

# CREDIT EDA CASE STUDY

ROHAN GHOGARE

JYOTHSI J

# PROBLEM STATEMENT

## **Introduction:**

- This study aims at applying EDA principals to understand the risk associated in any financial lending institutions
- This EDA study needs to identify the driving factors which can help to find out potential defaulters and also whether to offer a loan to clients. This way financial institution can cut back financial losses and make some profitable informed business decisions

## **Data Provided :**

- Spread sheet containing all the information of the client at the time of application. The data talks about attributes which has impact on clients payment difficulties.
- Spread sheet contains information about the client's previous loan data. This data shows whether previous application had been Approved, Cancelled, Refused or Unused offer.

# SOLUTION APPROACH

## **Approach for case study:**

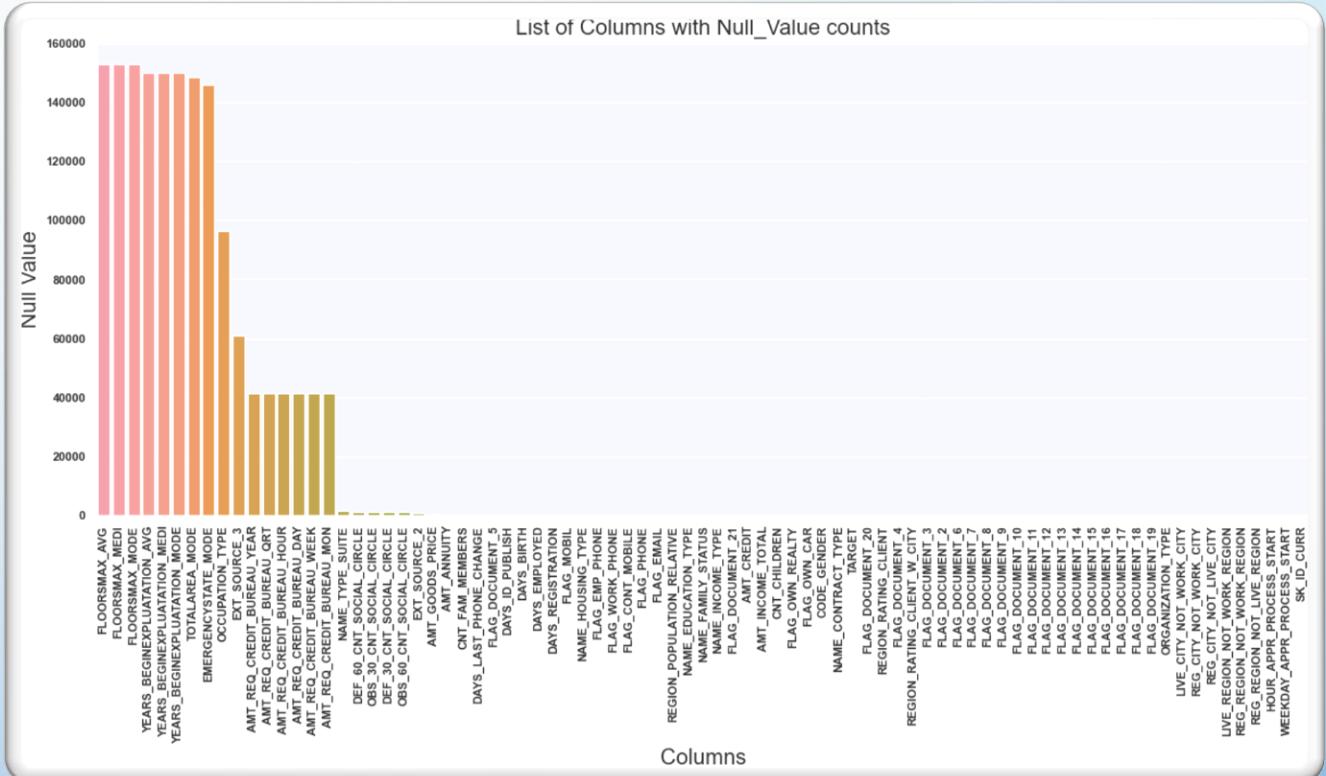
- Load/ Read all data provided and identify the data structure, information about data like data types, missing values.
- Data Cleaning Techniques :
  - ❖ Investigate the missing values and impute that with suitable method
  - ❖ Analysing if there is any outliers in data
  - ❖ Check data imbalance
- Data analysis
  - ❖ Cleaned data used to investigate the relationship of variables with that of TARGET variable
  - ❖ Univariate and bivariate analysis to understand the patterns in data and drew conclusions
  - ❖ Data visualization techniques to summarize the results

# DATA ANALYSIS

# DATA CLEANING

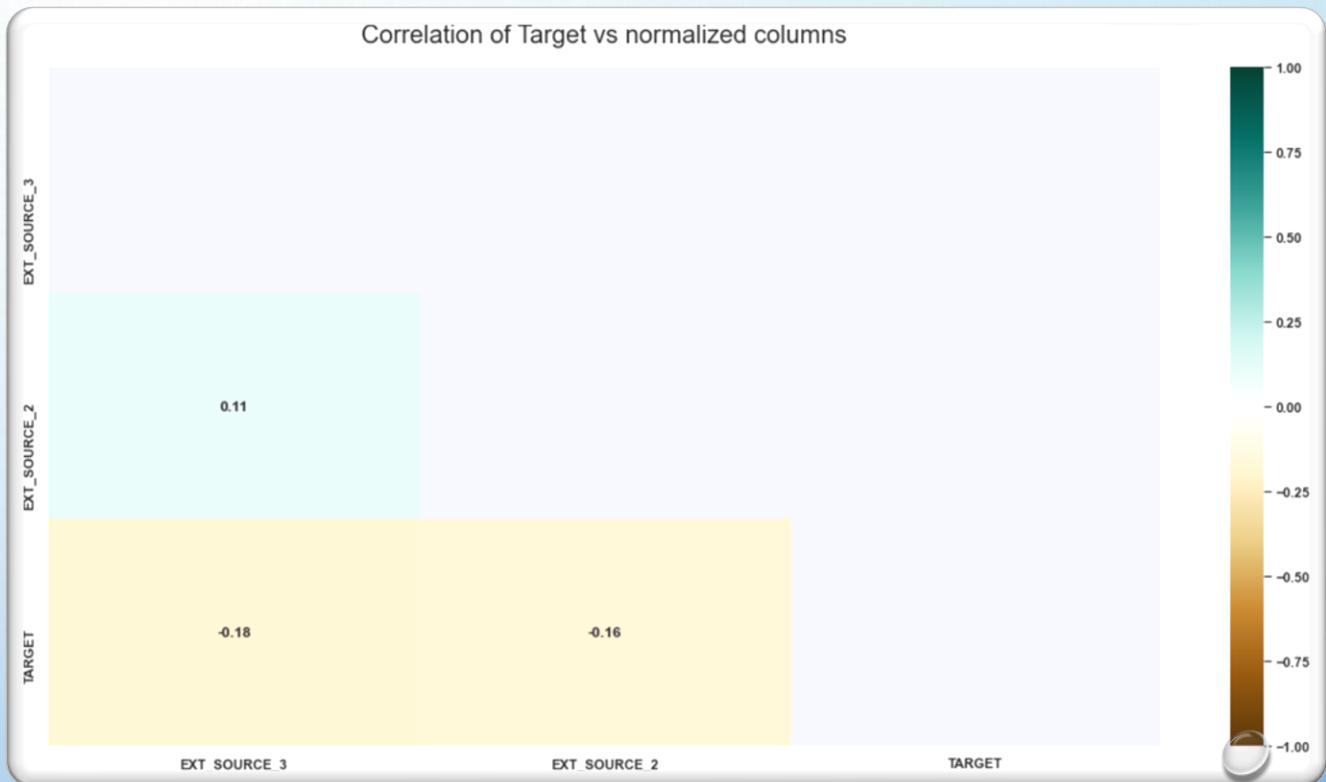
- ***INFERENCES :***

- Initially there are total 122 columns & 307511 rows
- Null values with higher % identified(say >50%) and dropped
- The columns with lesser % of missing values are identified and imputed with suitable technique viz. mean, mode
- Next the columns which have +ve and -ve values require fixing



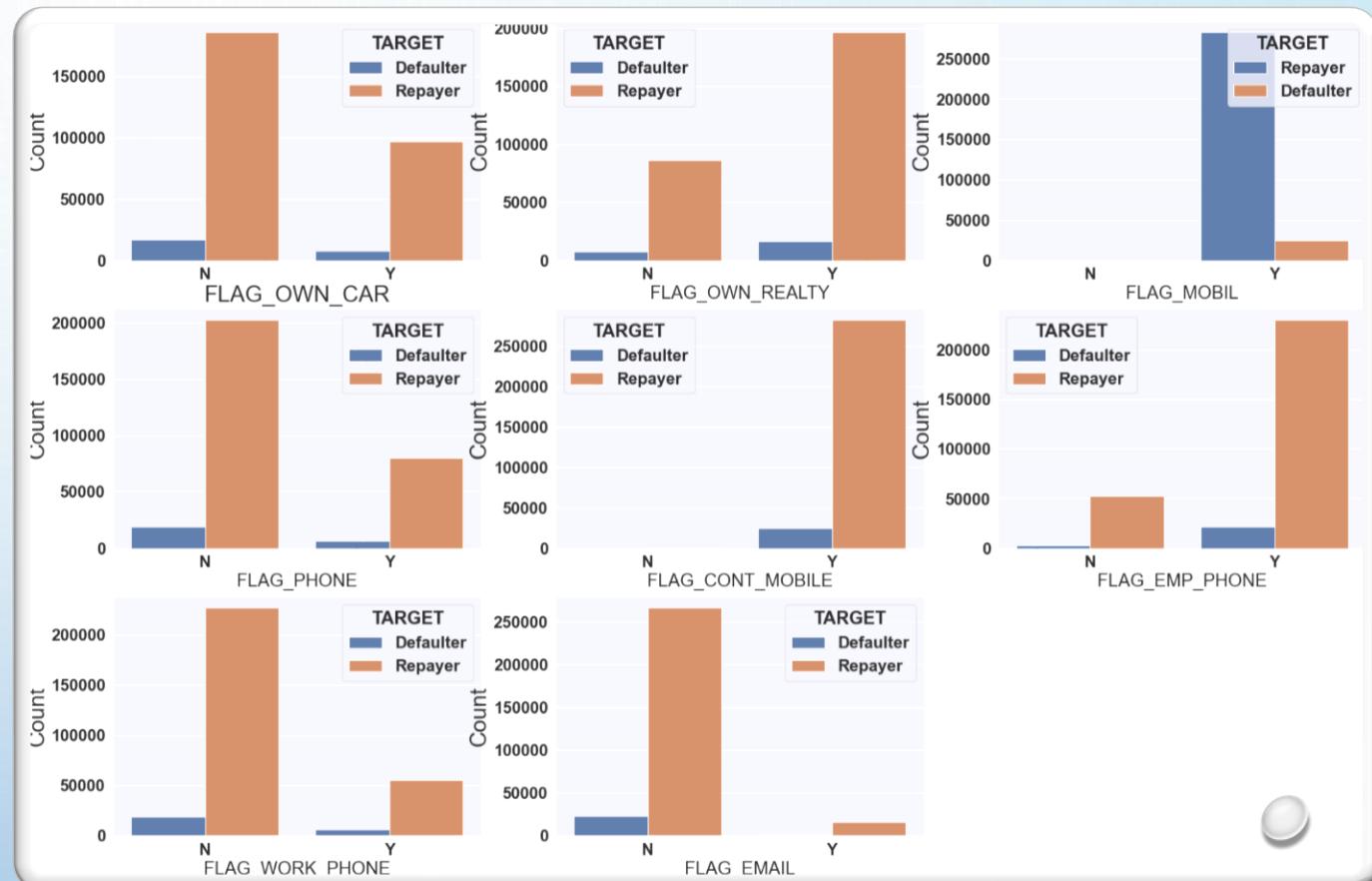
# DATA CLEANING

- **INFERENCES** : Correlation is used to find relationship with TARGET variable and if its not can be eliminated
- Correlation used above to represent the statistical measure of linear relationship between two variables and target variable. This can help to measure of dependence between two different variables. The Correlation between Target and other normalized column is Weak negative suggesting no causal relationship. Hence we drop these two columns



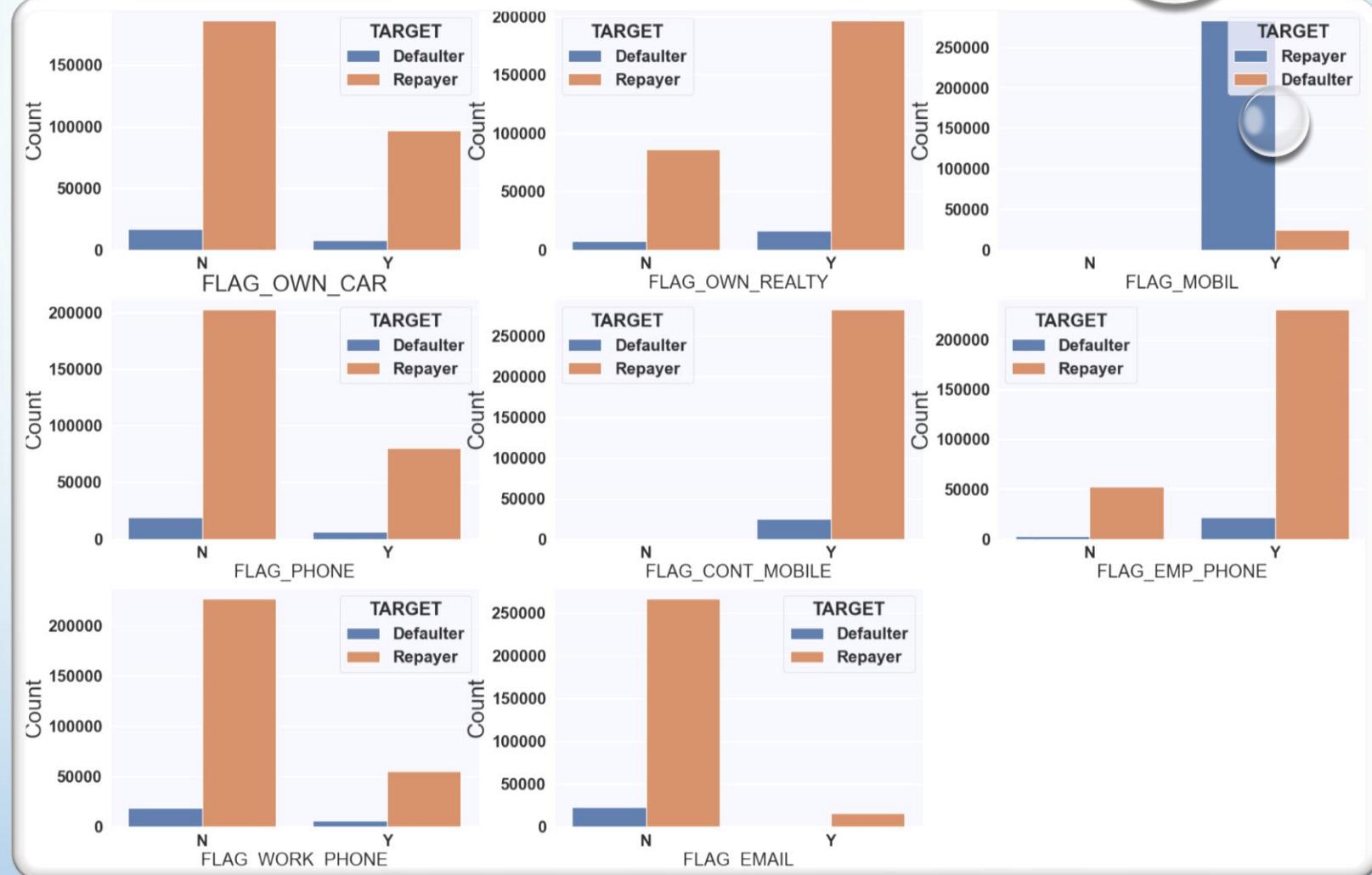
# DATA CLEANING

- **INFERENCES** : Plot all flag columns with respect to Target variable to get insight
- Columns like FLAG\_OWN\_REALTY, FLAG\_MOBIL , FLAG\_CONT\_MOBILE, FLAG\_EMP\_PHONE, shows there are more repayor's than defaulter and hence keep these columns and remove all other columns for further analysis



# DATA CLEANING

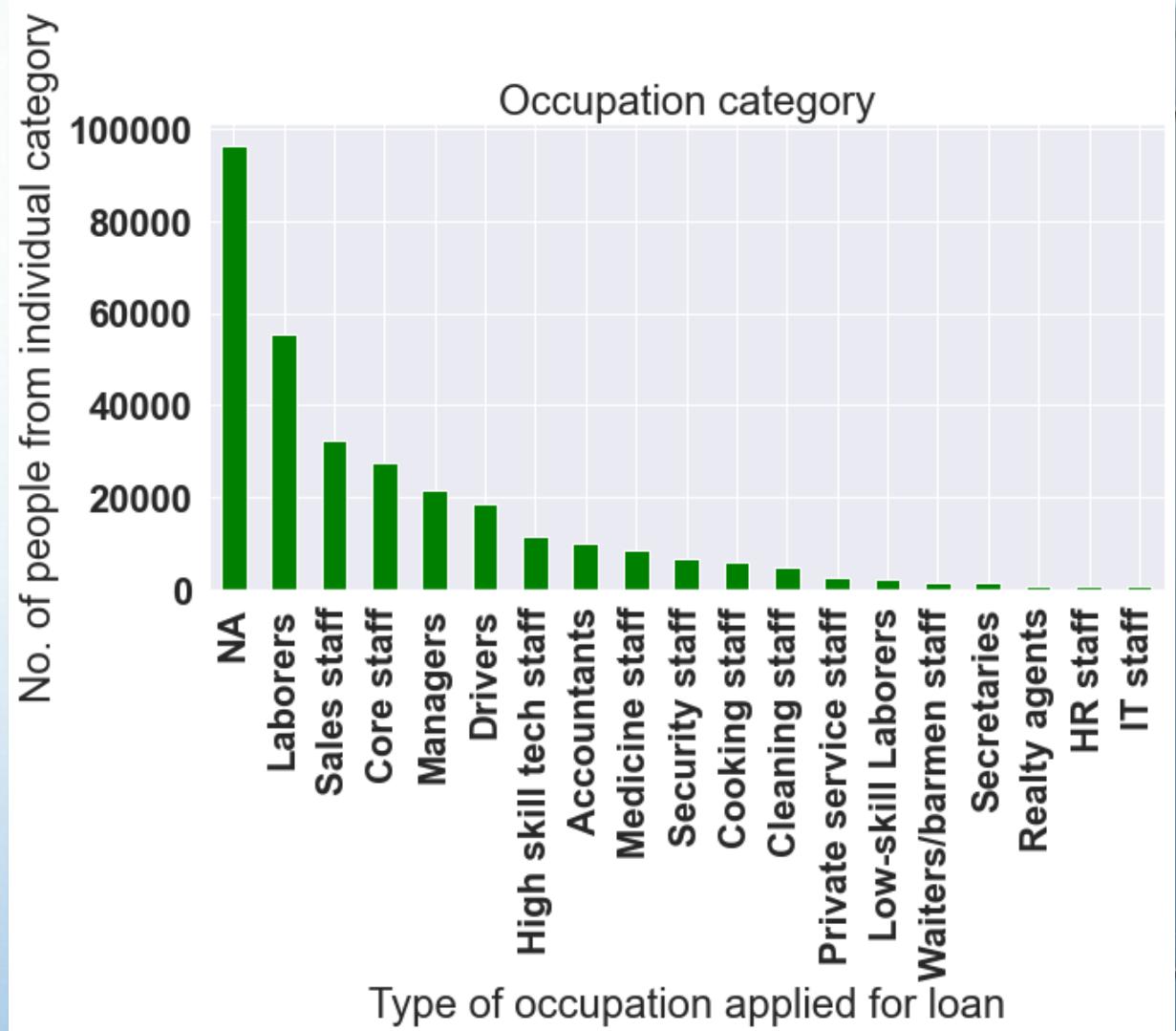
- **INFERENCES** : Plot all flag columns with respect to Target variable to get insight
- Columns like FLAG\_OWN\_REALTY, FLAG\_MOBIL, FLAG\_CONT\_MOBILE, FLAG\_EMP\_PHONE, shows there are more repayers than defaulter and hence keep these columns and remove all other columns for further analysis



**Note :** Similar approach has been followed with other variable and shown in python file

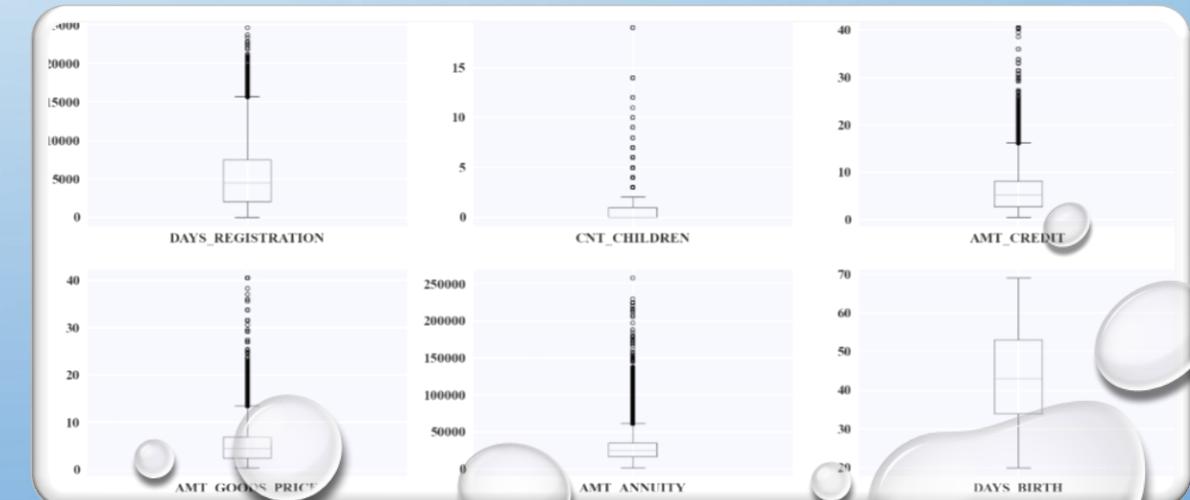
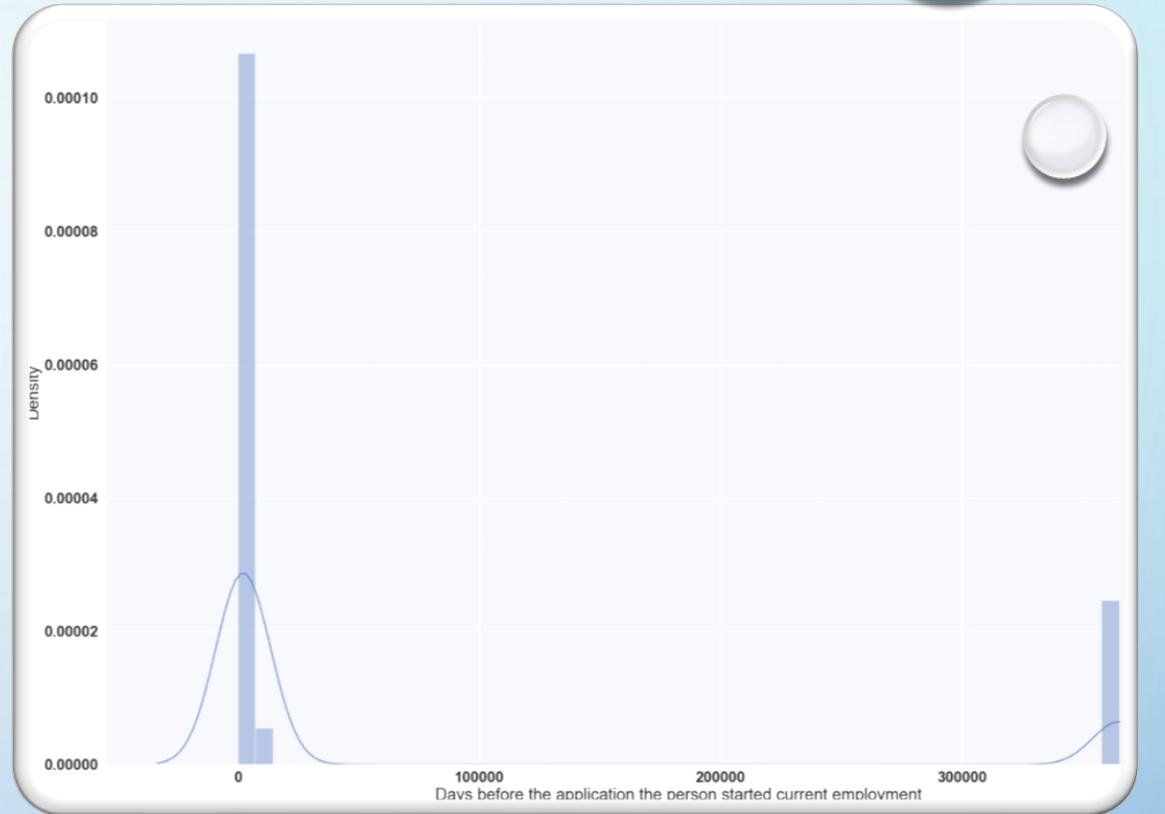
# DATA IMPUTATION

- columns having null values greater than 1% , will impute that lets do it one by one
- **INFERENCES:**
  - NA are the highest % amongst all people applied for loan
  - The columns imputed with mode and median



# INVESTIGATING OUTLIERS

- **INFERENCES :**
- The column 'DAYS\_EMPLOYED' shows Days before the application the person started current employment. Hence if it is more than 20,000 that can be outlier
- Does not see any outlier in DAYS\_BIRTH
- In column 'AMT\_ANNUITY' a single value of > than 250000 is an outlier.
- In the column 'DAYS\_REGISTRATION' a value greater than 24000 is an outlier.

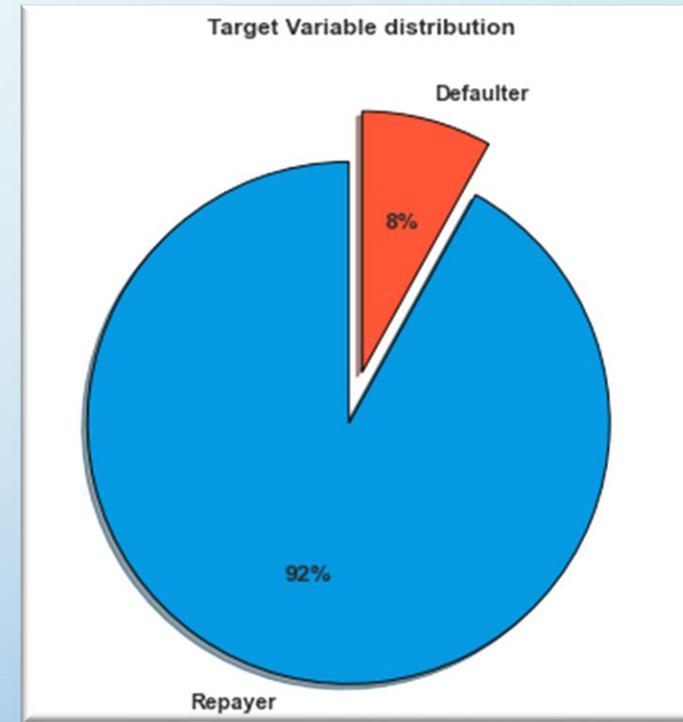


# **IDENTIFYING DRIVING FACTORS FOR BUSINESS DECISIONS USING DATA ANALYSIS**

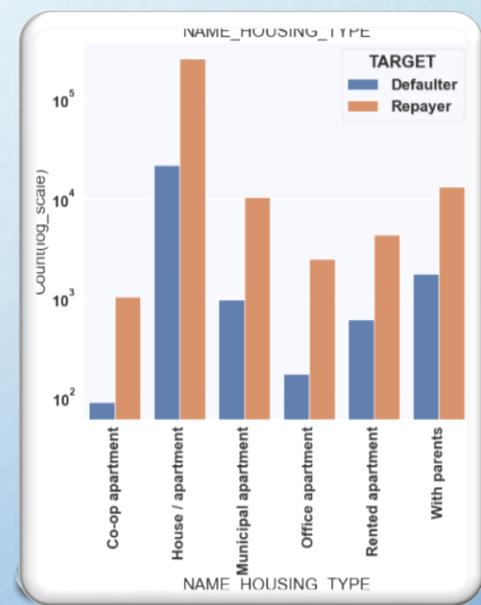
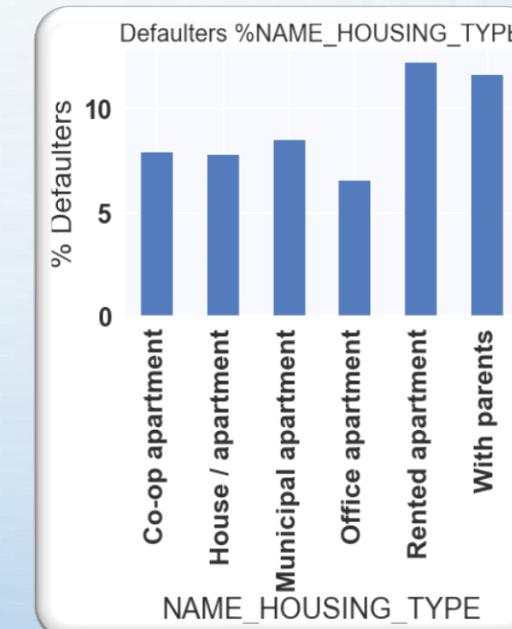
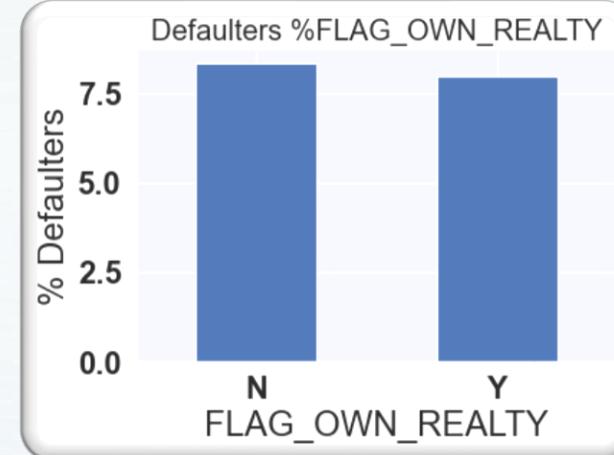
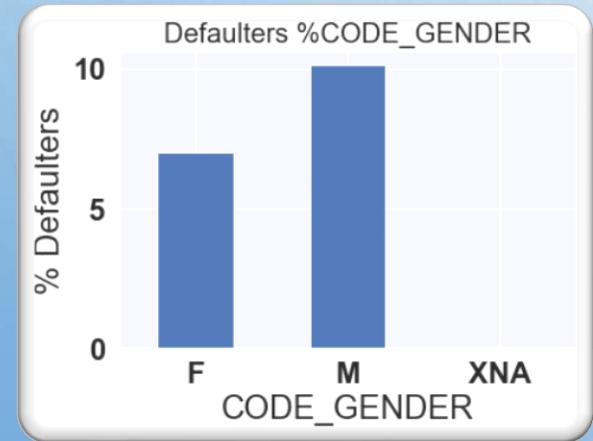
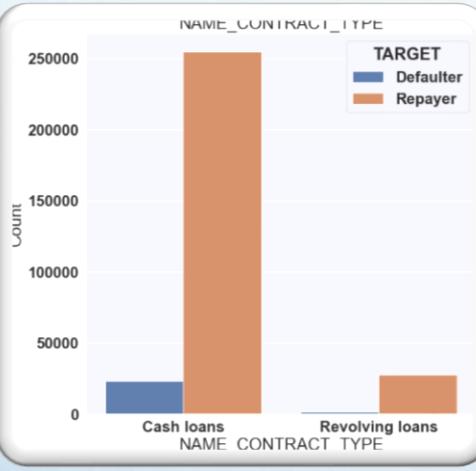
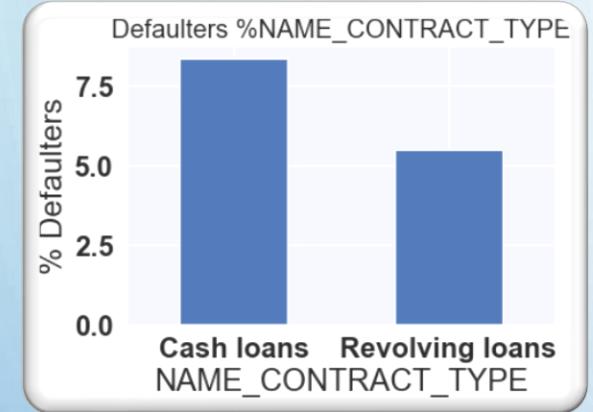
# DATA IMBALANCE

## INFERENCES :

- The data when analysed w.r.t Target variable it clearly indicates imbalance. (Ratio of Data Imbalance 11.4)
- The univariate and bivariate analysis will make things more clear

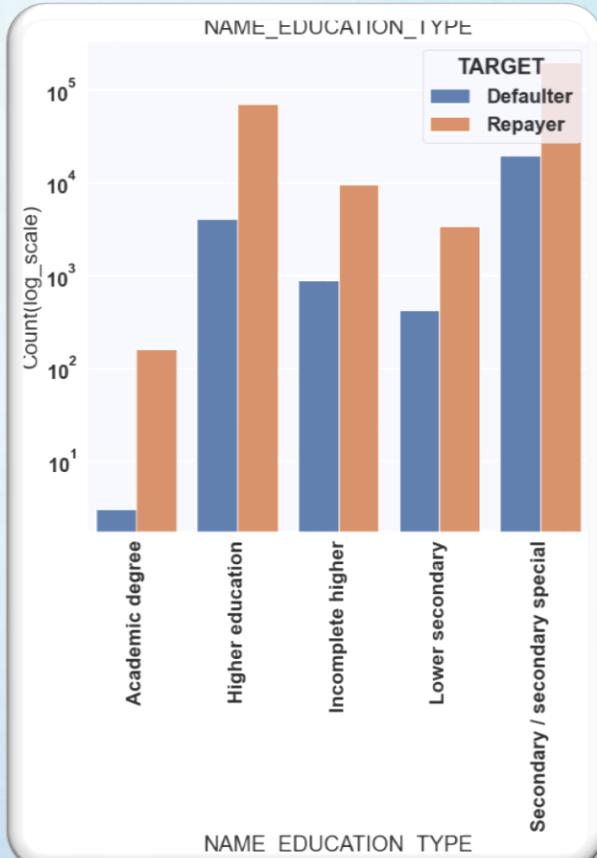
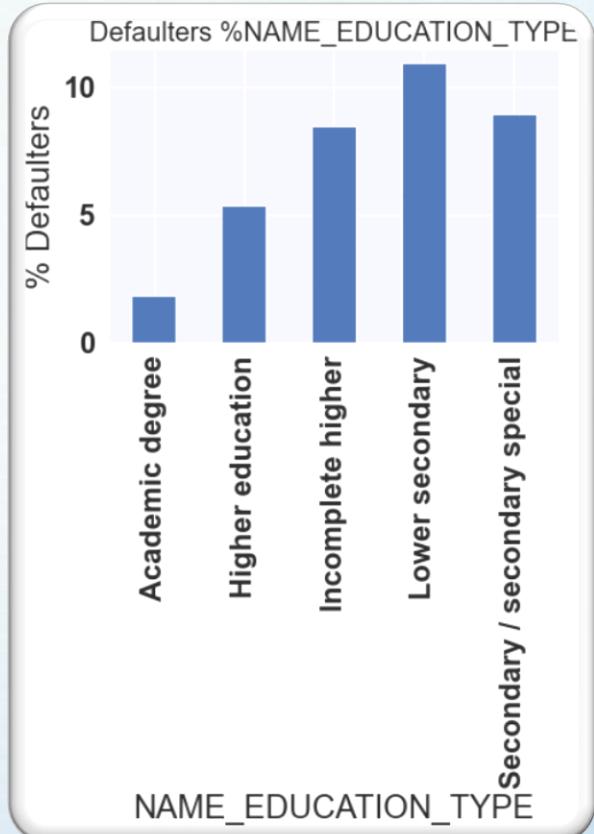
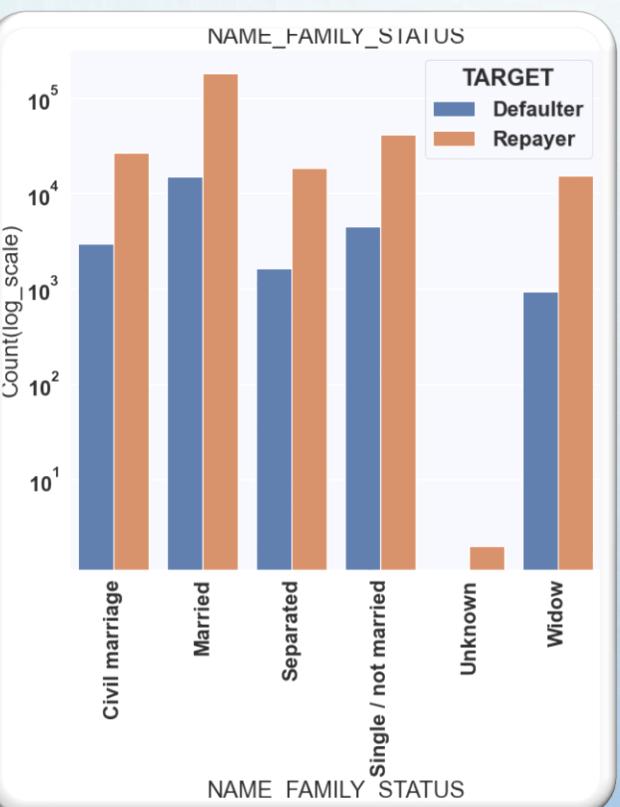
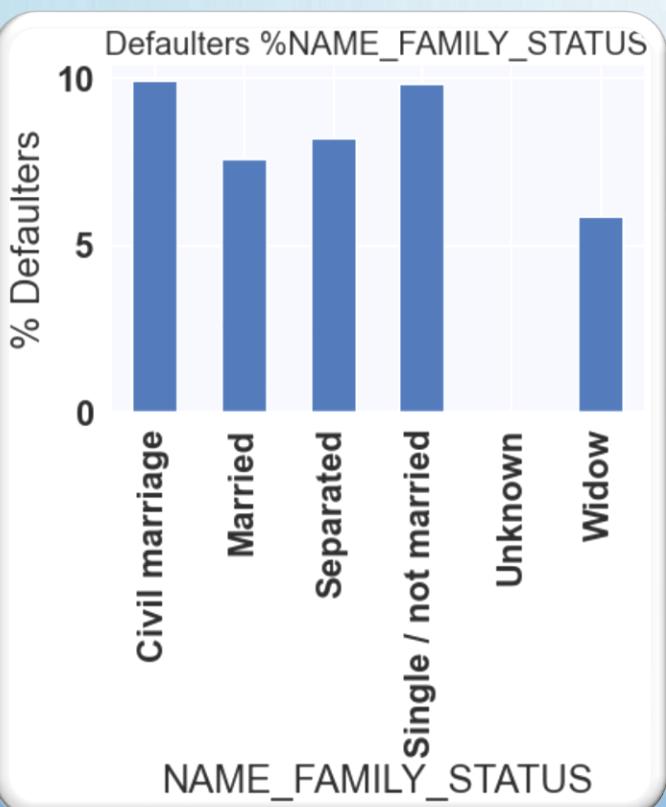


# UNIVARIATE ANALYSIS



**Inferences** : Columns likes FLAG\_OWN\_REALTY, FLAG\_MOBIL , FLAG\_CONT\_MOBILE, FLAG\_EMP\_PHONE, shows there are more repayers than defaulter and hence keep these columns

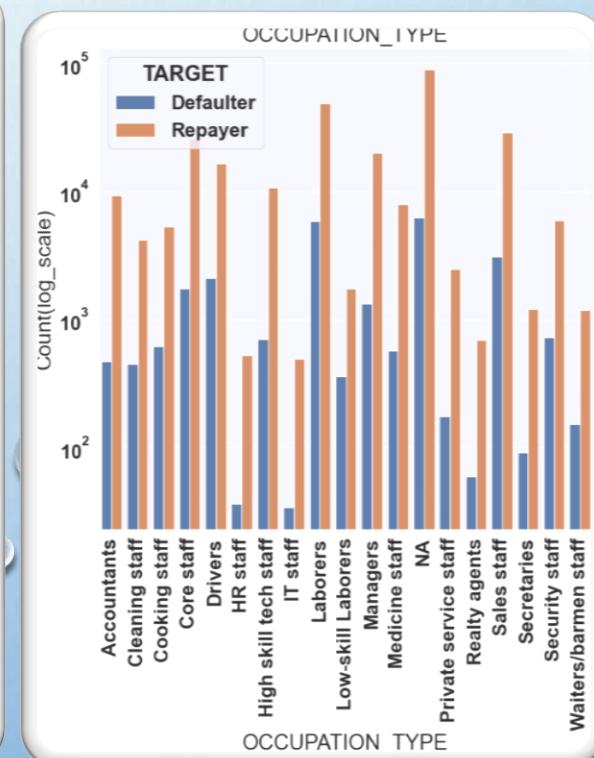
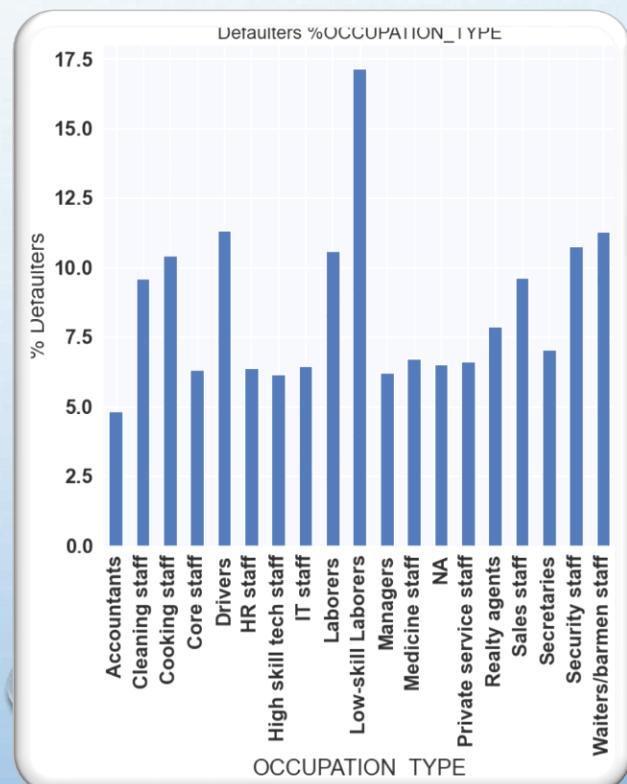
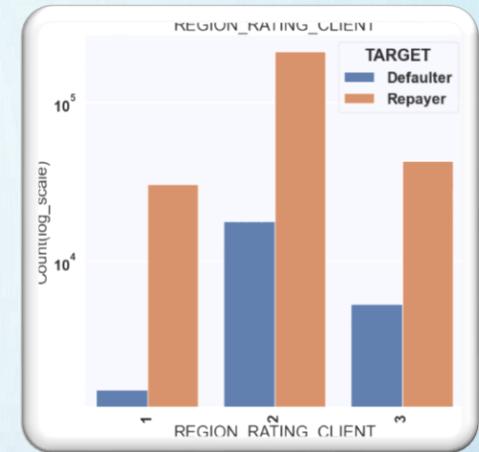
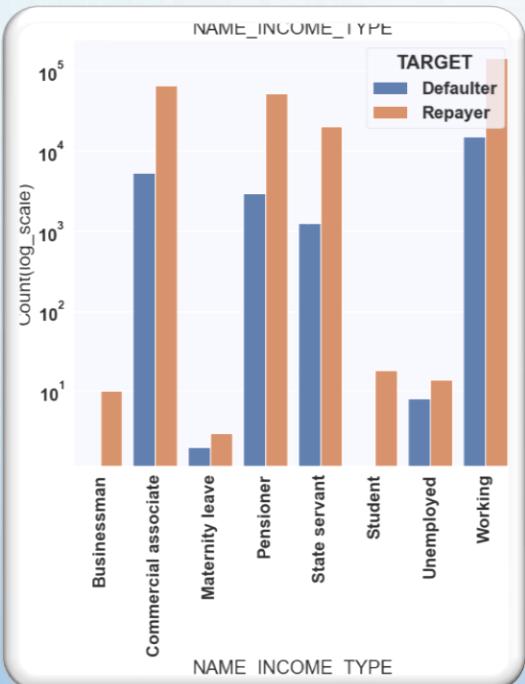
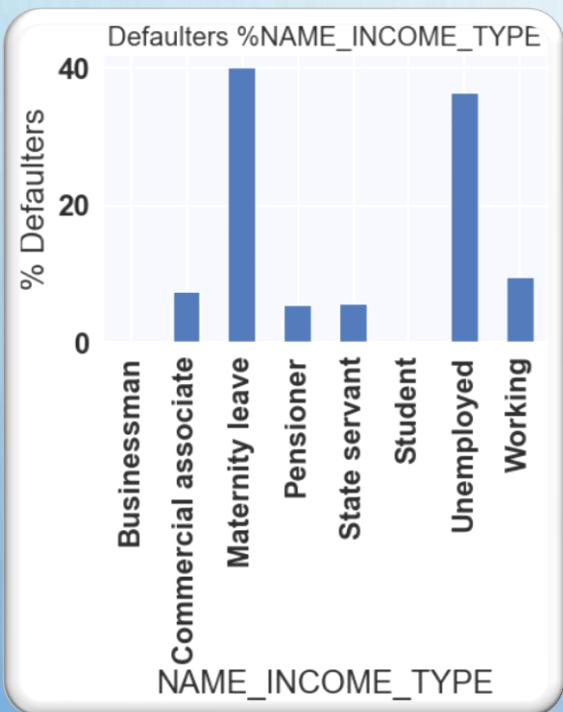
# UNIVARIATE ANALYSIS



## Inferences:

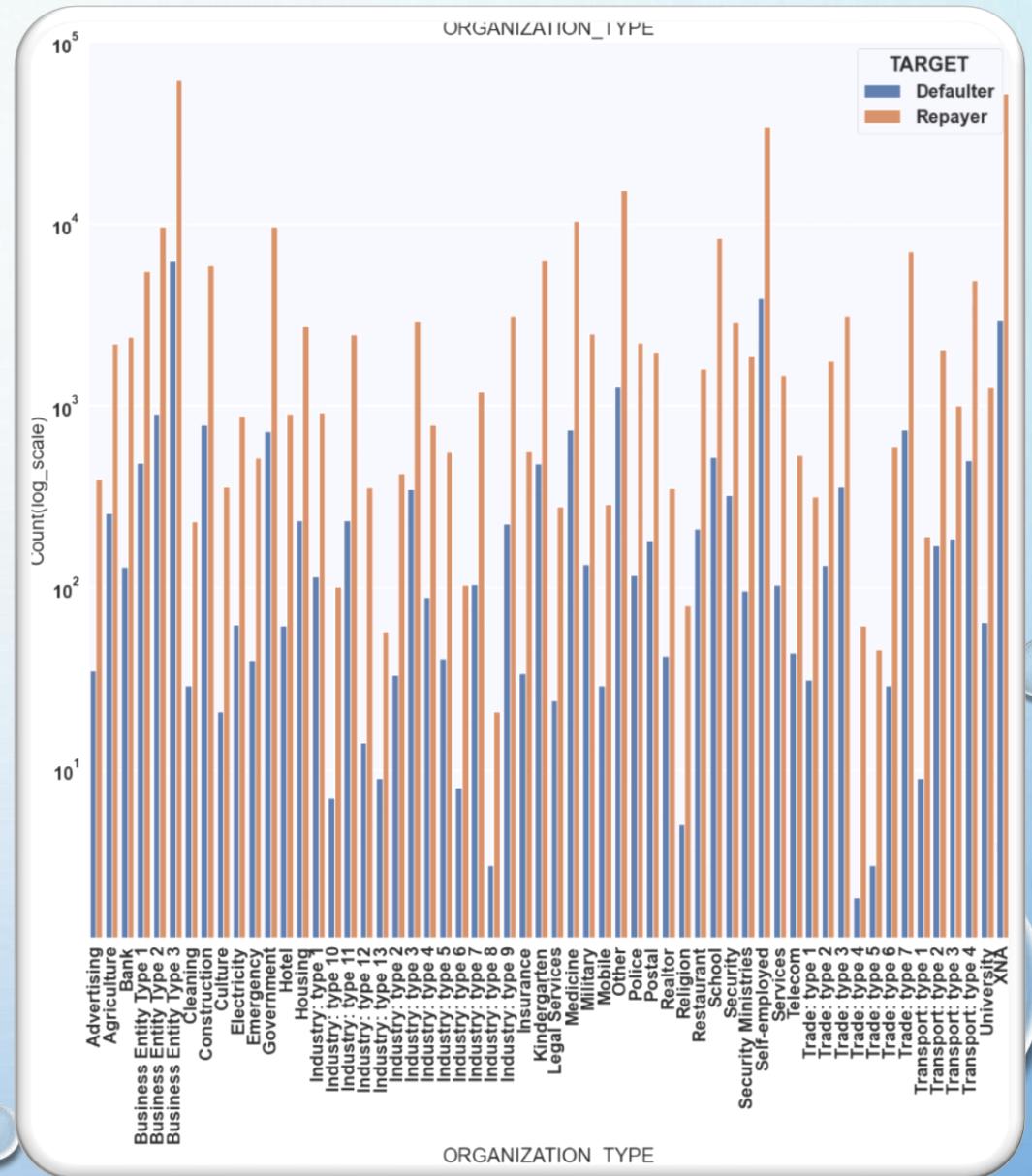
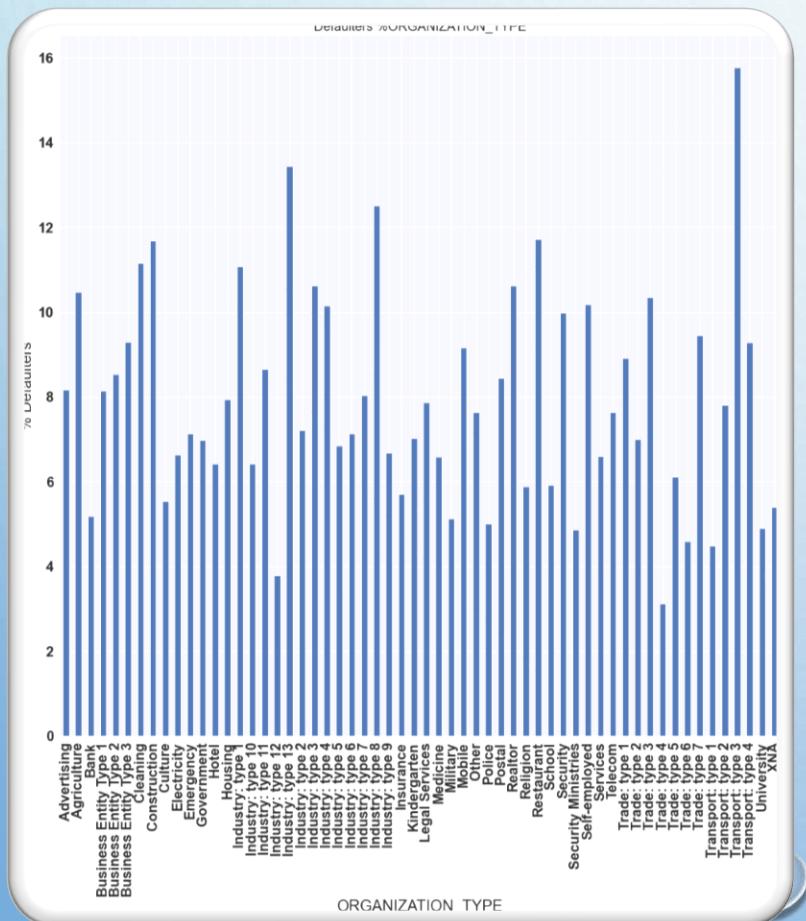
1. Clients with Single/not married and Civil marriage status have highest chances of being defaulter. But Married people has least chances of being defaulter
2. Clients with Lower secondary education have more chances than any client with any other education type and but there are more client with secondary/secondary special education Type

# UNIVARIATE ANALYSIS



**Inferences :** Low skills labourer have highest chance of being defaulter and IT staff have less chances of being applied for loan

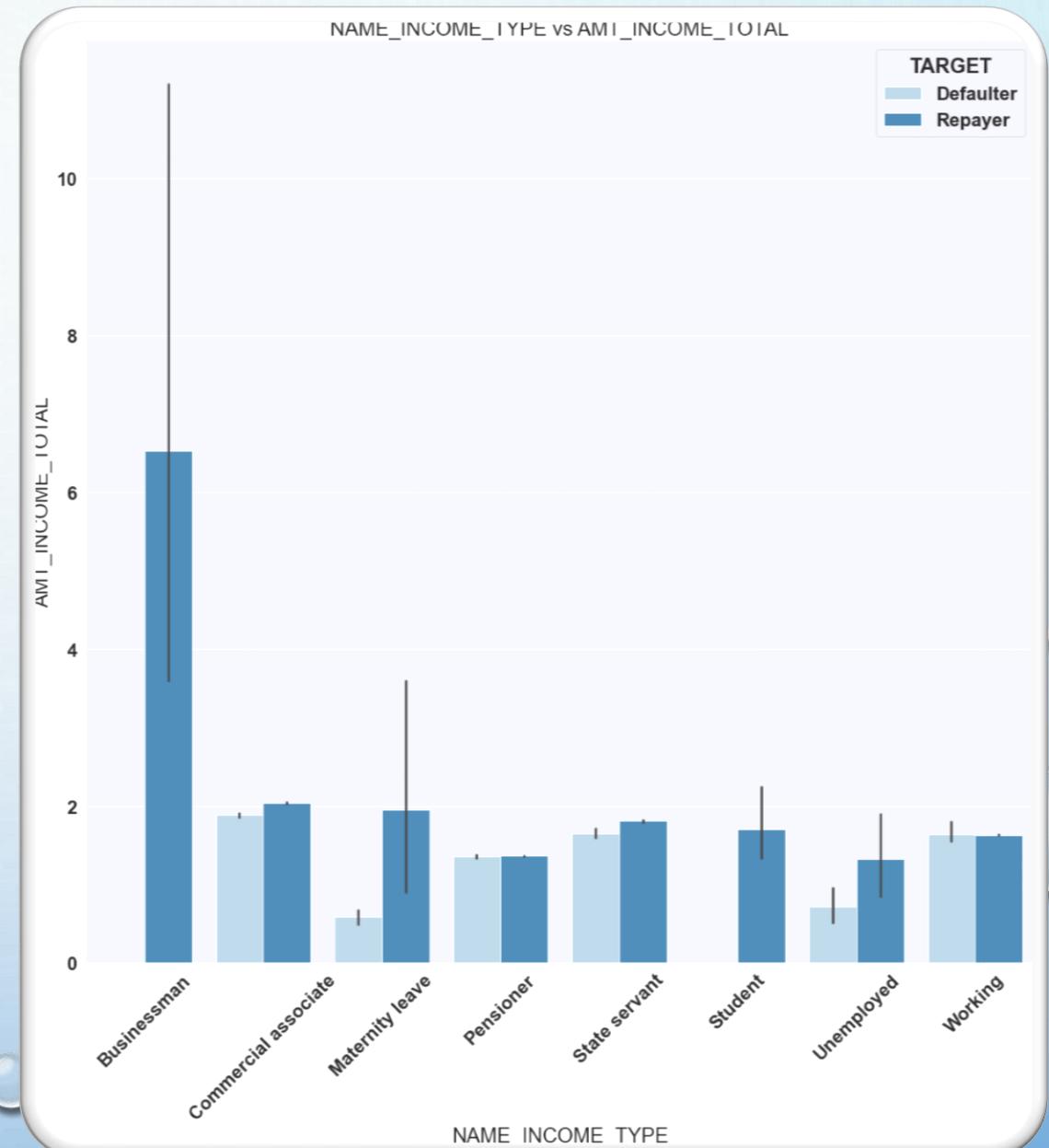
# UNIVARIATE ANALYSIS



**Inferences** : Clients with type 3 are highest defaulter and most application are from XNA and Business entity 3

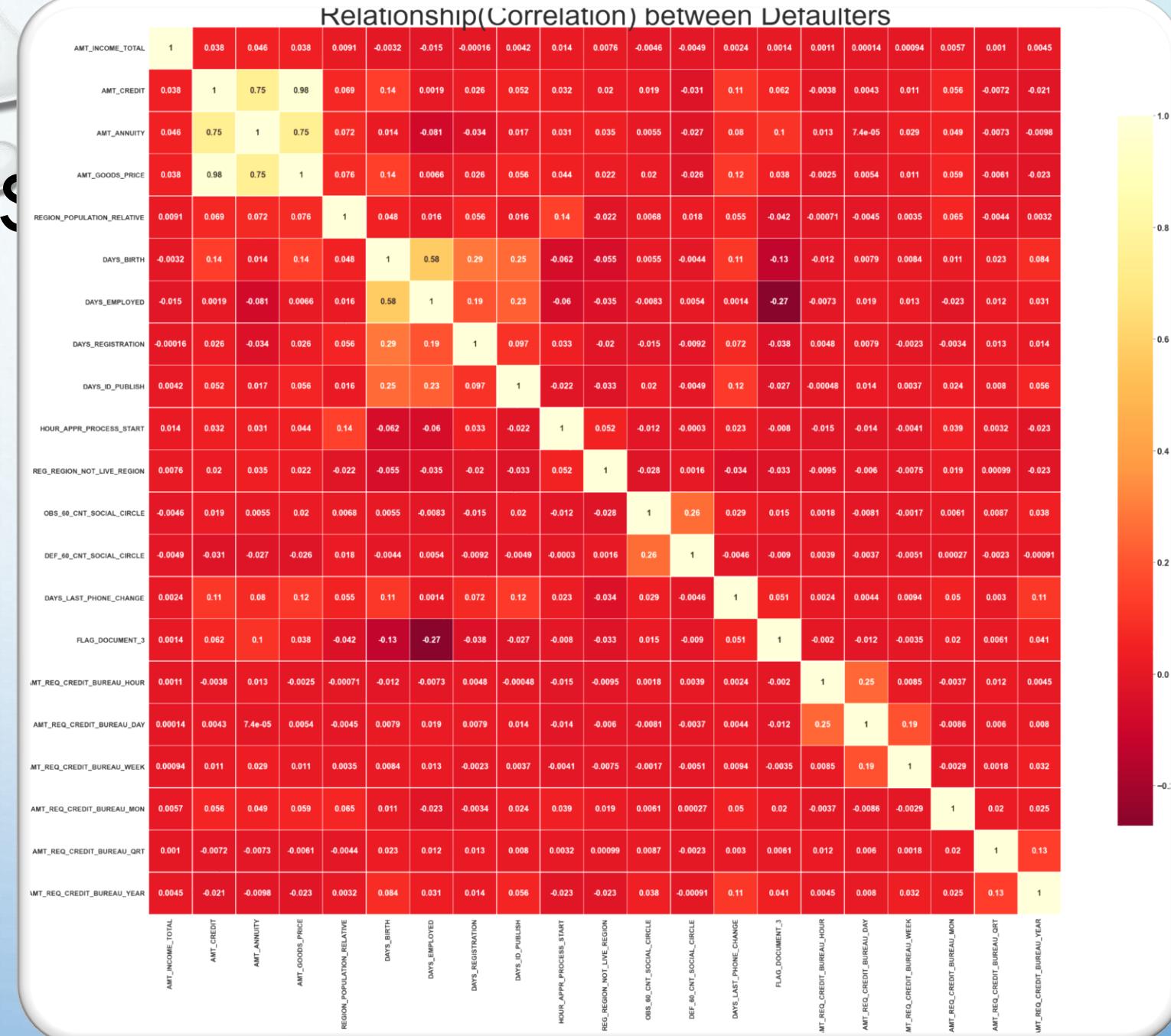
# BIVARIATE ANALYSIS

**Inference:** Businessmen category has highest income with range of ~ 3.5 - 10.5 and more likely to apply for loan



# BIVARIATE ANALYSIS DEFALTERS

**Inferences** : Credit amount of the loan is correlated well with loan annuity and the price of the goods for which the loan is given. Also number of days employed carry more weightage with loan repayment

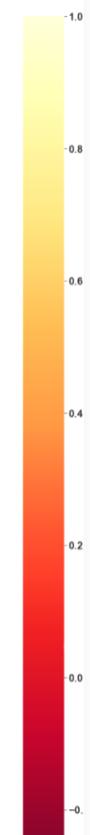


# BIVARIATE ANALYSIS REPAYERS

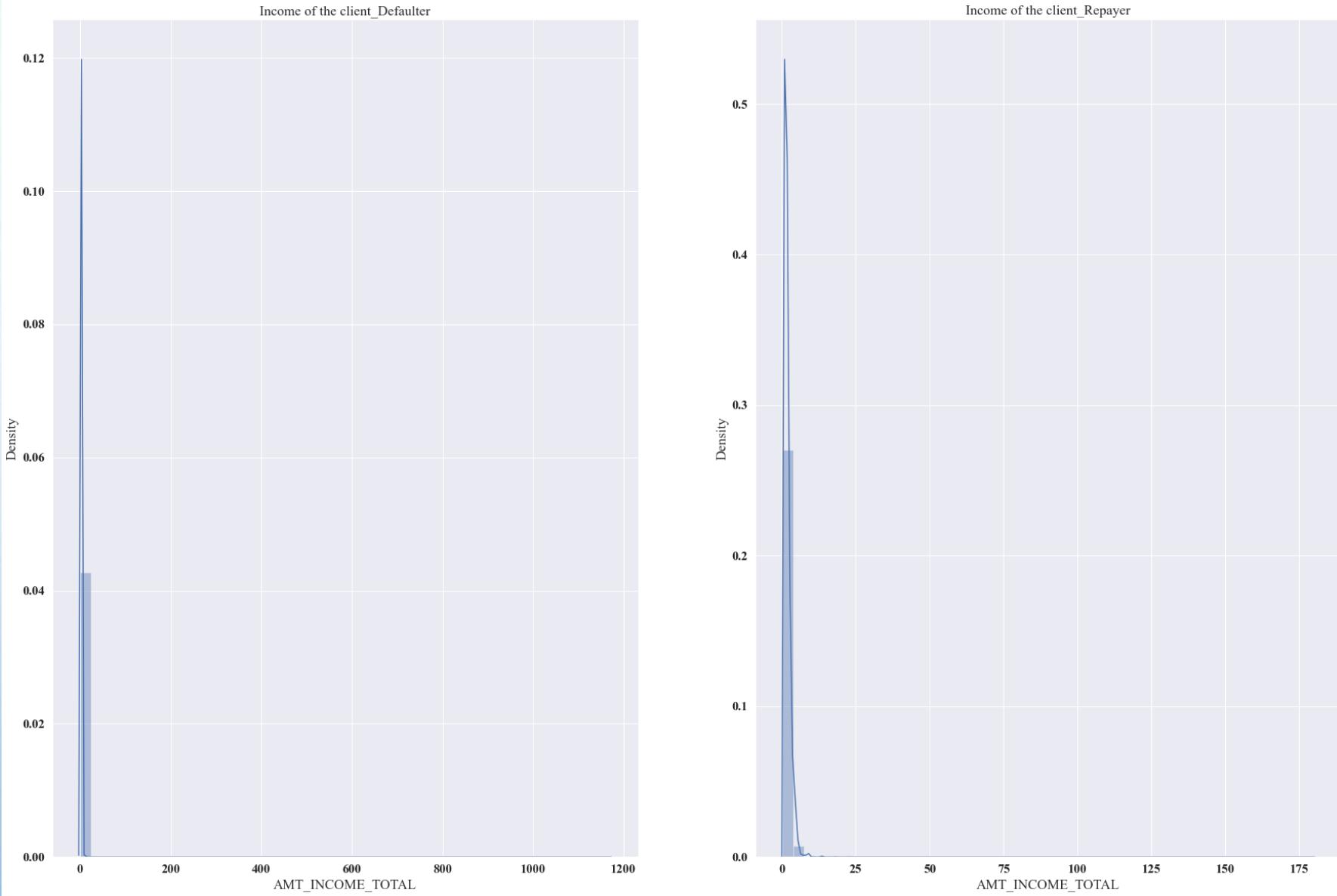
**Inferences :** Credit amount is strongly correlated with loan annuity and the price of the goods for which the loan is given and Total Income. Also repayers have good correlation in number of days client employed

Relationship(Correlation) between Repayers

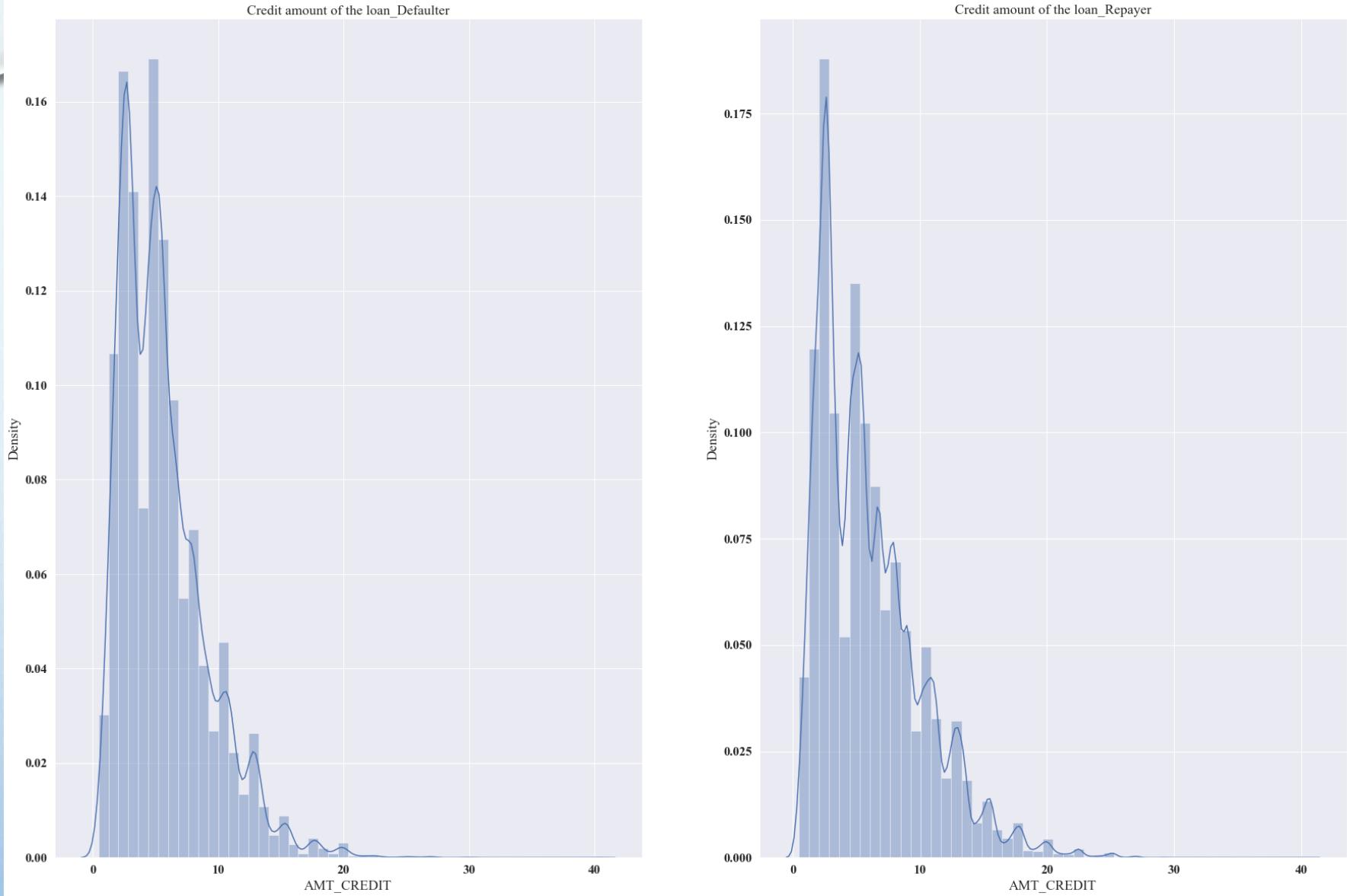
	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE	DAY_BIRTH	DAY_EMPLOYED	DAY_REGISTRATION	DAY_ID_PUBLISH	HOUR_APPR_PROCESS_START	REG_REGION_NOT_LIVE_REGION	OBS_60_CNT_SOCIAL_CIRCLE	DEF_60_CNT_SOCIAL_CIRCLE	DAY_LAST_PHONE_CHANGE	FLAG_DOCUMENT_3	MT_REQ_CREDIT_BUREAU_HOUR	AMT_REQ_CREDIT_BUREAU_DAY	MT_REQ_CREDIT_BUREAU_WEEK	AMT_REQ_CREDIT_BUREAU_MON	AMT_REQ_CREDIT_BUREAU_QRT	AMT_REQ_CREDIT_BUREAU_YEAR
AMT_INCOME_TOTAL	1	0.34	0.42	0.35	0.17	-0.062	-0.14	-0.065	-0.023	0.077	0.069	-0.028	-0.028	0.041	-0.039	0.0027	0.008	0.0086	0.059	0.018	0.034
AMT_CREDIT	0.34	1	0.77	0.99	0.1	0.047	-0.07	-0.013	0.0015	0.054	0.025	-0.00089	-0.022	0.07	0.1	-0.0023	0.0051	0.00094	0.055	0.022	-0.038
AMT_ANNUITY	0.42	0.77	1	0.78	0.12	-0.012	-0.1	-0.039	-0.014	0.054	0.042	-0.013	-0.023	0.062	0.1	0.0032	0.0025	0.012	0.036	0.012	-0.008
AMT_GOODS_PRICE	0.35	0.99	0.78	1	0.1	0.045	-0.069	-0.016	0.0036	0.063	0.027	-0.0072	-0.023	0.071	0.079	-0.0017	0.0055	0.0012	0.057	0.022	-0.04
REGION_POPULATION_RELATIVE	0.17	0.1	0.12	0.1	1	0.025	-0.0072	0.052	0.0011	0.17	0.0043	-0.012	0.0023	0.041	-0.086	-0.0023	0.0016	-0.0028	0.071	-0.002	0.00015
DAY_BIRTH	-0.062	0.047	-0.012	0.045	0.025	1	0.63	0.33	0.27	-0.096	-0.066	-0.00673	0.001	0.076	-0.1	-0.0029	-0.0016	0.0037	0.002	0.015	0.073
DAY_EMPLOYED	-0.14	-0.07	-0.1	-0.069	-0.0072	0.63	1	0.21	0.28	-0.095	-0.038	0.0075	0.016	-0.023	-0.24	-0.0043	-0.00093	0.0017	-0.033	0.013	0.047
DAY_REGISTRATION	-0.065	-0.013	-0.039	-0.016	0.052	0.33	0.21	1	0.1	0.008	-0.029	-0.0082	-0.0027	0.054	-0.032	0.0025	9.3e-06	0.0013	0.011	0.00036	0.024
DAY_ID_PUBLISH	-0.023	0.0015	-0.014	0.0036	0.0011	0.27	0.28	0.1	1	-0.034	-0.035	0.013	-0.0025	0.083	-0.05	-0.0019	0.0022	0.0069	0.017	0.017	0.048
HOUR_APPR_PROCESS_START	0.077	0.054	0.054	0.063	0.17	-0.096	-0.095	0.008	-0.034	1	0.055	-0.008	-0.0088	0.013	-0.013	-0.014	0.0039	-0.0015	0.036	0.0012	-0.025
REG_REGION_NOT_LIVE_REGION	0.069	0.025	0.042	0.027	0.0043	-0.066	-0.038	-0.029	-0.035	0.055	1	-0.02	-0.009	-0.038	-0.034	-0.0016	-0.0012	0.00078	-0.0029	-0.004	-0.018
OBS_60_CNT_SOCIAL_CIRCLE	-0.028	-0.00089	-0.013	-0.00072	-0.012	-0.0073	0.0075	-0.0082	0.013	-0.008	-0.02	1	0.25	0.015	0.027	0.00058	-0.0017	0.001	0.0025	0.0047	0.032
DEF_60_CNT_SOCIAL_CIRCLE	-0.028	-0.022	-0.023	-0.023	0.0023	0.001	0.016	-0.0027	-0.0025	-0.0088	0.009	0.25	1	0.00016	0.012	-0.002	-0.0016	-0.0022	-0.0014	-0.0032	0.016
DAY_LAST_PHONE_CHANGE	0.041	0.07	0.062	0.071	0.041	0.076	-0.023	0.054	0.083	0.013	-0.038	0.015	0.00016	1	0.065	0.0028	-0.00081	0.0074	0.045	0.01	0.12
FLAG_DOCUMENT_3	-0.039	0.1	0.1	0.079	-0.086	-0.1	-0.24	-0.032	-0.05	-0.013	-0.034	0.027	0.012	0.065	1	-0.0025	0.0021	0.0093	0.011	0.011	0.046
MT_REQ_CREDIT_BUREAU_HOUR	0.0027	-0.0023	0.0032	-0.0017	-0.0023	-0.0029	-0.0043	0.0025	-0.0019	-0.014	-0.0016	0.00058	-0.002	0.0028	-0.00025	1	0.23	0.0062	0.0034	-3.7e-05	-2.8e-06
AMT_REQ_CREDIT_BUREAU_DAY	0.008	0.0051	0.0025	0.0055	0.0016	-0.0016	-0.00093	9.3e-06	0.0022	0.0039	-0.0012	-0.0017	-0.0016	-0.00081	0.0021	0.23	1	0.22	-0.0024	-0.002	0.00011
MT_REQ_CREDIT_BUREAU_WEEK	0.0086	0.00094	0.012	0.0012	-0.0028	0.0037	0.0017	0.0013	0.0069	-0.0015	0.00078	0.001	-0.0022	0.0074	0.0093	0.0062	0.22	1	-0.0078	-0.0081	0.029
AMT_REQ_CREDIT_BUREAU_MON	0.059	0.055	0.036	0.057	0.071	0.002	-0.033	0.011	0.017	0.036	-0.0029	0.0025	-0.0014	0.045	0.011	0.0034	-0.0024	-0.0078	1	0.0045	0.013
AMT_REQ_CREDIT_BUREAU_QRT	0.018	0.022	0.012	0.022	-0.002	0.015	0.013	0.00036	0.017	0.0012	-0.004	0.0047	-0.00032	0.01	0.011	-3.7e-05	-0.002	-0.0081	0.0045	1	0.093
AMT_REQ_CREDIT_BUREAU_YEAR	0.034	-0.038	-0.008	-0.04	0.00015	0.073	0.047	0.024	0.048	-0.025	-0.018	0.032	0.016	0.12	0.046	-2.8e-06	0.00011	0.029	0.013	0.093	1



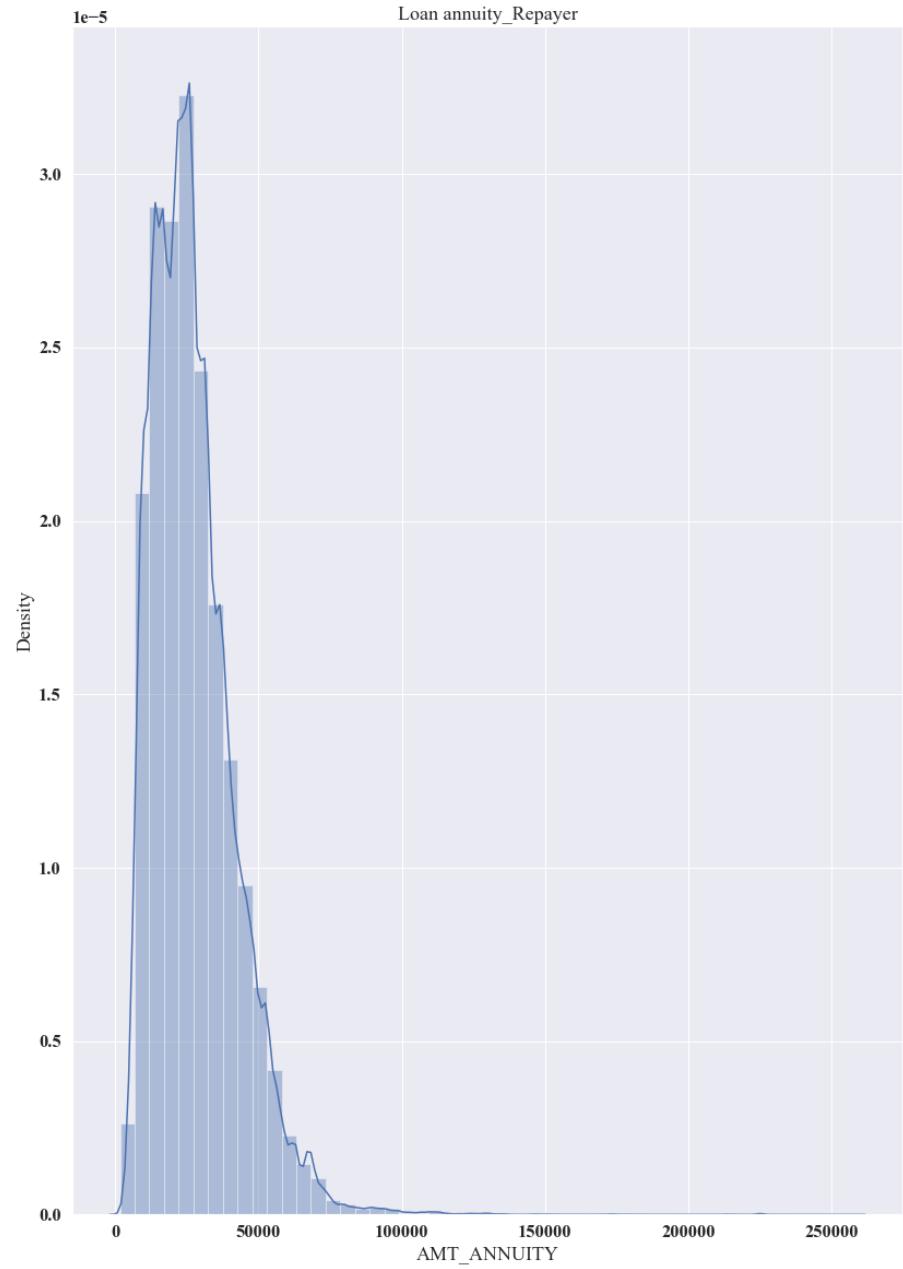
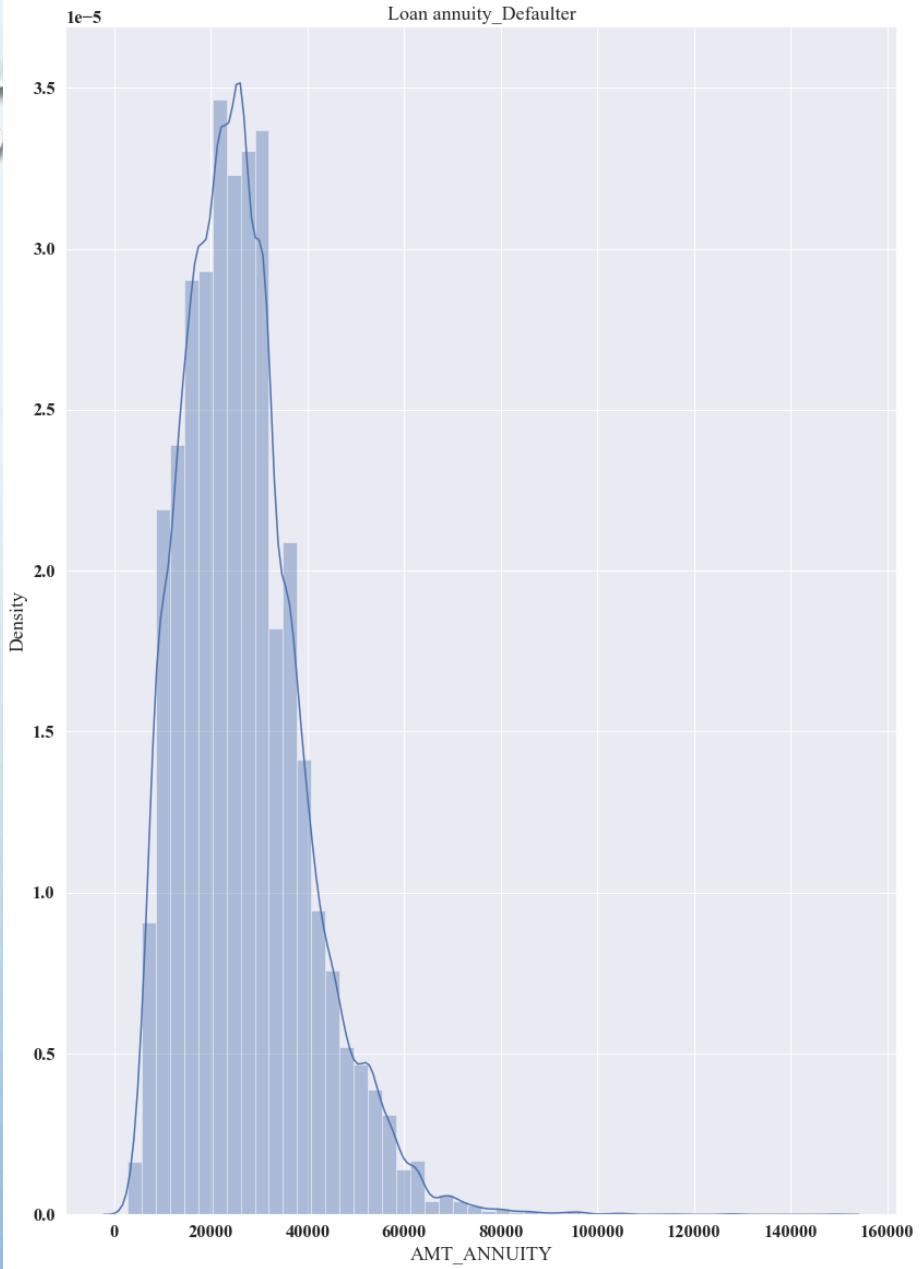
# BIVARIATE ANALYSIS



# BIVARIATE ANALYSIS



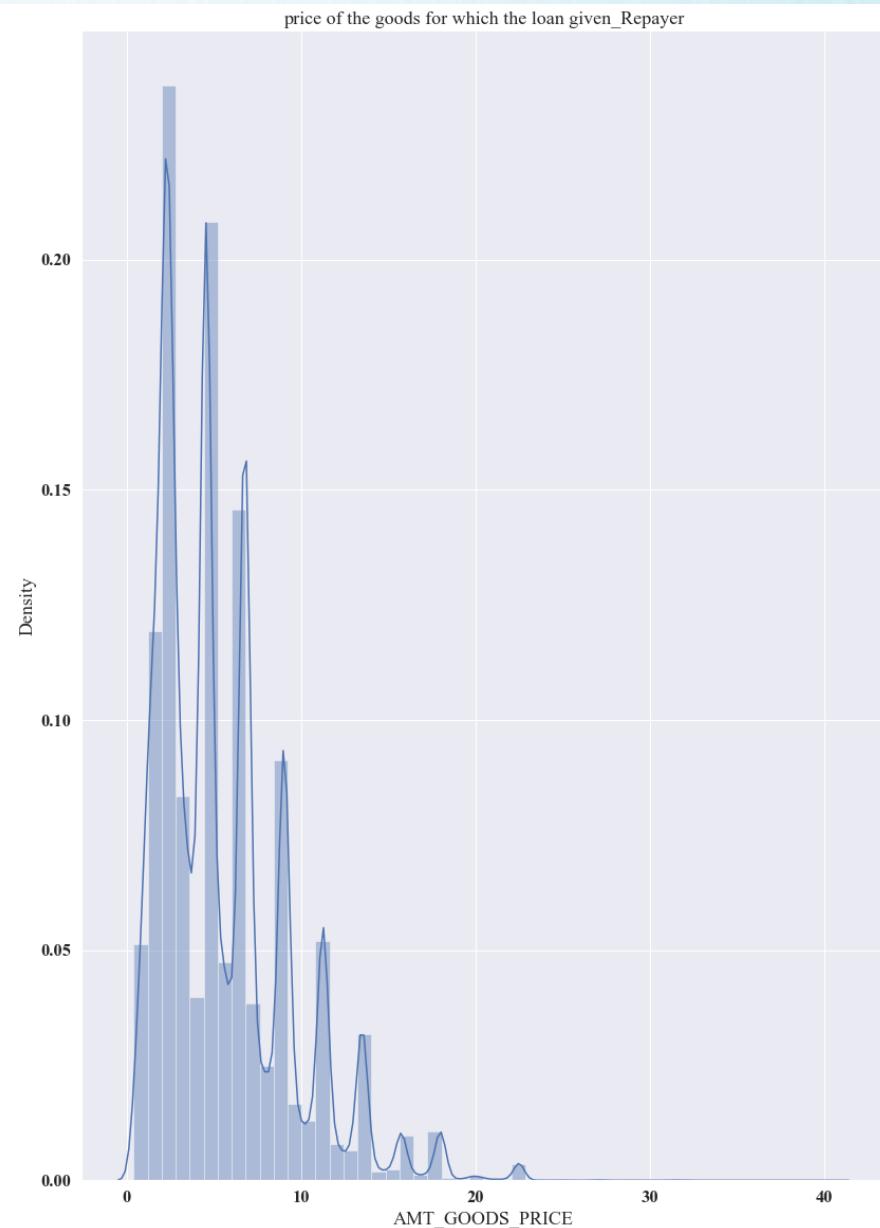
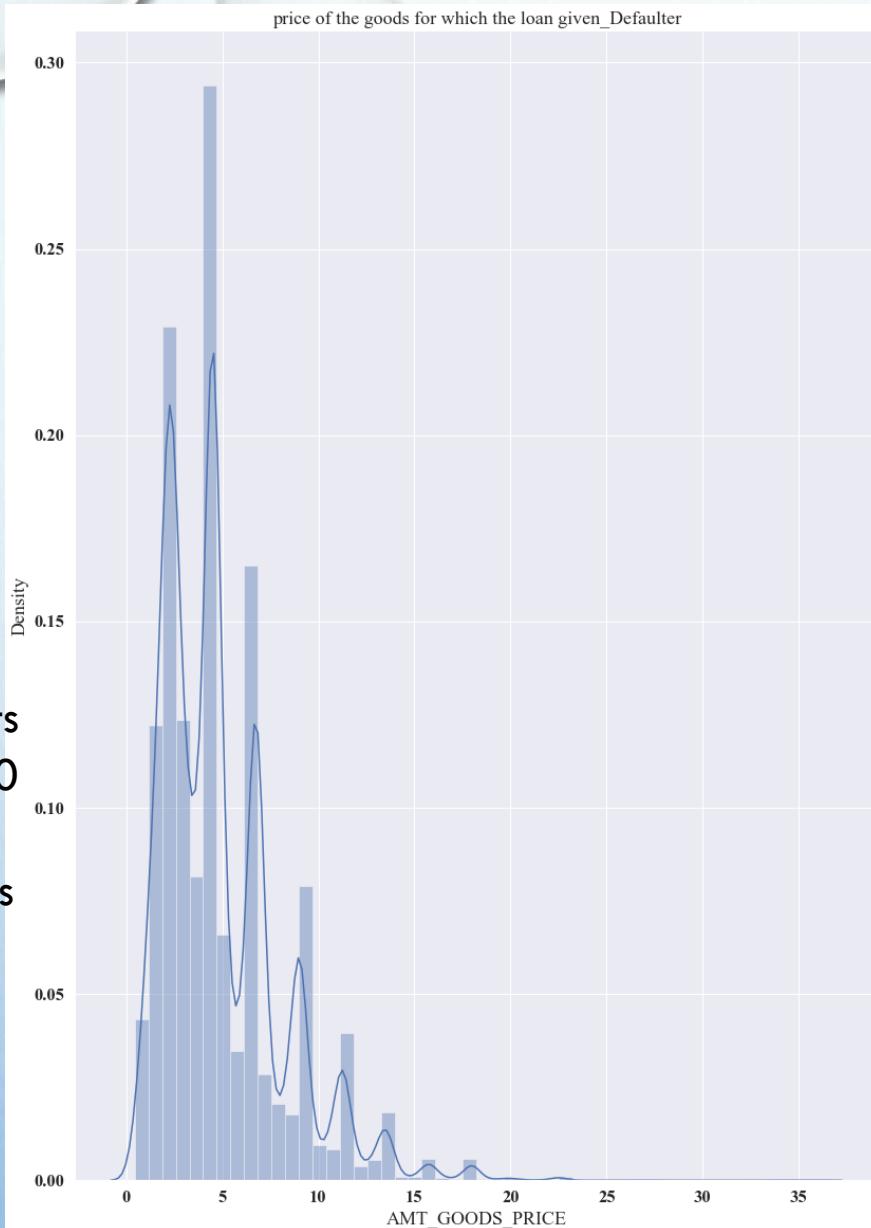
# BIVARIATE ANALYSIS REPAYERS



# BIVARIATE ANALYSIS

## Inferences :

- 1) Credit amount given to clients is approximately less than 10 lakhs
- 2) The repayers and defaulters shows similar distribution for all above four graphs and hence it may not suitable to make judgement based on this single varibale



# BIVARIATE ANALYSIS

## Inferences :

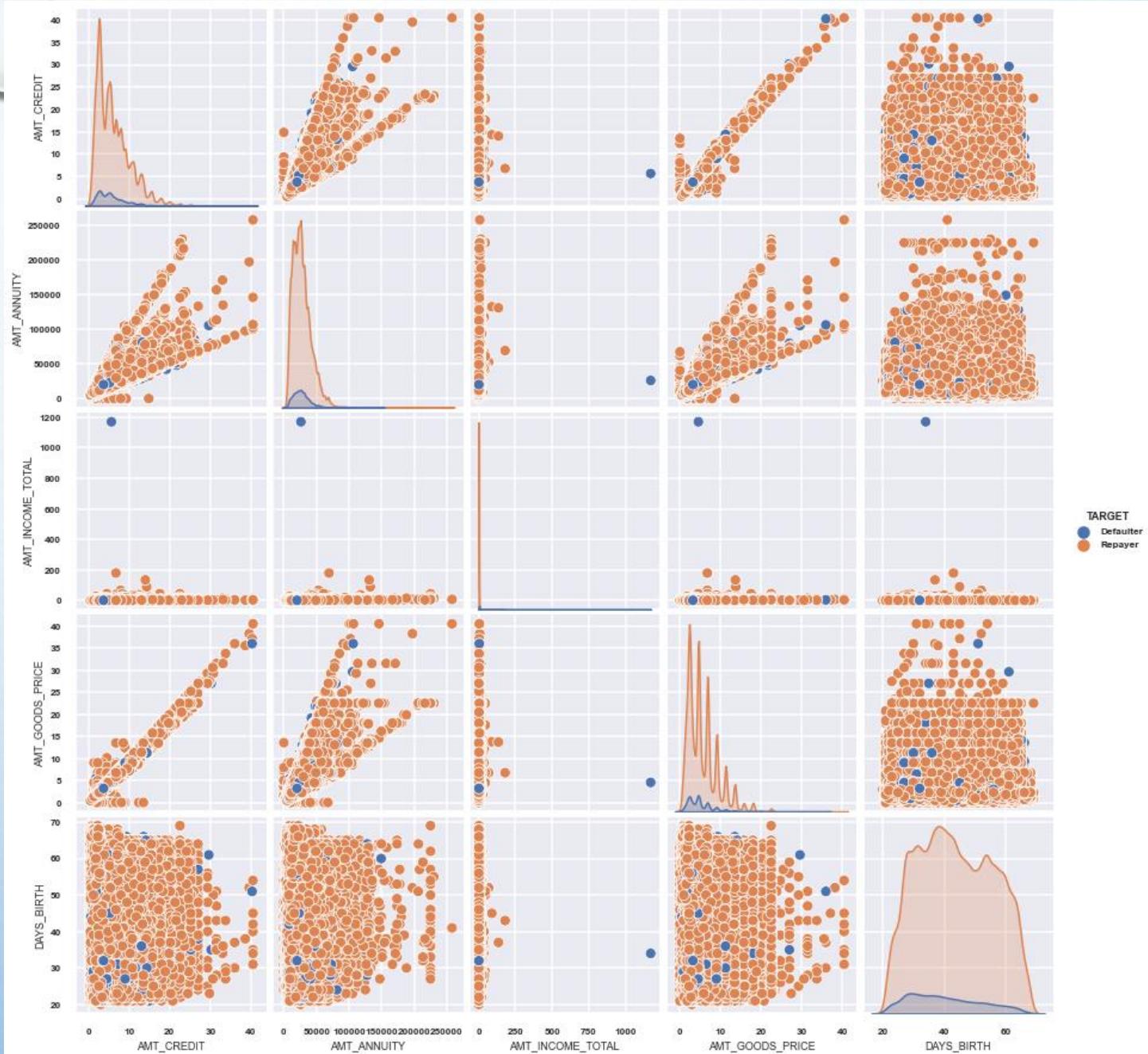
With loan amount more than 30 lakh the defaulter people increases



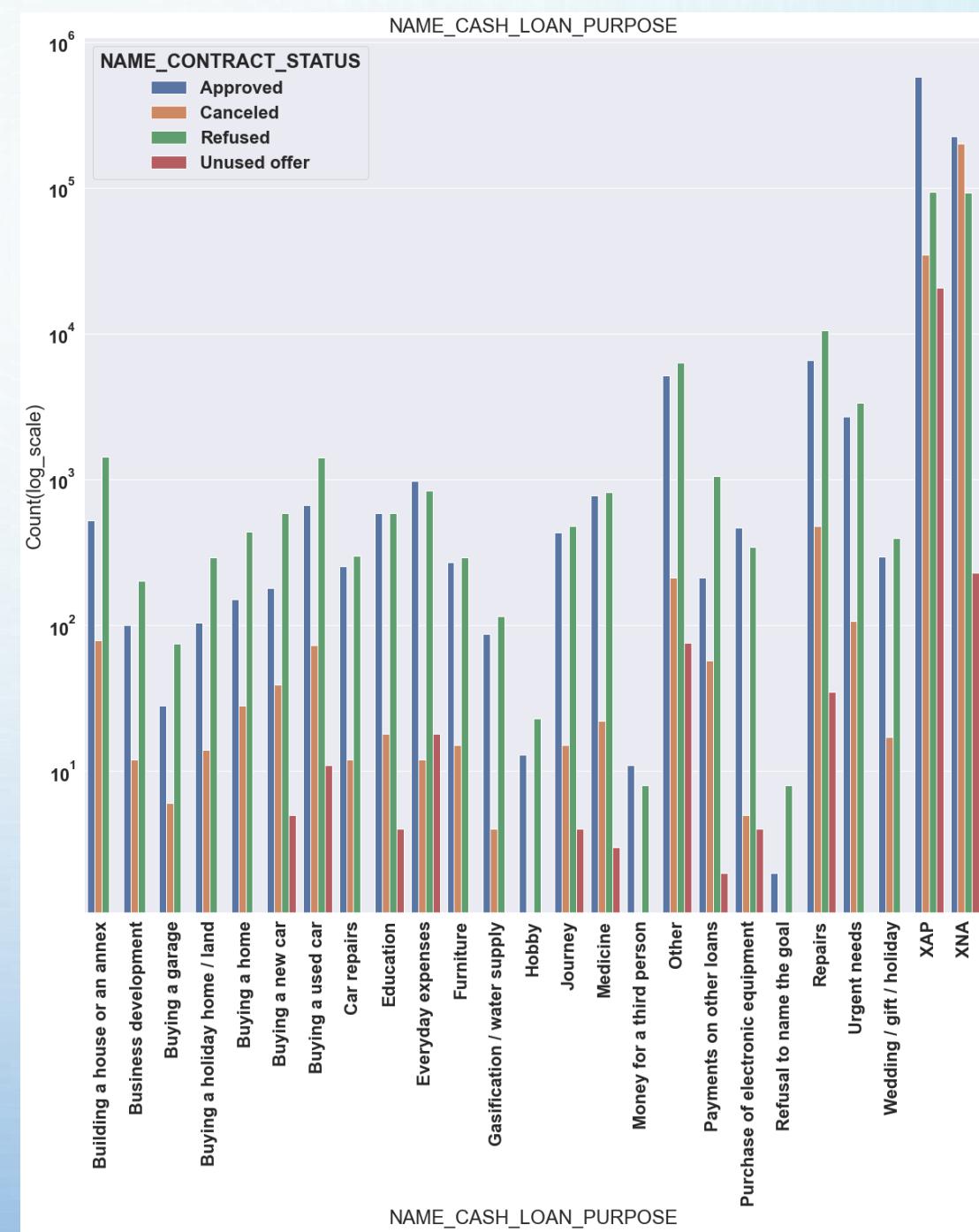
# BIVARIATE ANALYSIS

## Inferences :

1. Loan Amount and Goods price shows positive correlation
2. AMT\_INCOME\_TOTAL does not give any clear indication about defaulters and repayers



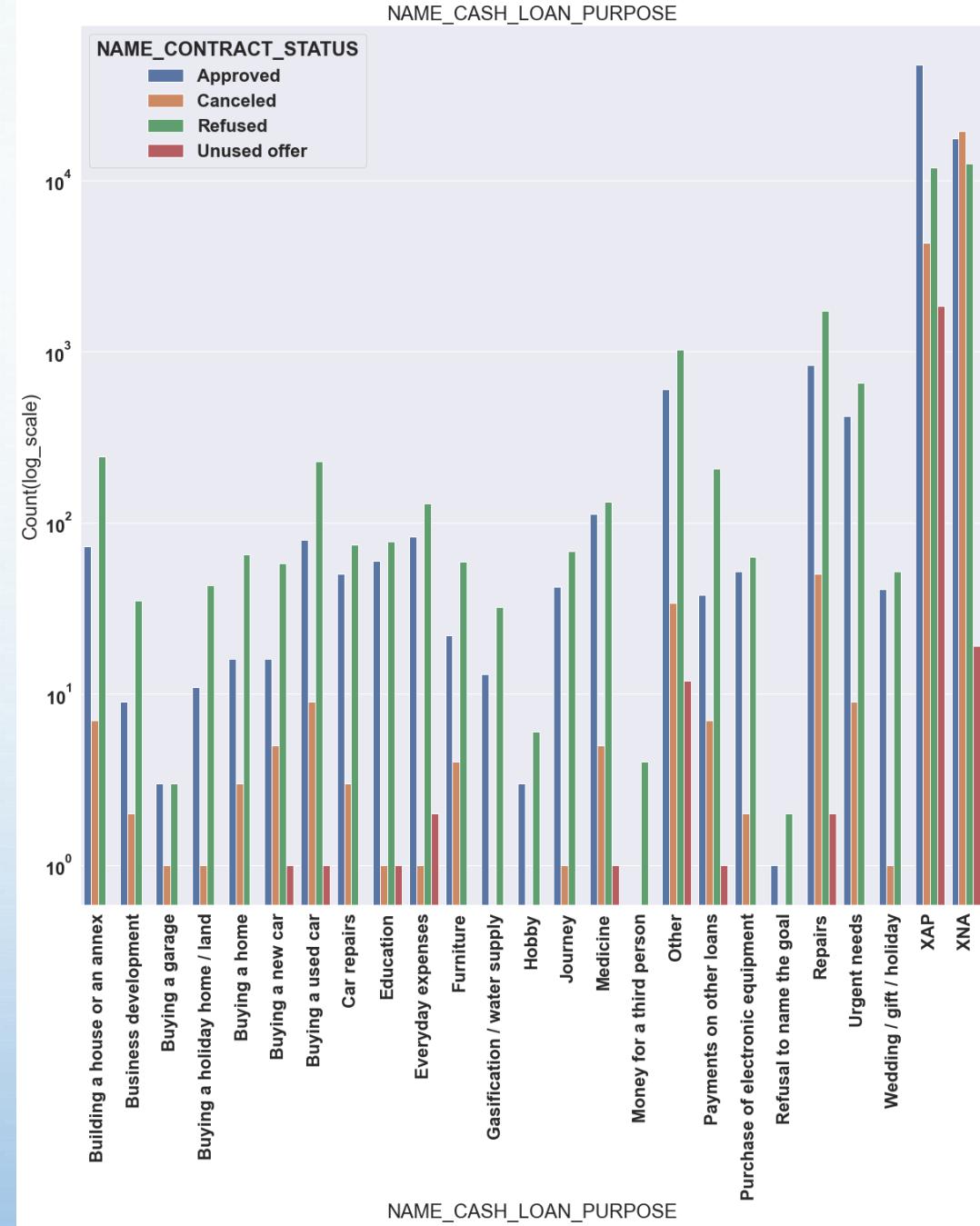
# BIVARIATE ANALYSIS



# BIVARIATE ANALYSIS

## Inferences :

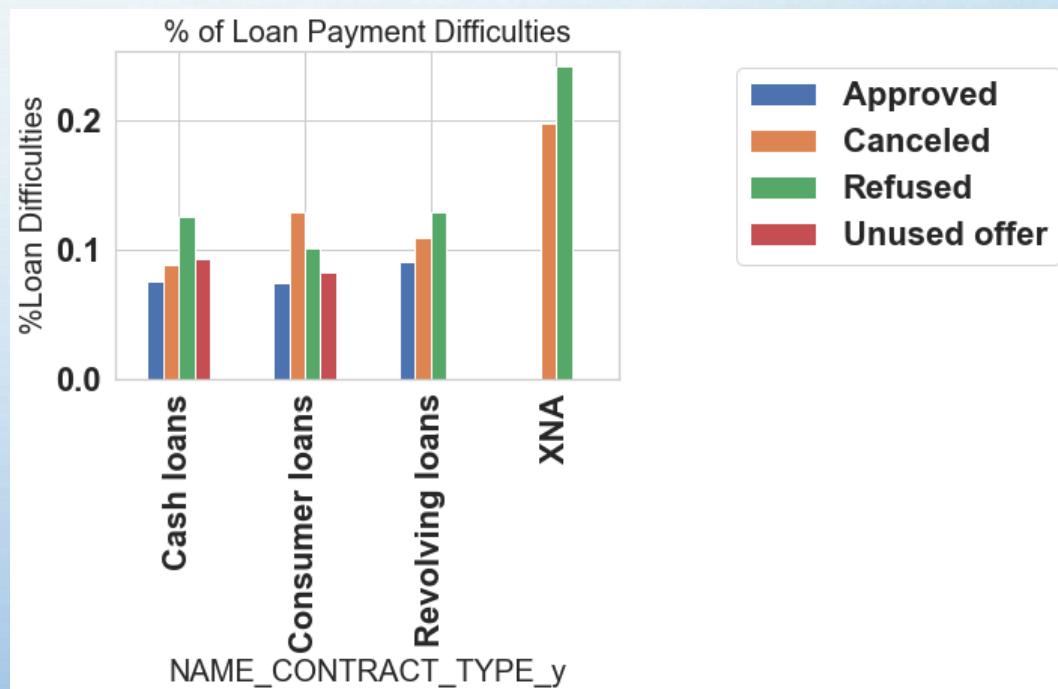
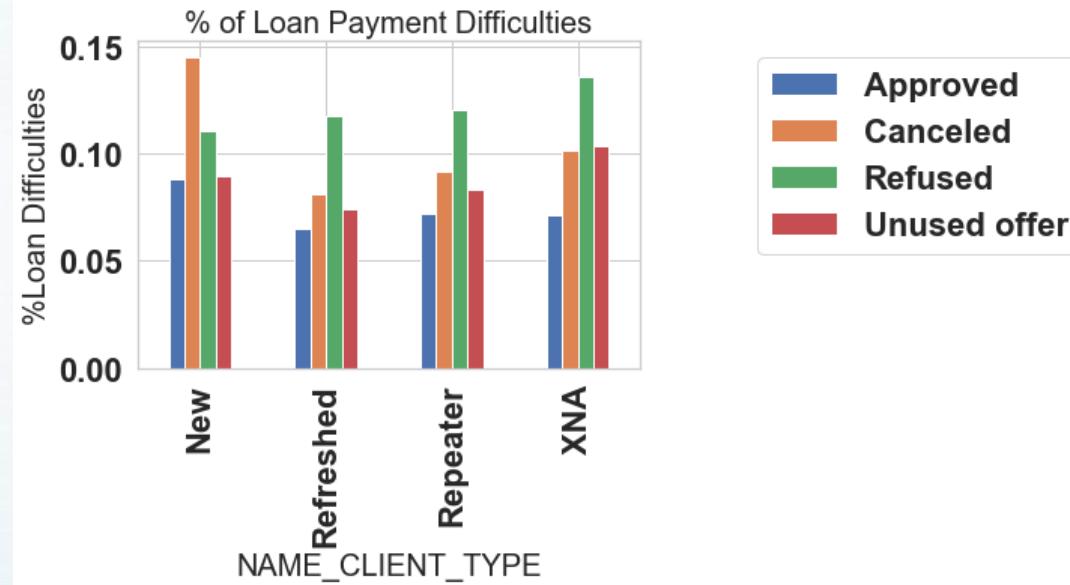
- 1) Purpose of the cash loan for repair and other has substantial applicant and in this category more applications are rejected
- 2) Also XNA and XAP which are the purpose of loan have the most application for



# BIVARIATE ANALYSIS REPAYERS

## Inferences :

1. Client who are "NEW and could secure loan application previously ( 'Cancelled')tend to be more % of Loan-Payment Difficulties
2. Clients with 'Revolving loans' and also 'Refused' status previously have more difficulties in % of Loan-Payment Difficulties. More over no conclusions can be drawn for XNA loan application which is unknown



# CONCLUSION

## FACTORS WHICH PREDICT THE CLIENT CAN BE DEFULTER

1. LOWER SECONDARY' IN 'NAME\_EDUCATION\_TYPE' HAS MORE DIFFICULTY IN REPAYING LOANS
2. LOW SKILLED LABORERS' IN 'OCCUPATION\_TYPE' ALSO HAS MAXIMUM DIFFICULTY IN LOAN REPAYMENT
3. 'MATERNITY LEAVE' IN 'NAME\_INCOME\_TYPE' ALSO HAS MAXIMUM DIFFICULTY IN LOAN REPAYMENT.
4. 0-5 YEARS OF EMPLOYMENT ARE MORE LIKELY TO DEFAULT LOAN
5. AMT\_INCOME\_TOTAL:HIGHER INCOME PEOPLE LESS LIKELY TO DEFULTERS
6. CNT\_CHILDREN: MORE NUMBER OF CHILDREN >3 ARE MORE LIKELY TO DEFULTERS

# CONCLUSION

## **FACTORS WHICH PREDICT THE CLIENT CAN BE A POTENTIAL REPAYER AND CAN BE HELPFUL FOR BUSINESS PROFIT**

1. PEOPLE WHO ARE STUDENT AND BUSINESSMEN ARE MOST LIKELY TO REPAY LOANS
2. LOANS WHICH ARE PERSONAL IN NATURE LIKE HOBBY ARE LIKELY TO BE REPAY
3. PEOPLE WHO ARE TOWARDS THEIR RETIREMENT I.E. AGE > 50 HAVE GOOD TRACK RECORD OF REPAYMENT
4. WOMEN HAVE HIGHER CHANCES OF REPAYMENT OF LOAN
5. MARRIAGE ALSO IS DRIVING FACTOR WHILE DECIDING LOAN REPAYMENT I.E. APPLICANT WHO ARE SINGLE AND CIVIL MARRAIGE ARE MORE LIKELY TO REPAY LOAN
6. IF PRICE OF THE GOODS FOR WHICH THE LOAN GIVEN IS LESS THAN ~3 LAKH THERE ARE HIGHER CHANCES GETTING REPAYMENT

# CONCLUSION

## FACTORS WHICH CAN BE HELPFUL TO OFFER THE LOAN BUT REDUCE ITS AMOUNT & AT HIGHER INTEREST RATES

1. BASED ON NUMBER OF FAMILY MEMBERS (I.E. IF NUMBER IS LARGE) REDUCE AMOUNT OF LOAN THAN ELIGIBLE (SAY SOME % OF TAKE HOME SALARY) CAN BE OFFERED AT HIGHER INTEREST RATES
2. PEOPLE WHO HAVE GOODS PRICE AND LESS INCOME CAN BE POTENTIAL CUSTOMER FOR REDUCE LOAN AND HIGHER INTEREST RATES
3. IF CLIENT REFUSED TO MENTION THE PURPOSE OF LOAN CAN BE TARGETED WITH REDUCED LOAN AMOUNT AND HIGHER INTEREST RATE
4. CLIENT WHO HAVE REFUSED LOAN PREVIOUSLY AND REFUSED PURPOSE CAN BE TARGETED WITH EXTRA SECURITY ASSURANCE TO INCREASE THE BUSINESS

THANK YOU