters whereas it is 0.342 among Repayers.

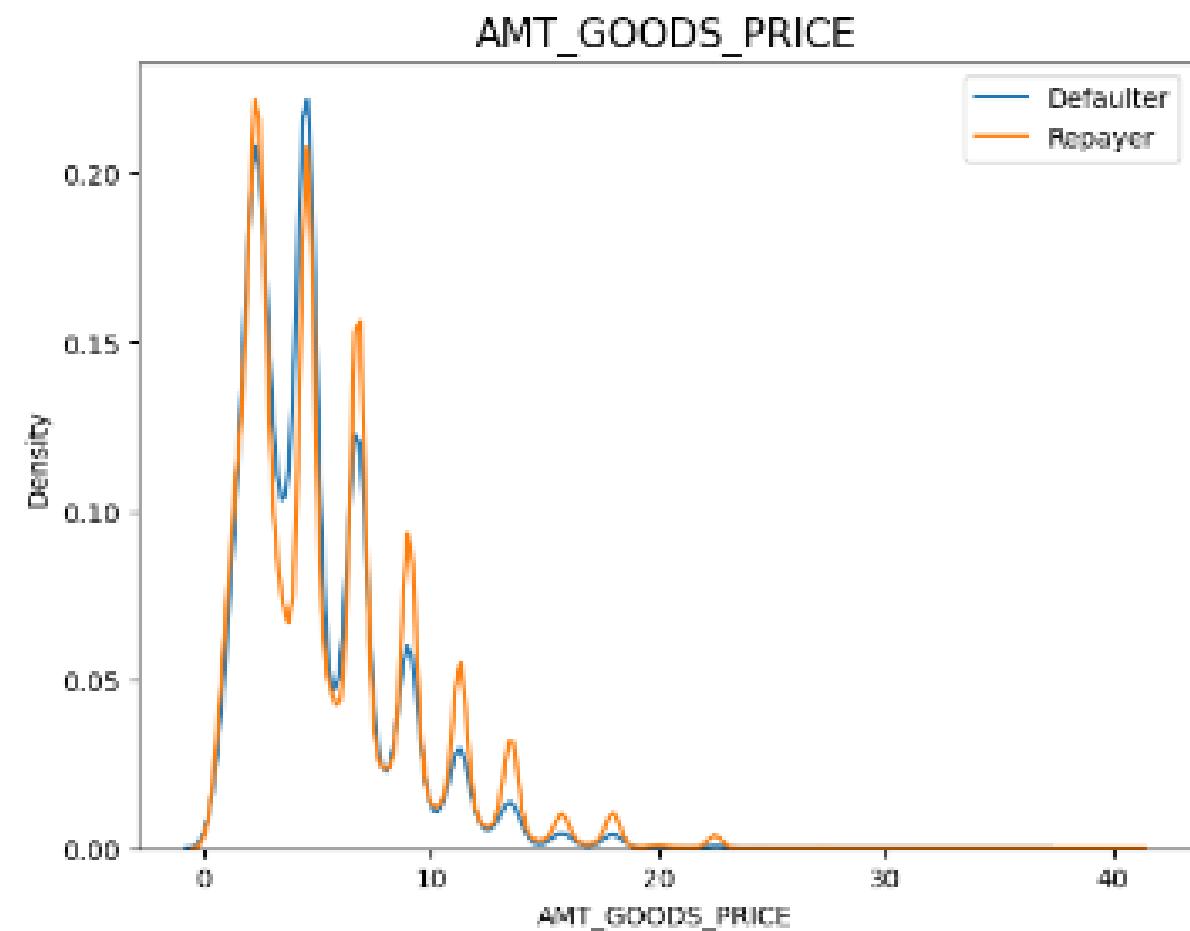
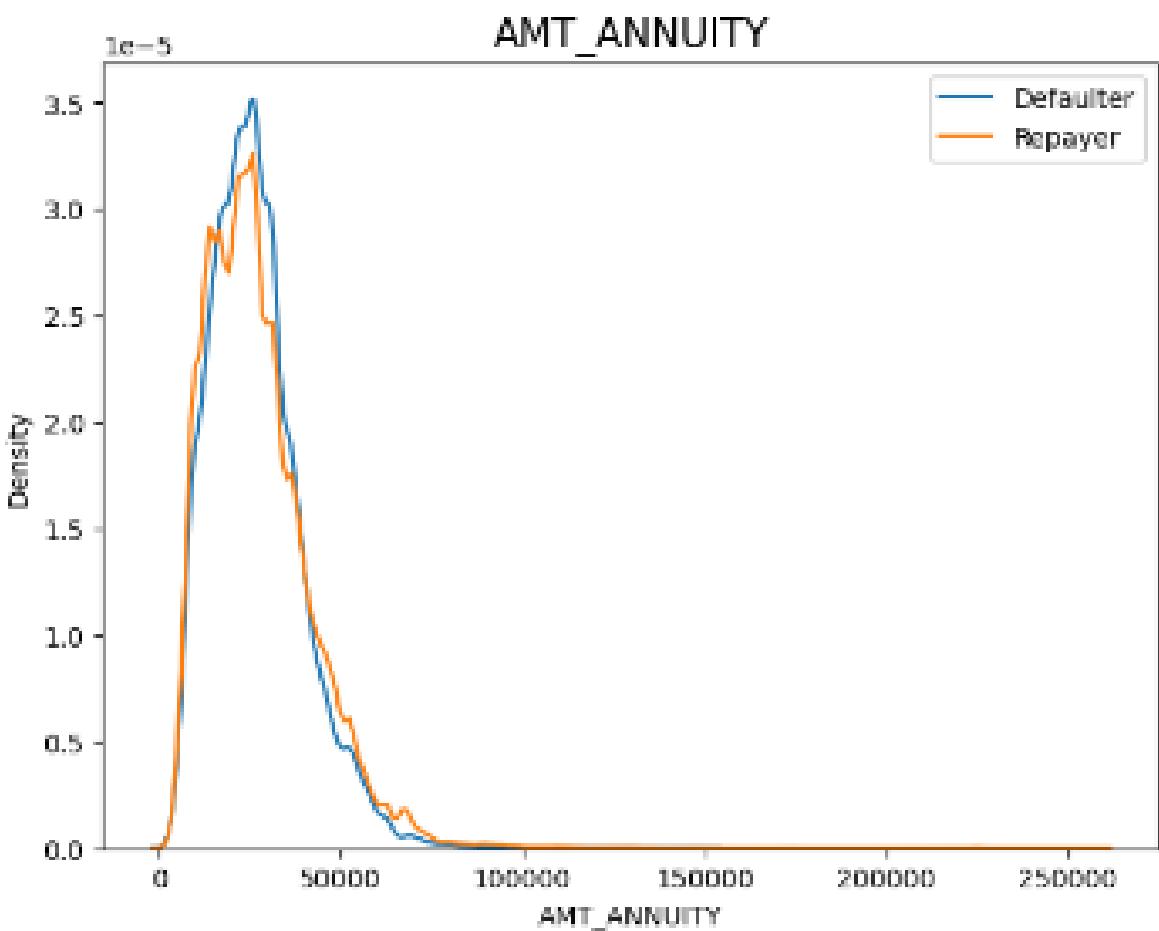
# Density distribution amount

AMT\_INCOME\_TOTAL



AMT\_CREDIT

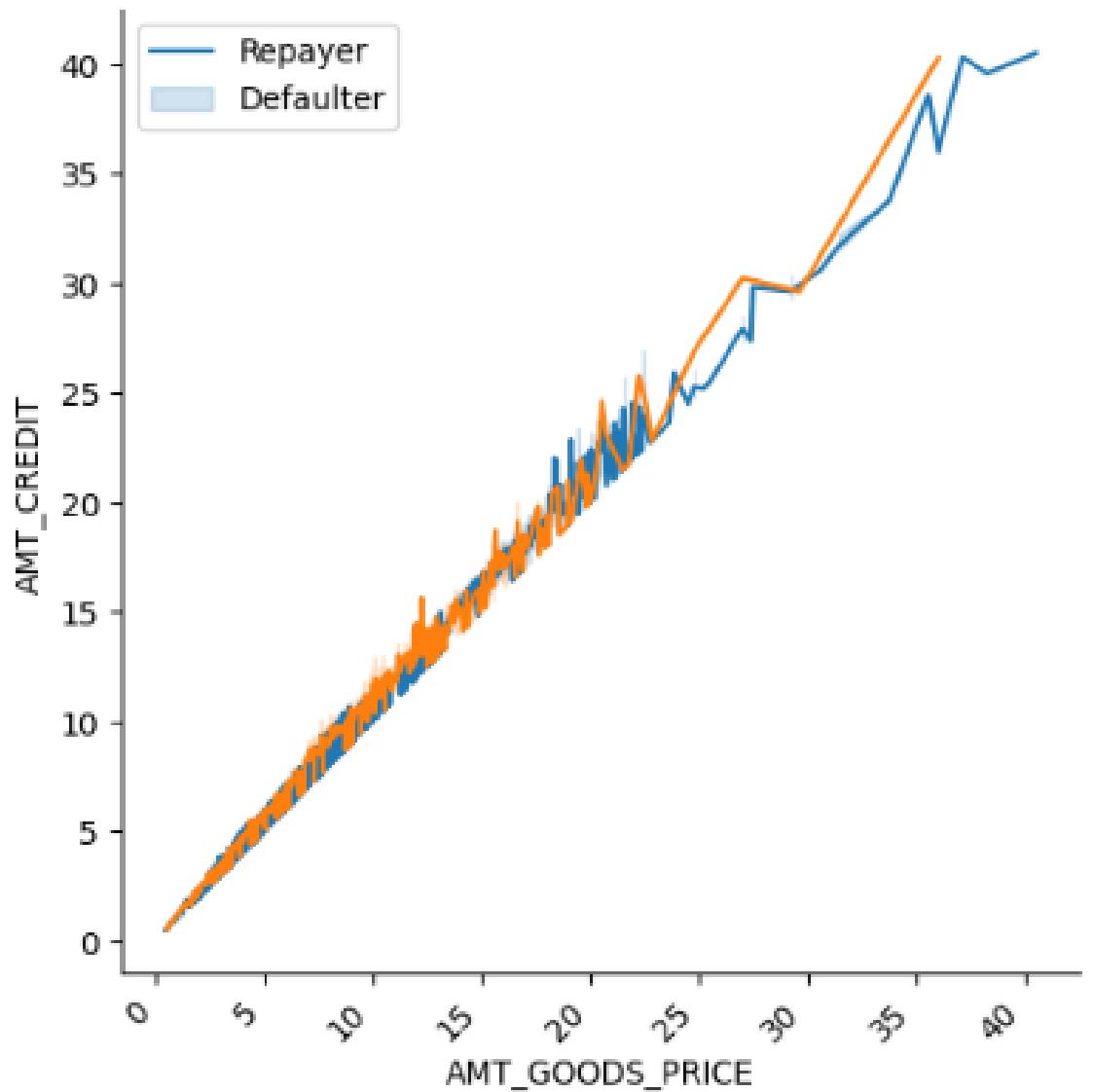




## Inferences

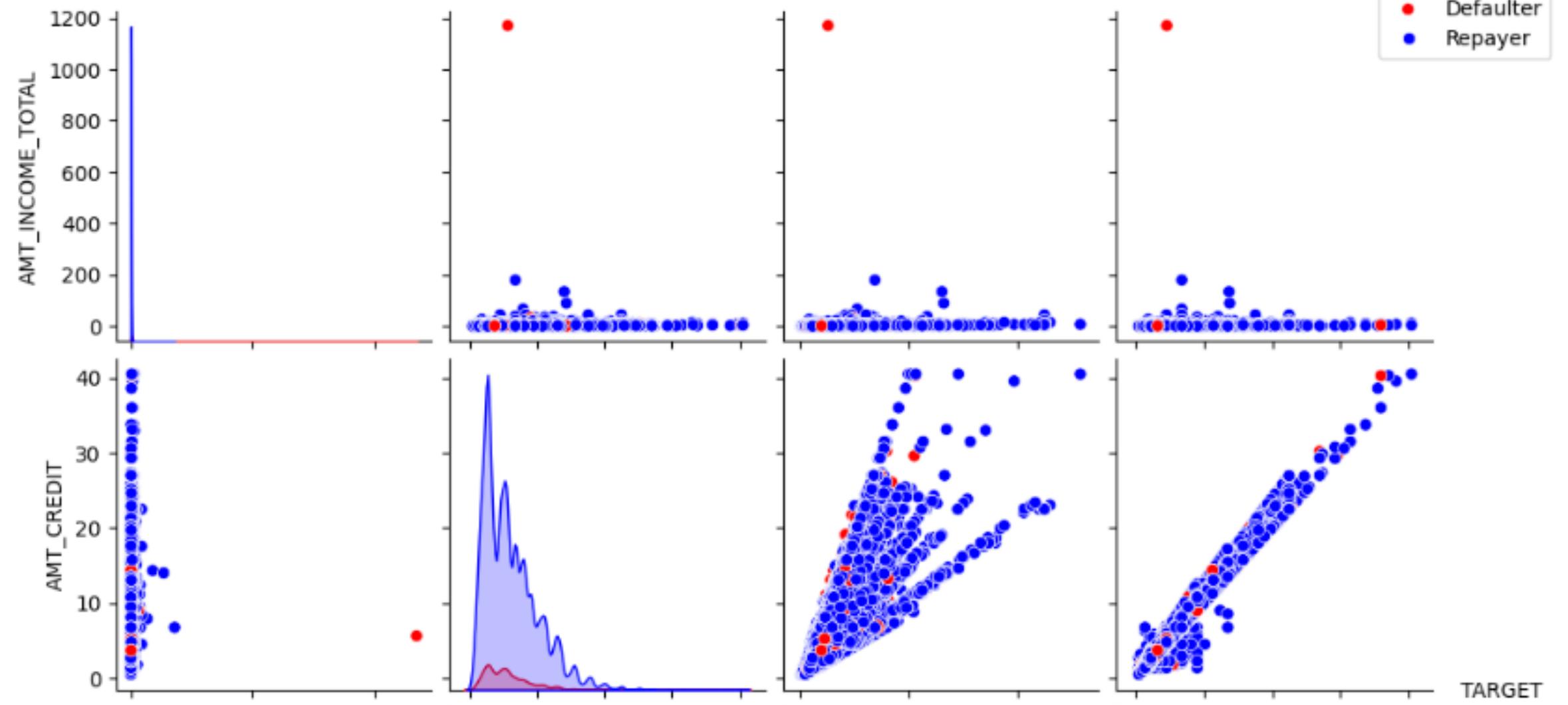
- The Repayers and Defaulters distribution overlap in all the plots and hence we cannot use any of these variables in isolation to make a decision.
- Most no of loans are given for goods price below 10L. Credit amount of the loan is mostly less than 10L
- Most people pay annuity below 50K for the credit loan.

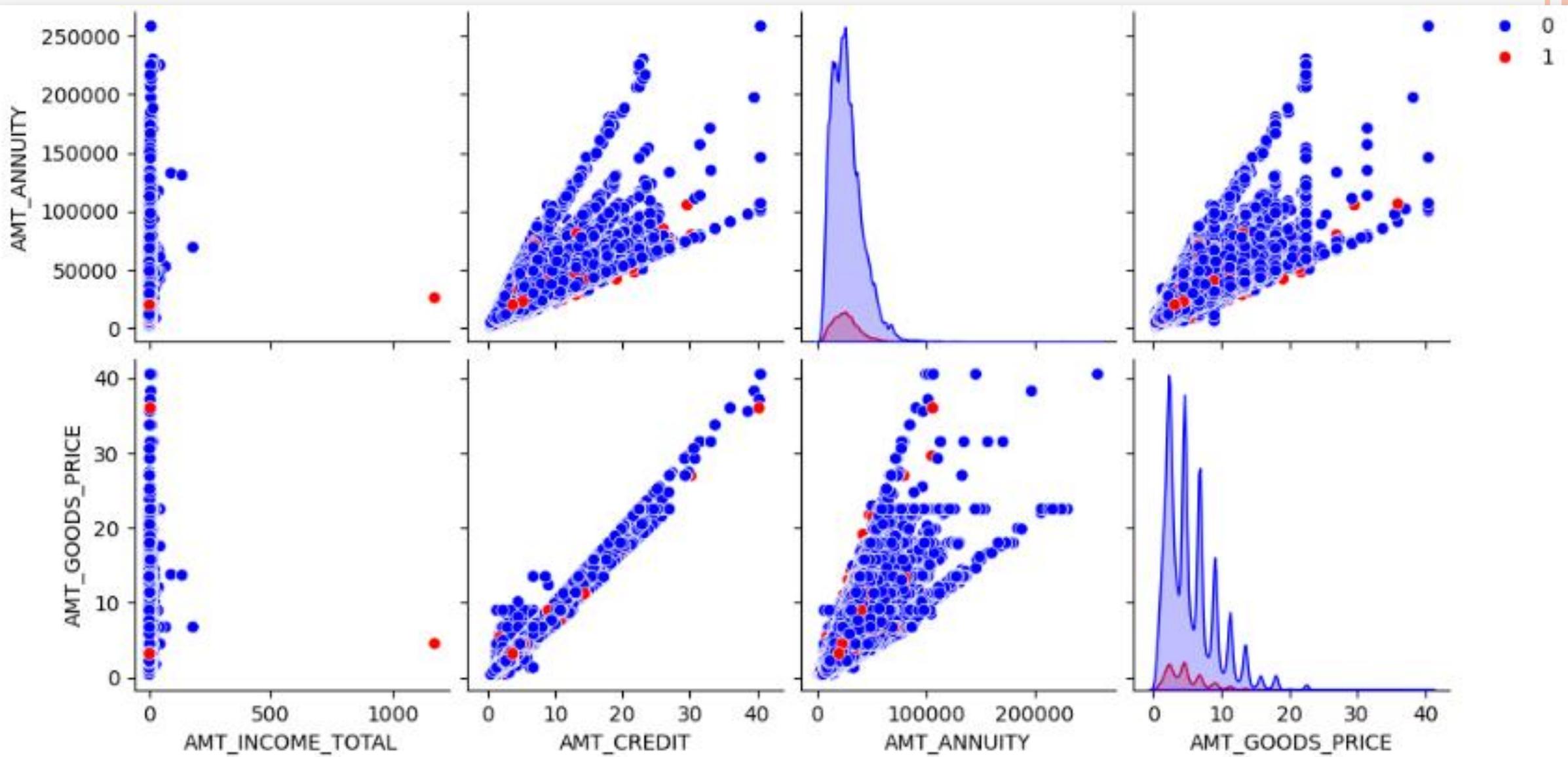
# Numerical Bivariate Analysis



## Inferences

- When the credit amount goes beyond 30L, there is an increase in defaulters.





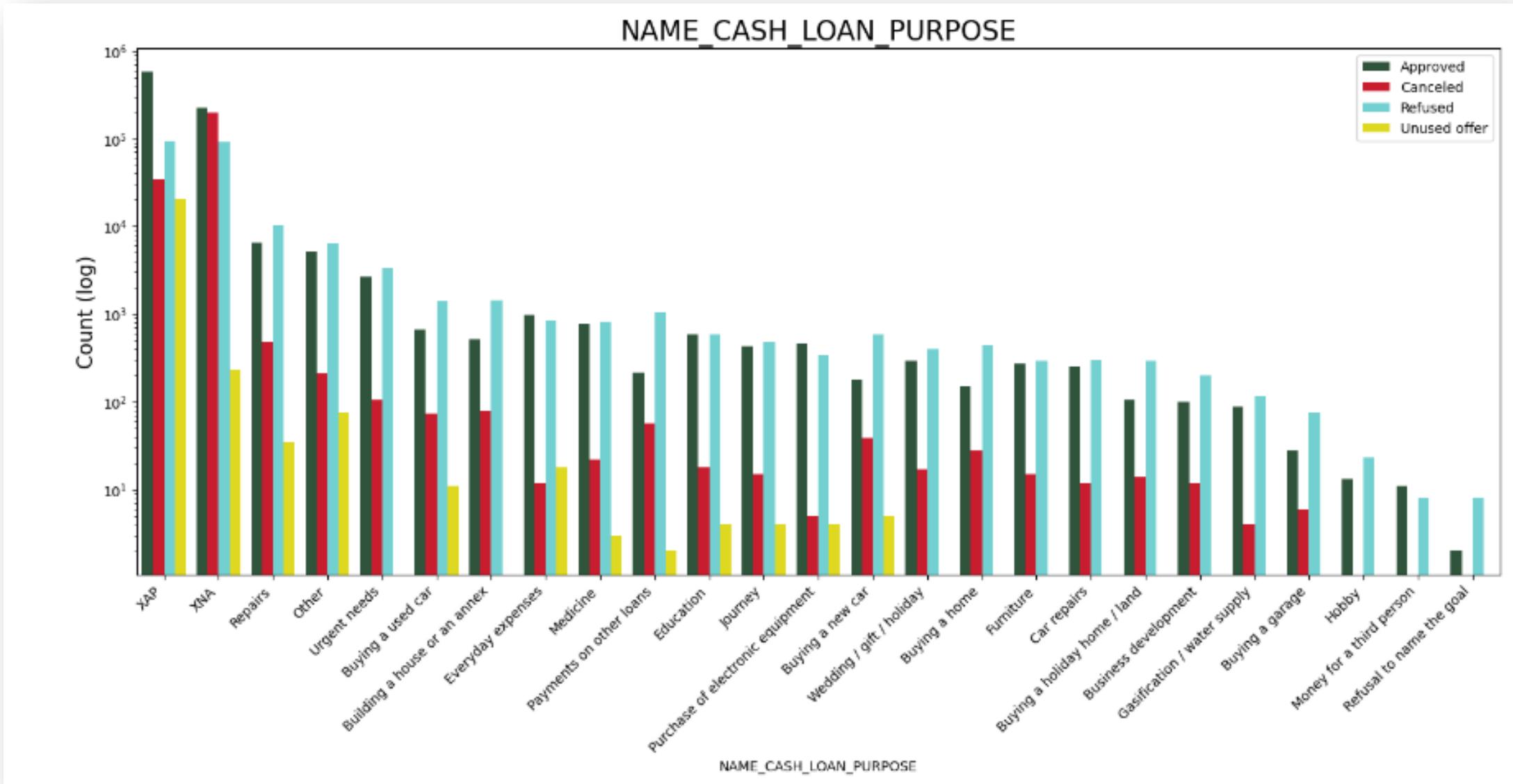
## Inferences

- When Annuity Amount > 15K and Good Price Amount > 20L, there is a lesser chance of Defaulters, and there are very less Defaulters for AMT\_CREDIT >20L
- Loan Amount(AMT\_CREDIT) and Goods price(AMT\_GOODS\_PRICE) are highly correlated as based on the scatter plot where most of the data are consolidated in form of a line

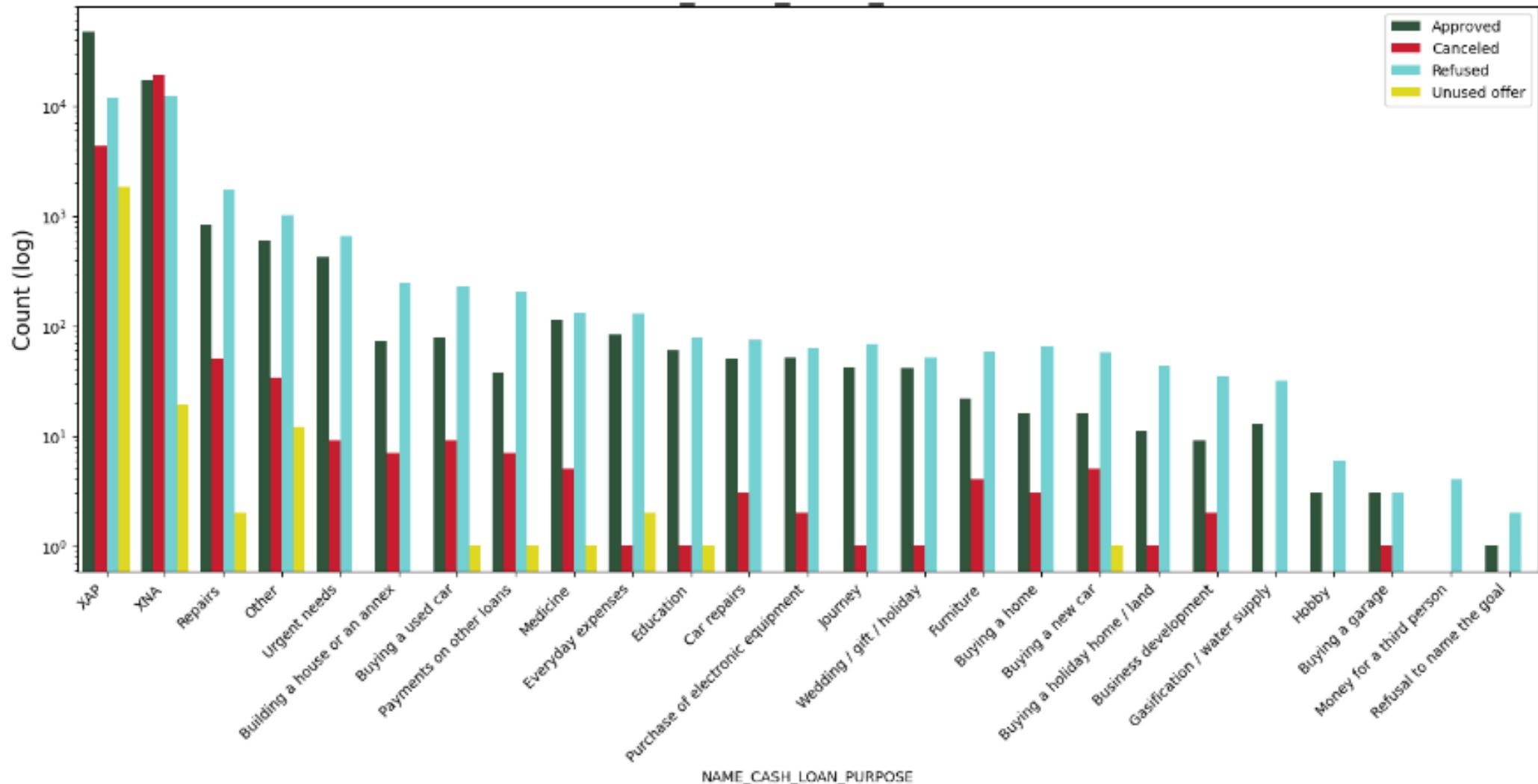
# Merged Data frames Analysis



# Bisecting the "loan\_data" dataframe based on Target value 0 and 1 for correlation and other analysis



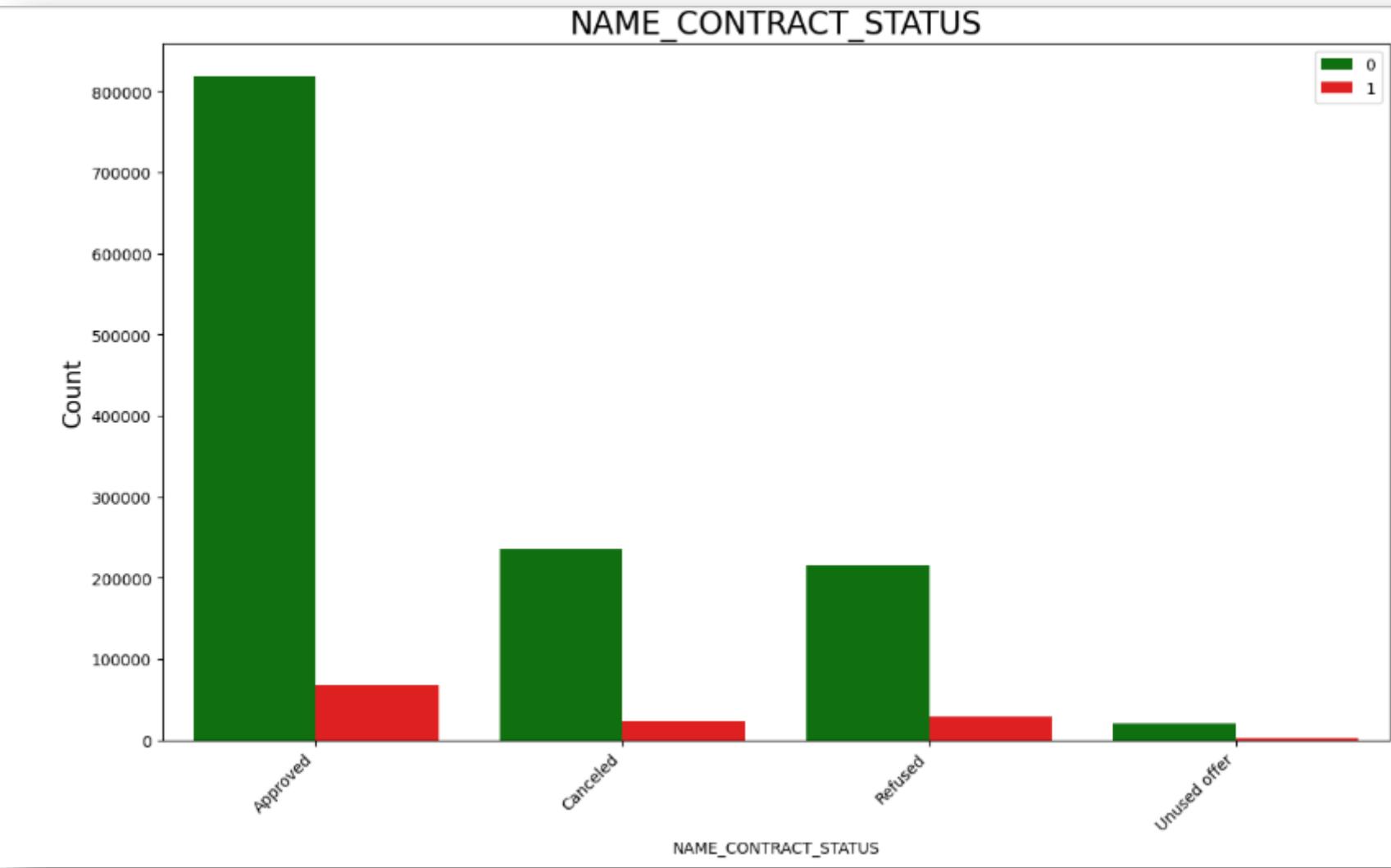
## NAME\_CASH\_LOAN\_PURPOSE



## Inferences

- Loan purpose has high number of unknown values (XAP, XNA) • Loan taken for the purpose of Repairs looks to have highest default rate • Huge number application have been rejected by bank or refused by applicant who is applied for Repair or Other.
- From this we can infer that repair is considered high risk by bank. Also, either they are rejected or bank offers loan on high interest rate which is not feasible by the clients and they refuse the loan.

# Univariate Merged Data

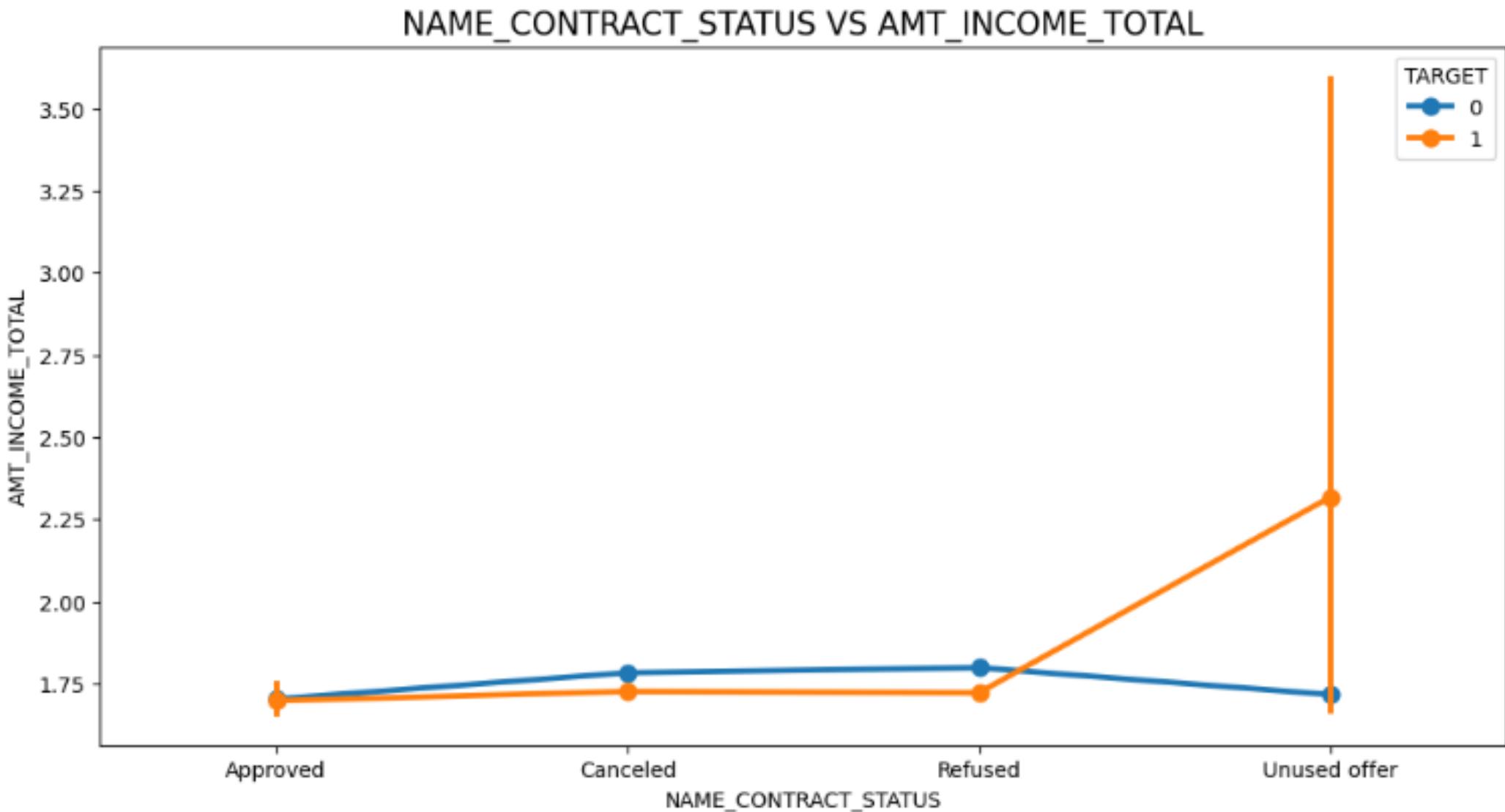


		Counts	Percentage
NAME_CONTRACT_STATUS	TARGET		
Approved	0	818856	92.41%
	1	67243	7.59%
Canceled	0	235641	90.83%
	1	23800	9.17%
Refused	0	215952	88.0%
	1	29438	12.0%
Unused offer	0	20892	91.75%
	1	1879	8.25%

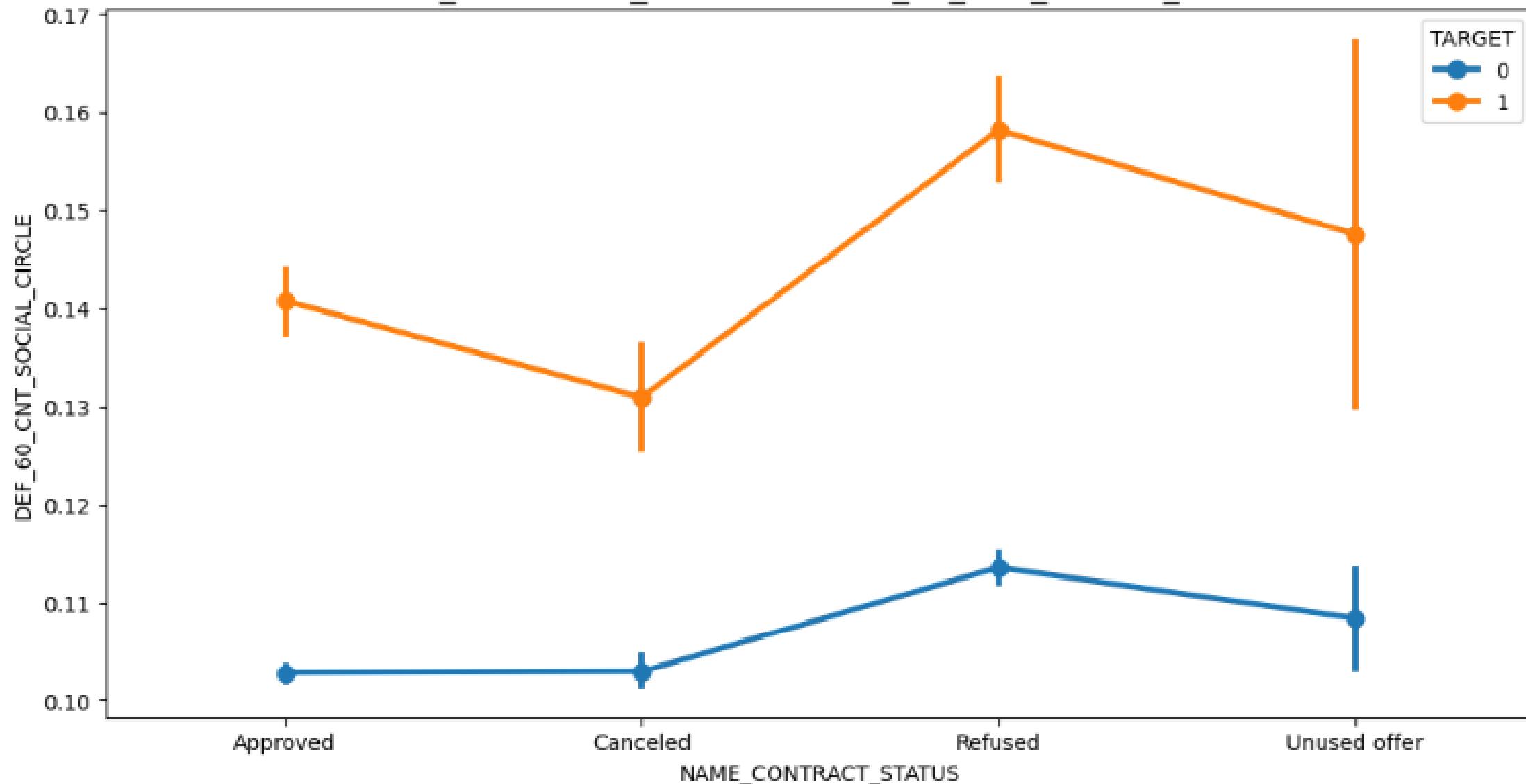
## Inferences

- Refusal reason should be recorded for further analysis as these clients could turn into potential repaying customer.
- 90% of the previously cancelled applicants have actually Repayed the loan. Revising the interest rates would increase business opportunity for these clients
- 88% of the clients who have been previously refused a loan has payed back the loan in current case.

# Bivariate Merged Data



## NAME\_CONTRACT\_STATUS VS DEF\_60\_CNT\_SOCIAL\_CIRCLE



## Inferences

- Applicants who have average of 0.13 or higher their DEF\_60\_CNT\_SOCIAL\_CIRCLE score tend to default more and thus analysing client's social circle could help in disbursement of the loan.
- The point plot show that the applicants who have not used offer earlier have defaulted even when there average income is higher than others

# Conclusions



## Decisive Factor whether an applicant will be Defaulter:

- REGION\_RATING\_APPLICANTS: Applicants who live in Rating 3 has highest defaults.
- CNT\_CHILDREN & CNT\_FAM\_MEMBERS: Applicants who have children => 9 default 100% and so their applications are to be rejected.
- AMT\_GOODS\_PRICE: When the credit amount goes beyond 3L, there is an increase in defaulters.
- OCCUPATION\_TYPE: Due to default rate is huge avoid Low-skill Laborers, Drivers and Waiters/barmen staff, Security staff, Laborers and Cooking staff.
- CODE\_GENDER: Females are good repayer where as men are at relatively higher default rate.
- NAME\_FAMILY\_STATUS : Civil marriage or who are single applicants who have default a lot.

## Decisive Factor whether an applicant will be Defaulter:

- NAME\_EDUCATION\_TYPE: Lower Secondary & Secondary education applicants are default a lot.
- NAME\_INCOME\_TYPE: Maternity leave OR Unemployed applicants are default a lot.
- ORGANIZATION\_TYPE: Organizations with highest percent of loans not repaid are Transport: type 3 (16%), Industry: type 13 (13.5%), Industry: type 8 (12.5%) and Restaurant (less than 12%). Self-employed applicants have relative high defaulting rate, and thus should be avoided to be approved for loan or provide loan with higher interest rate to might be the risk of defaulting.
- DAYS\_BIRTH: Avoid young applicants who are in age group of 20-40 as they have higher probability of defaulting
- DAYS\_EMPLOYED: Applicants who have less than 5 years of employment have high default rate.

## Decisive Factor whether an applicant will be Repayer:

- CNT\_CHILDREN: Applicants with 0 to 2 children tend to repay the loans.
- DAYS\_BIRTH: If age above of 50 have low probability of defaulting.
- DAYS\_EMPLOYED: 40+ year experience having less than 1% default rate.
- AMT\_INCOME\_TOTAL: Income more than 700,000 are less likely to default.
- NAME\_EDUCATION\_TYPE: Academic degree has fewer defaults.

## Decisive Factor whether an applicant will be Defaulter:

- NAME\_INCOME\_TYPE: Student and Businessmen have no defaults.
- REGION\_RATING\_APPLICANTS: RATING 1 is safer.
- ORGANIZATION\_TYPE: Trade Type 4 and 5 and Industry type 8 have defaulted less than 3%
- NAME\_CASH\_LOAN\_PURPOSE: Loans bought for Hobby, Buying garage are being repayed mostly.

## **Below are few conditions where interest rate is high might getting any default risk leading to business loss:**

- CNT\_CHILDREN & CNT\_FAM\_MEMBERS: higher interest should be imposed on their loans., who have 4 to 8 children has a very high default rate.
- AMT\_INCOME: Those applicants could be offered loan with higher interest compared to other income category, since 90% of the applications have Income total less than 3 and they have high probability of defaulting.
- NAME\_HOUSING\_TYPE: Offering the loan would might the loss if any of applicants are from the category of people who live in rented apartments & living with parents.
- AMT\_CREDIT: Having higher interest specifically for Applicants who get loan for 3L - 6L credit range would be ideal because they tend to default more than others.
- NAME\_CASH\_LOAN\_PURPOSE: Repairs purpose seems to have highest default rate. Very high number applications have been rejected by bank or refused by client in previous applications as well which has purpose as repair or other. This shows that purpose repair is taken as high risk by bank and either bank offers very high loan interest rate which they are rejected.

## **Business Advices after this EDA as below:**

- 88% of the applicants who were refused by bank for loan earlier have now turned into a repaying applicants. Hence documenting the reason for rejection could moderate the business loss and these clients could be contacted for further loans.
  
- 90% of the previously cancelled applicants have actually repayed the loan. Record the reason for cancellation which might help the bank to determine and negotiate terms with these repaying customers in future for increase business opportunity.



Thank You!!

