

# Project Name: **Credit Risk Analysis**

## **Project Objective**

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 the loan, lending (to risky applicants) at a higher interest rate, etc.

## **Approach/ Steps to be followed**

1. Upload
2. remove heavy null columns (>50%)
3. Go through columns\_description and develop data understanding.
4. Start individual column analysis and Cleaning if necessary
5. Merge the dataframes
6. Identify Object and Numeric data of dataframe (Created 2 different dataframes to clearly identify)
7. Identify categorical columns (having  $\leq 10$  unique values) and their cleaning in the merged dataframe.
8. Segmenting the Merged dataframe based on Target group ( 0 or 1)
9. Running Univariate / Bivariate / multivariate analysis to figure out correlation
10. Conclusion : Summarise the outcome with Correlations identified.

Note : to facilitate the ease of usage, common functions are created before step 3.  
Calling it **step-Functions**

Step 1 : Usual upload

Step 2 : There are many columns with too much null values.. we have 41 dropped cols with null values from current Applications, while such 4 columns from Previous applications data were dropped. Threshold considered for deletion : 50%

Step 3 : Study of Columns Description :

This step was the most time consuming as data understanding was built from here.

Step 4 : Individual columns analysis (Univariate) and cleaning if necessary:

Approach:

Parameters to check :

- How many null values ?
- does column have any unexpected entry ? Definition of unexpected would be different for each column
- what inference could be drawn from column ?
- Finally - what could be done with this column ?

Outcome:

- Null handling: For Numerical columns, nulls are mostly replaced with 0 or mean value depending on column nature.
- Outliers handling: For Numerical columns, Outliers were handled in some cases by replacing with Upper limit based on IQR.
- Categorical columns were cleaned separately after merge.

Step 5 : Merge:

The two data frames were merged after both data frames were individually cleaned.

Step 6 : Identify Object and Numeric data of dataframe

While keeping the first 2 columns as common, numeric and object columns are separated.

Step 7 : Identify categorical columns

A total of 23 columns were filtered out first. Their nulls were searched and found that 21 out of them have less than 2% nulls, which were replaced by Mode treatment. Other 2 columns were heavily loaded on nulls (~49%) which were dropped.

## Step 8 : Segmenting the Merged dataframe based on Target group ( 0 or 1)

2 separate data frames created from merged and cleaned dataframe. This would help us run bivariate analysis as now – we have a totally cleaned data.

In this step – we have taken care of cleaning required in categorical columns as well.

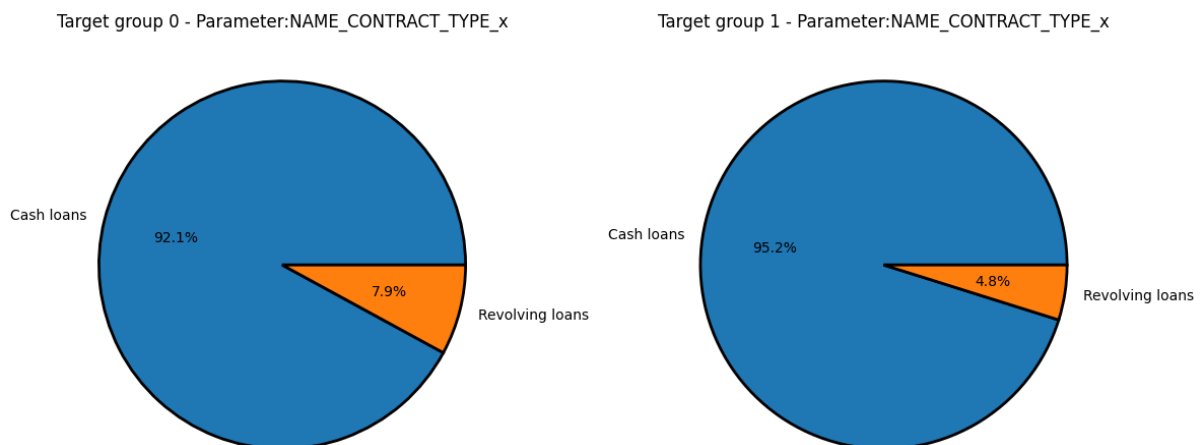
## Step 9 : Running Uni/Bivariate / multivariate analysis to figure out correlation

Following type of tools were used to identify the correlations :

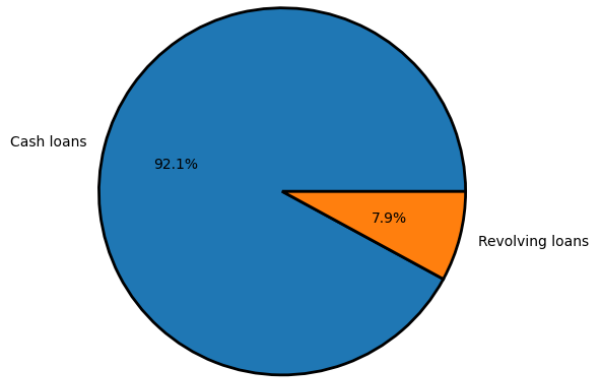
- Pie plots – Target group wise
- Heat map – On numerical data columns
- Scatter Plots – On numerical data columns
- Bar plots - Combination of Numerical and categorical data : Target group wise

Here are some of the snippets observed during the project study:

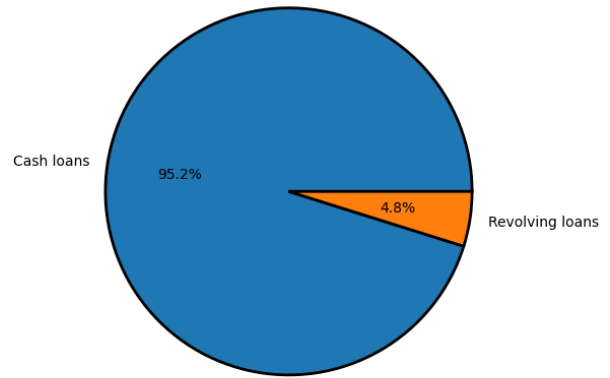
### Pie Plots:



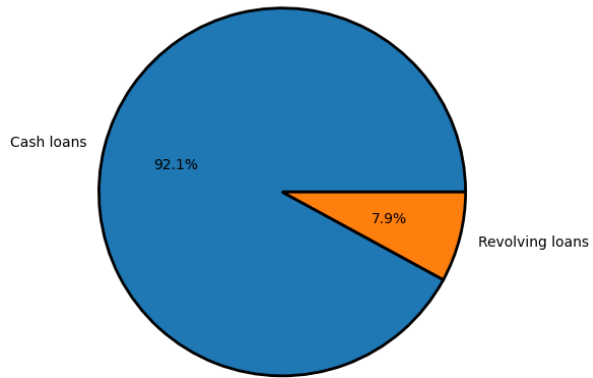
Target group 0 - Parameter:NAME\_CONTRACT\_TYPE\_x



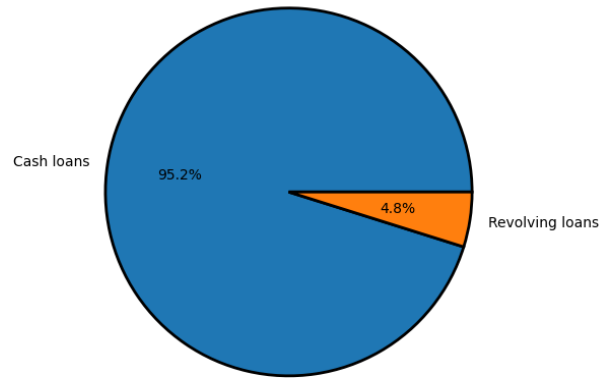
Target group 1 - Parameter:NAME\_CONTRACT\_TYPE\_x



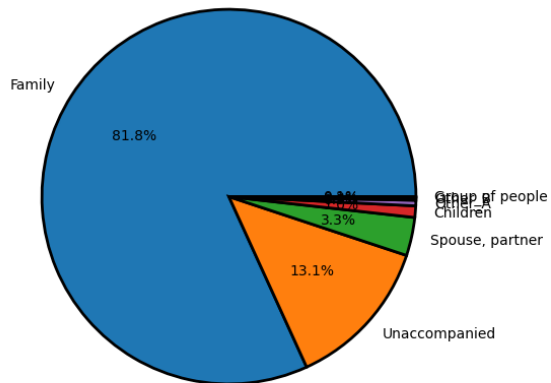
Target group 0 - Parameter:NAME\_CONTRACT\_TYPE\_x



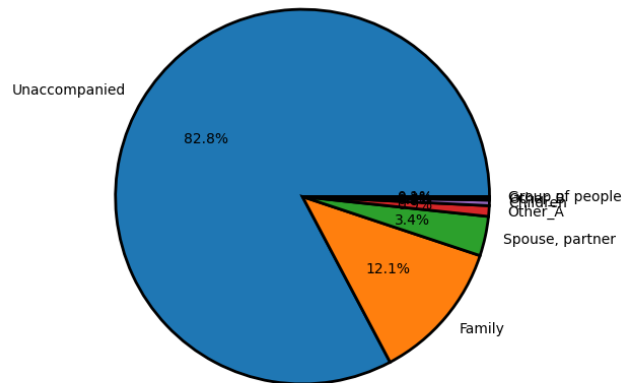
Target group 1 - Parameter:NAME\_CONTRACT\_TYPE\_x



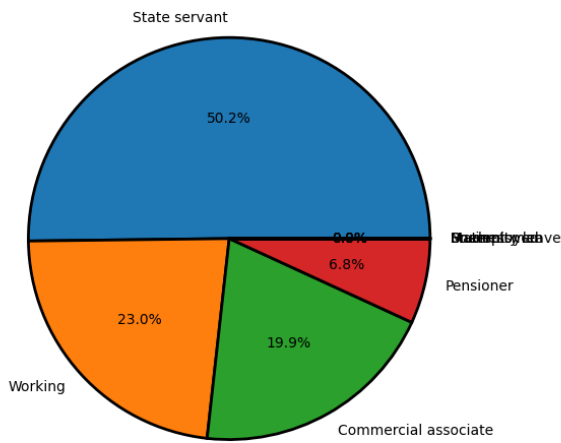
Target group 0 - Parameter:NAME\_TYPE\_SUITE\_x



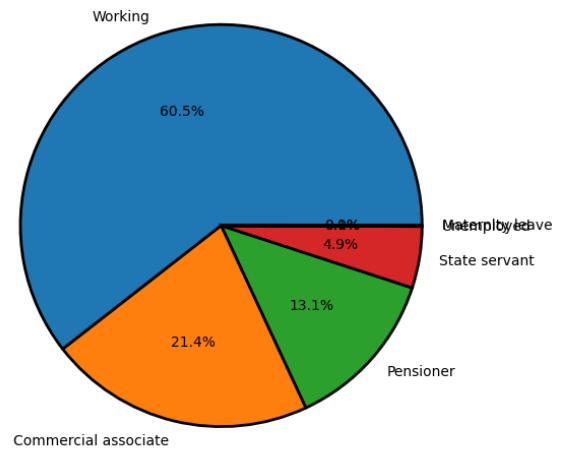
Target group 1 - Parameter:NAME\_TYPE\_SUITE\_x



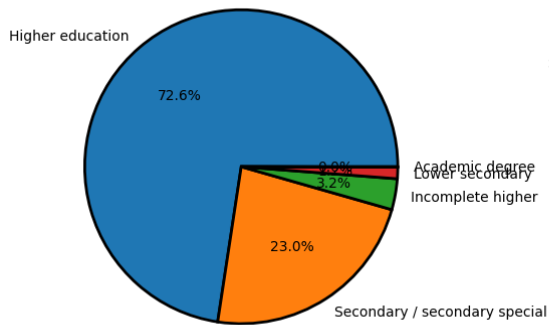
Target group 0 - Parameter:NAME\_INCOME\_TYPE



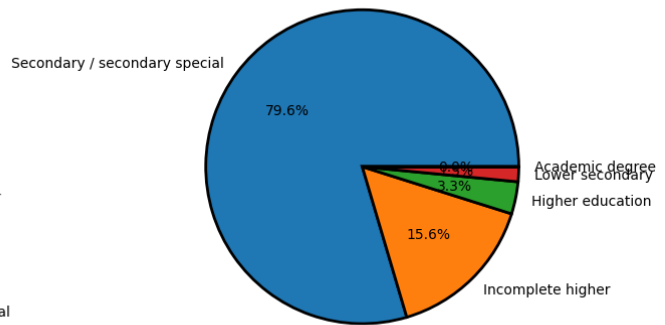
Target group 1 - Parameter:NAME\_INCOME\_TYPE



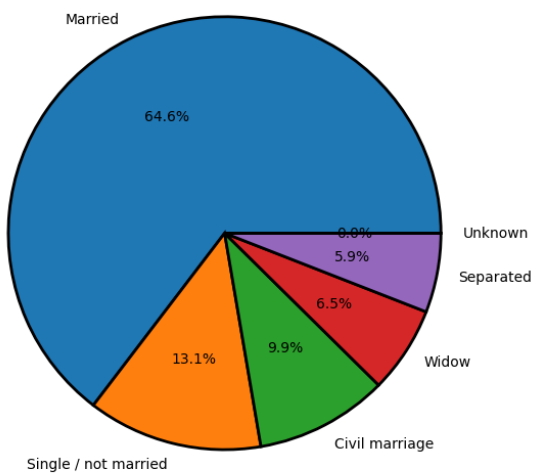
Target group 0 - Parameter:NAME\_EDUCATION\_TYPE



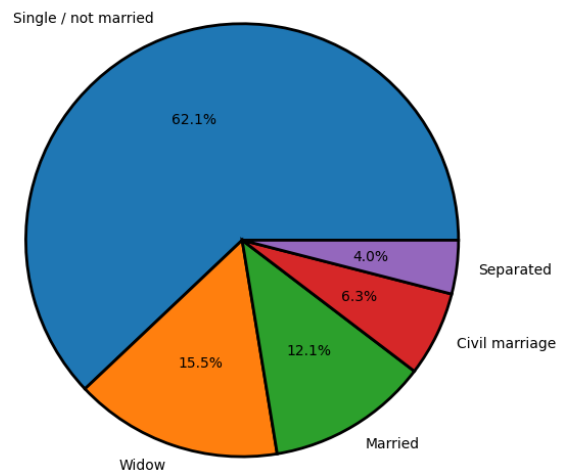
Target group 1 - Parameter:NAME\_EDUCATION\_TYPE



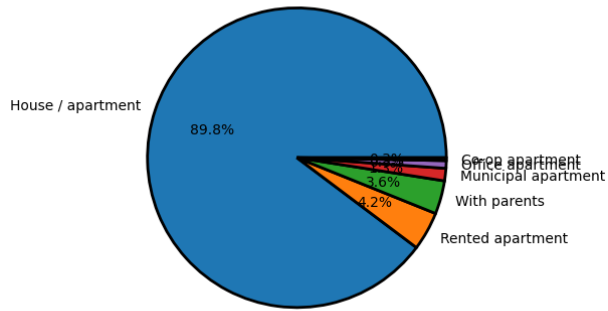
Target group 0 - Parameter:NAME\_FAMILY\_STATUS



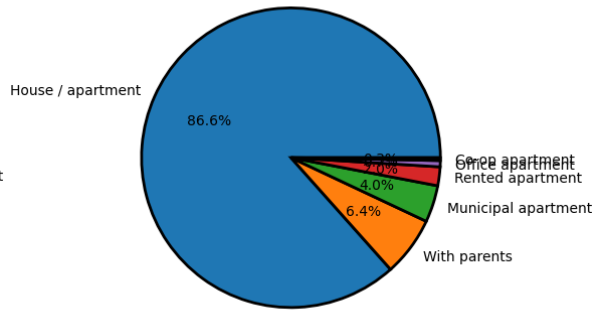
Target group 1 - Parameter:NAME\_FAMILY\_STATUS



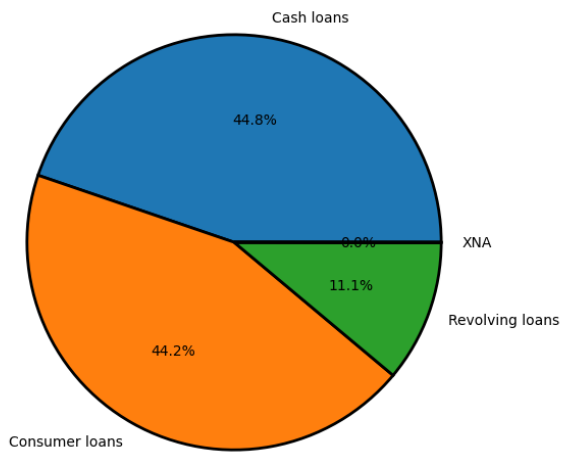
Target group 0 - Parameter:NAME\_HOUSING\_TYPE



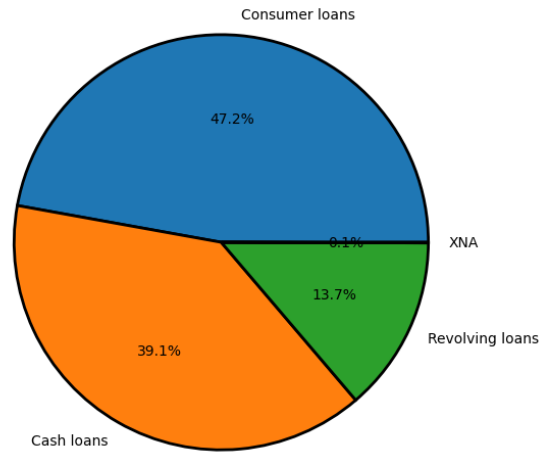
Target group 1 - Parameter:NAME\_HOUSING\_TYPE



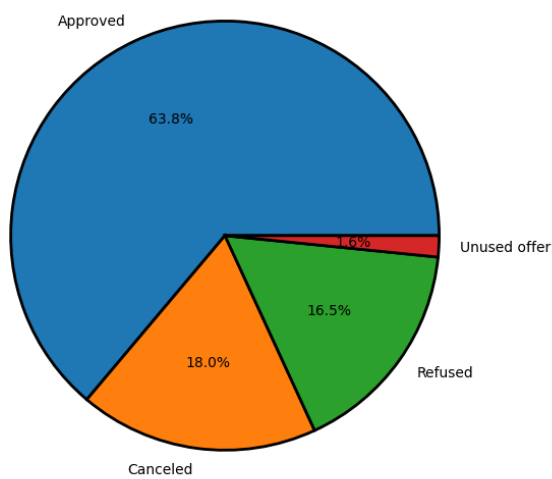
Target group 0 - Parameter:NAME\_CONTRACT\_TYPE\_y



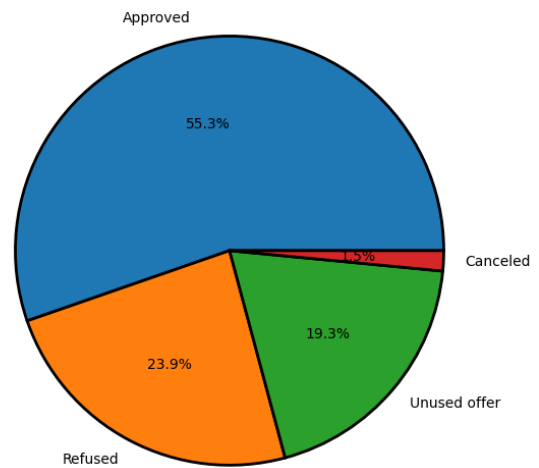
Target group 1 - Parameter:NAME\_CONTRACT\_TYPE\_y



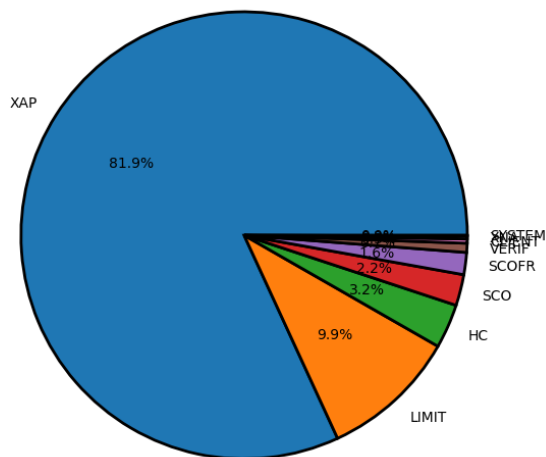
Target group 0 - Parameter:NAME\_CONTRACT\_STATUS



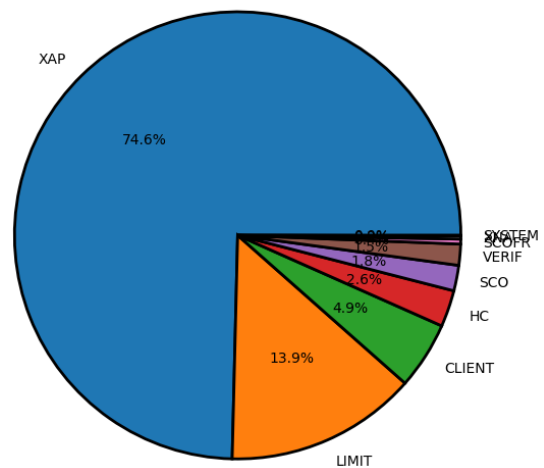
Target group 1 - Parameter:NAME\_CONTRACT\_STATUS



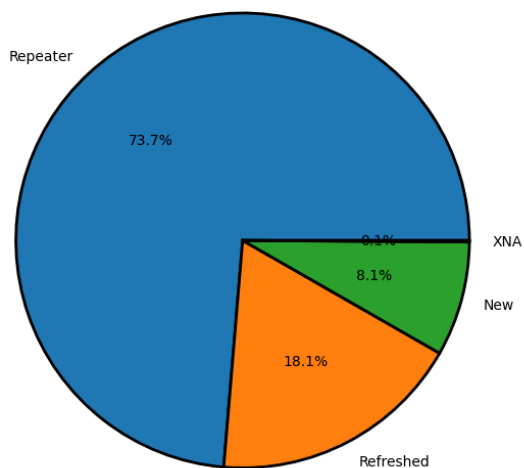
Target group 0 - Parameter:CODE\_REJECT\_REASON



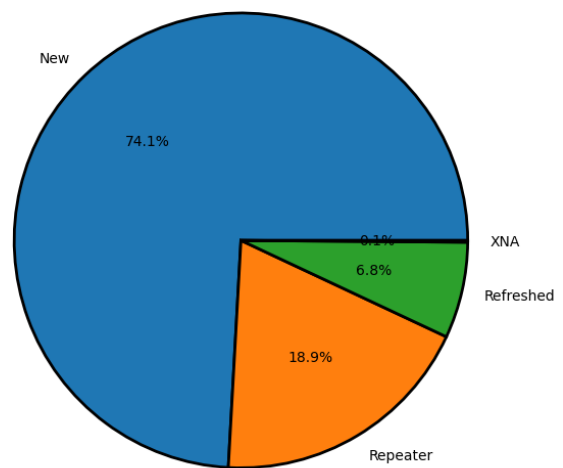
Target group 1 - Parameter:CODE\_REJECT\_REASON



Target group 0 - Parameter:NAME\_CLIENT\_TYPE



Target group 1 - Parameter:NAME\_CLIENT\_TYPE



## Inferences from pie chart

### 1. Loan type

Revolving loans have lesser payment problems compared to Cash loans : 4.8% / 7.9%

### 2. **\*\*Gender\*\***

Females are better at paying loans compared to males.

### 3. Car ownership

those who own a car have slightly better chances of being able to pay loans - nearly 3% better chances

### 4. Real estate ownership

Those who own a real estate : averages seem same for both target groups

### 5. **\*\*Name\_suite (accompanied by)\*\***

Those who were accompanied by Family are much better at paying loans compared to others. Hence, "NAME\_TYPE\_SUITE" is an important parameter to check while deciding if loan should be awarded.

### 6. **\*\*Income type : \*\***

State servernts are pretty good at payments followed by Working class. Working class people TOP in the category of those who experience problem in payments.

### 7. **\*\*Education : \*\***

Even though we saw, Secondary educated applicants were the most, But their biggest chunk falls in the category of those who have had problems in repayment.

Applicants with Higher education are the best at loan re-payments. hence, the most suitable candidates for awarding loan.

### 8. **\*\*Family Status: \*\***

Married people are best at loan repayment while Singles experience difficulties while paying back in most cases. Very important.

### 9. Housing type:

same in both cases.



#### 10. Weekday:

Same in both cases

#### 11. Contract status of previous application:

Those who were approved loan in previous application are good at paying back in current loan also.

#### 12. **Client Type**

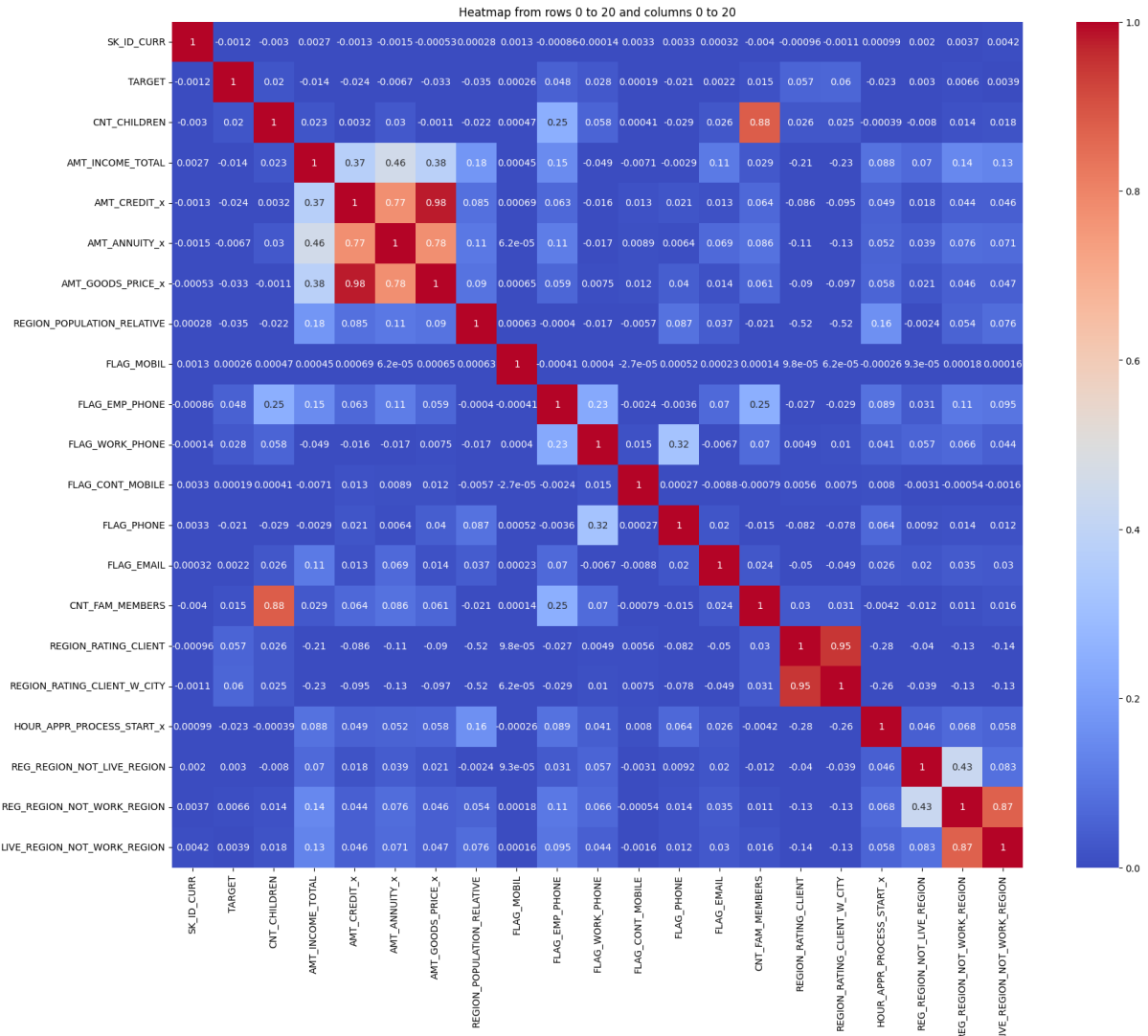
Repeaters are the best at paying loans while New applicants are the riskiest as they experience the maximum cases of payment difficulties.

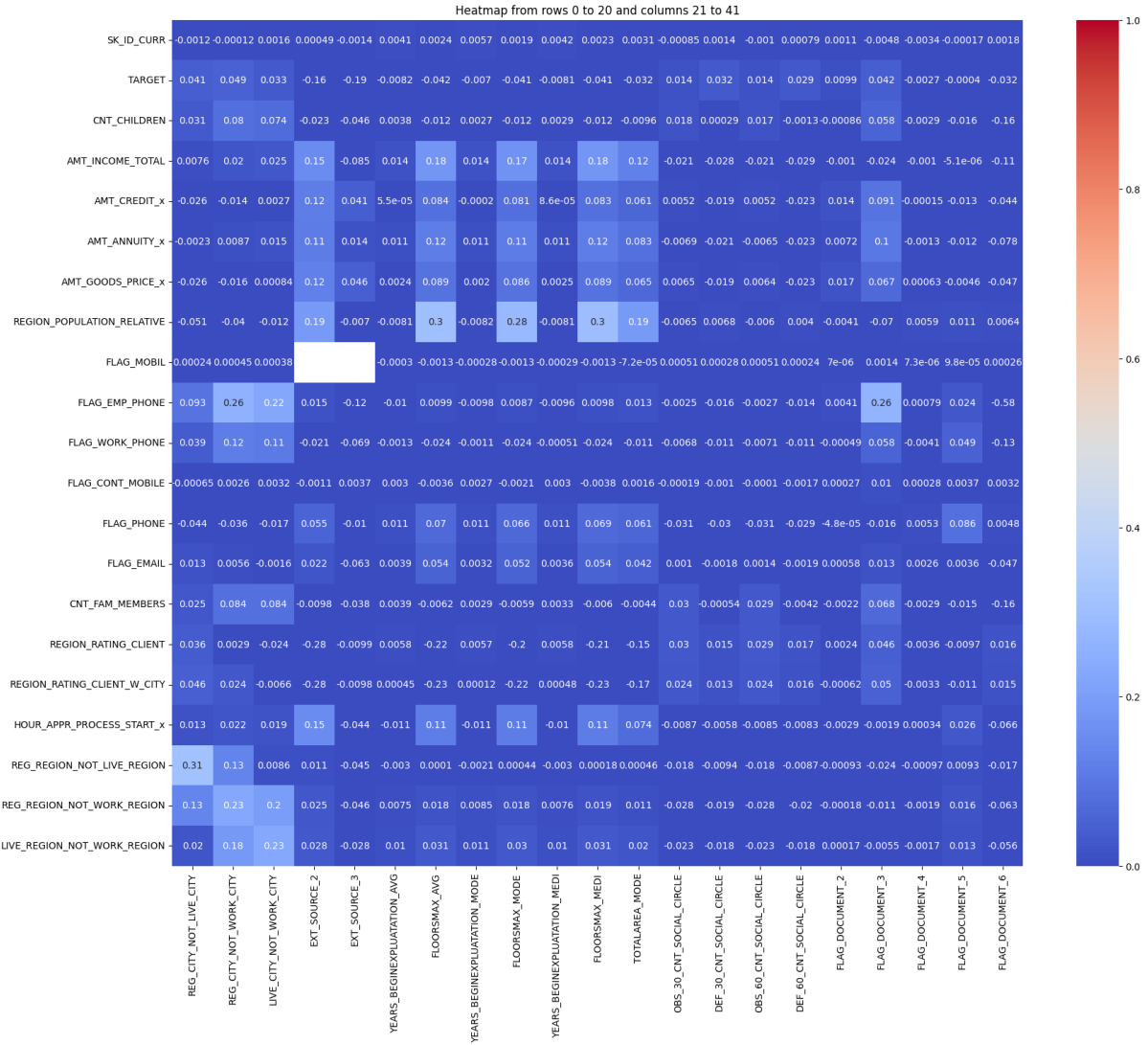
#### 13. Product type

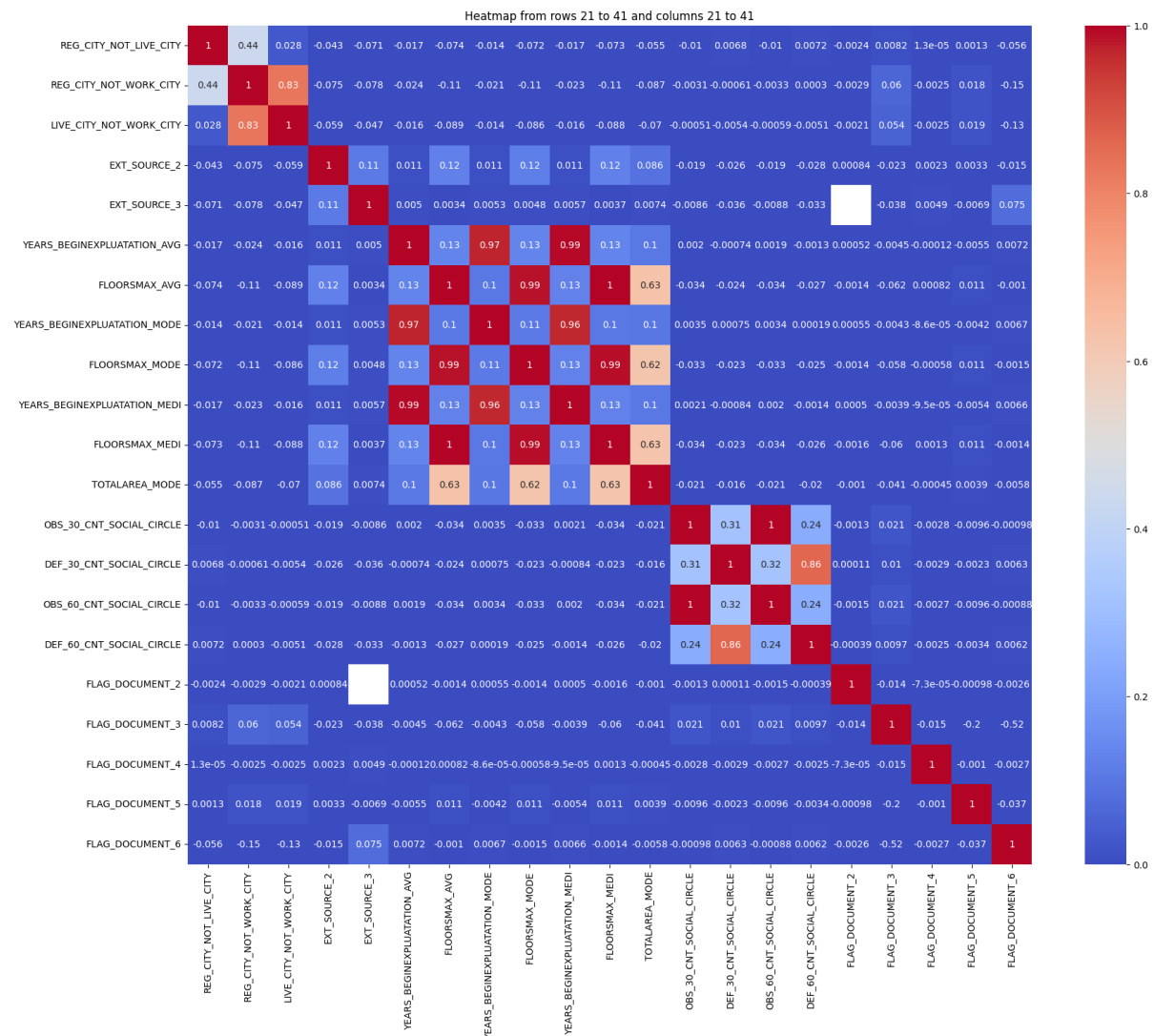
those who were using Product type x-sell in previous applications have good payment record. those with XNA are experiencing most payment problems

All other charts are not giving much - so ignoring them.

Heatmaps:







Inferences from heatmaps:

Strong correlations found between:

1. AMT\_GOODS\_PRICE\_x & AMT\_CREDIT\_X

2. AMT\_CREDIT\_x & AMT\_ANNUITY\_x

3. AMT\_ANNUITY\_x & AMT\_GOODS\_PRICE\_x

4. CNT\_CHILDREN & CNT\_FAM\_MEMBERS

5. DAYS\_LAST\_DUE vs DAYS\_TERMINATION

6. REGION\_RATING\_CLIENT\_W\_CITY & REGION\_RATING\_CLIENT

7. REG\_REGION\_NOT\_WORK\_REGION & LIVE\_REGION\_NOT\_WORK\_REGION

Moderately effective relations

1. CNT\_PAYMENT vs AMT\_APPLICATION

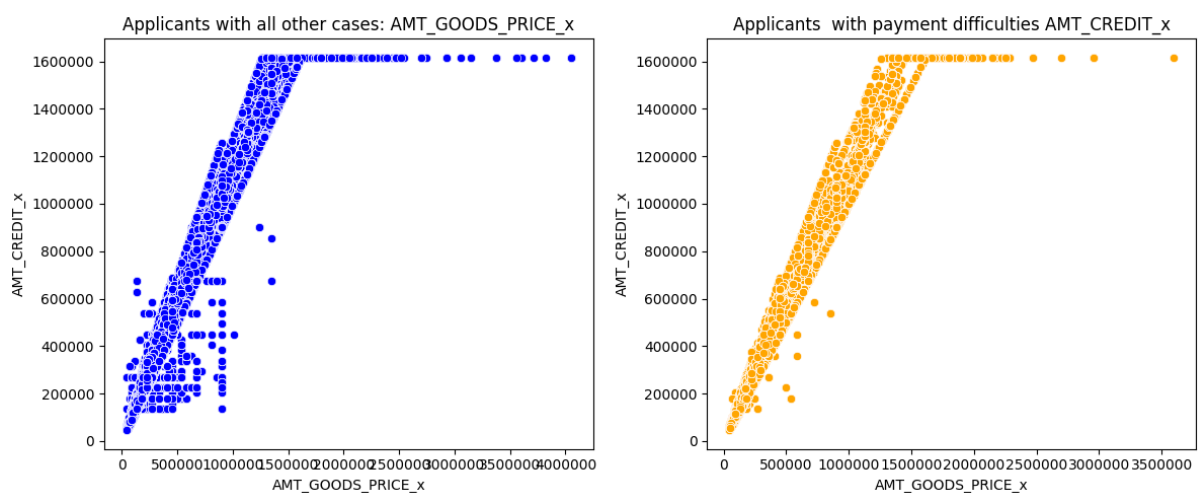
2. AMT\_CREDIT\_y vs CNT\_PAYMENT

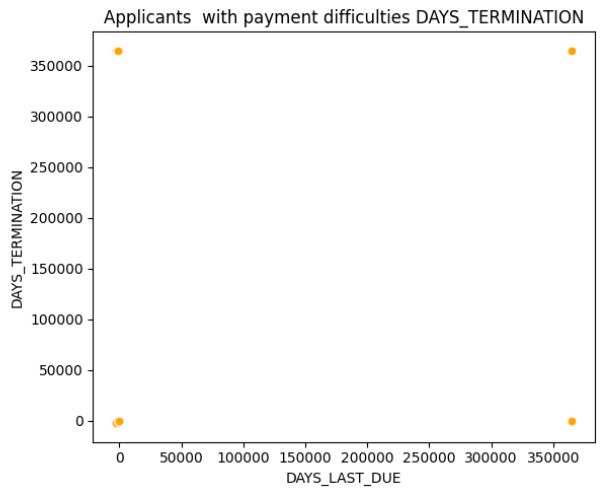
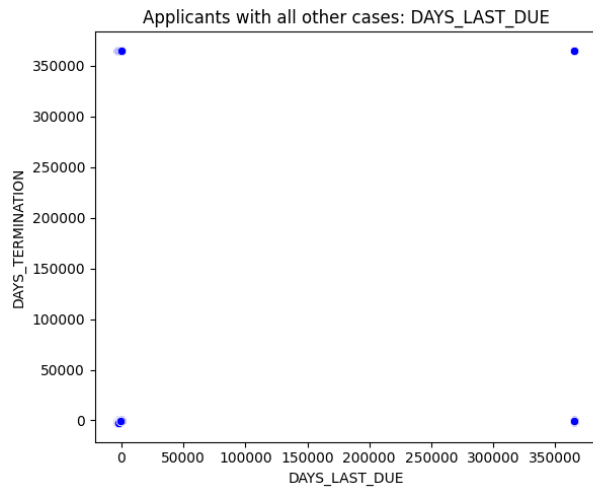
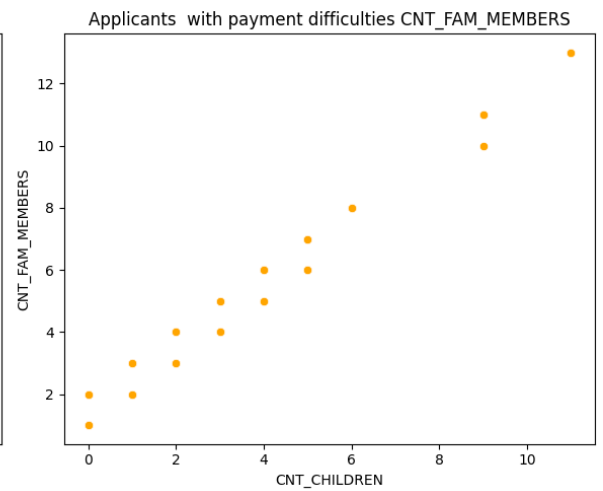
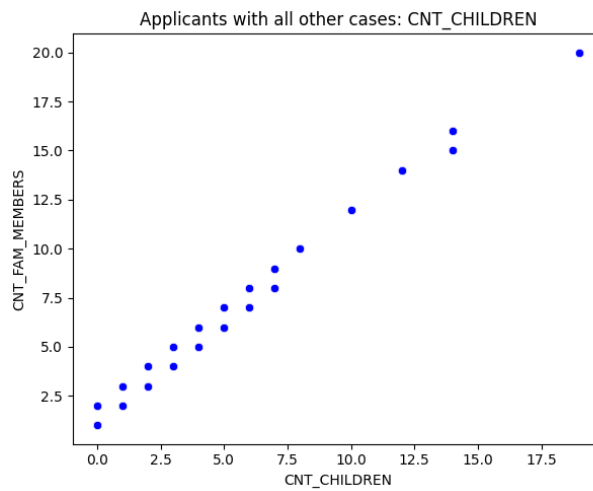
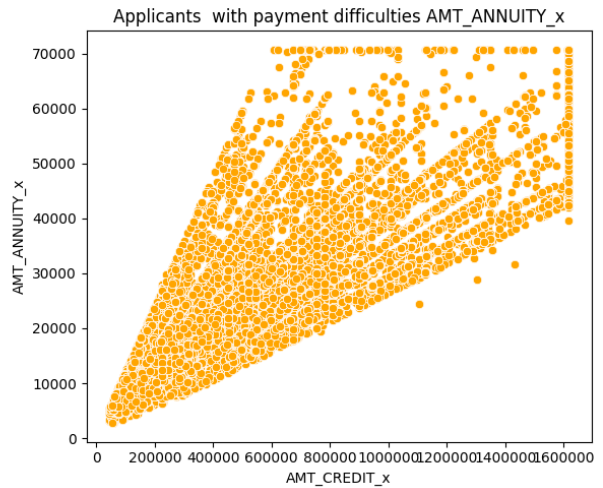
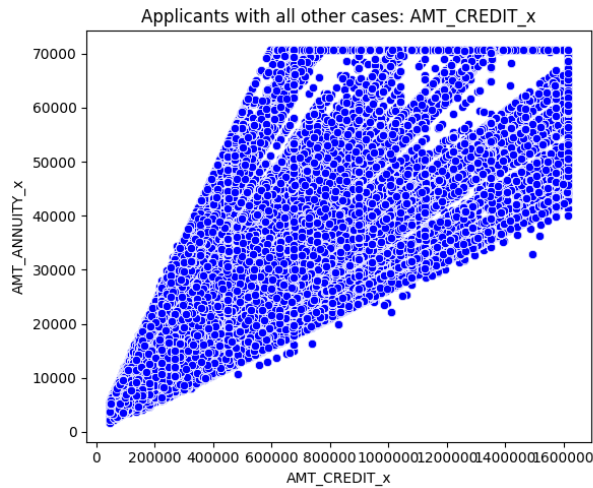
3. AMT\_GOODS\_PRICE\_y = CNT\_PAYMENT

**\*\*\*This is only for understanding - as their direct correlation with TARGET GROUP is not seen here\*\*\***

**\*\*Running scatterplot on all with hue of TARGET\*\***

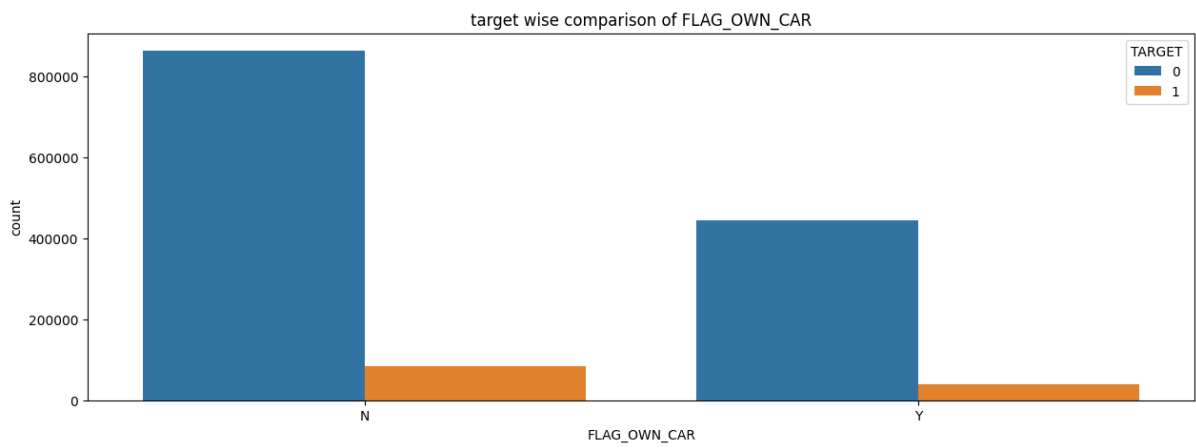
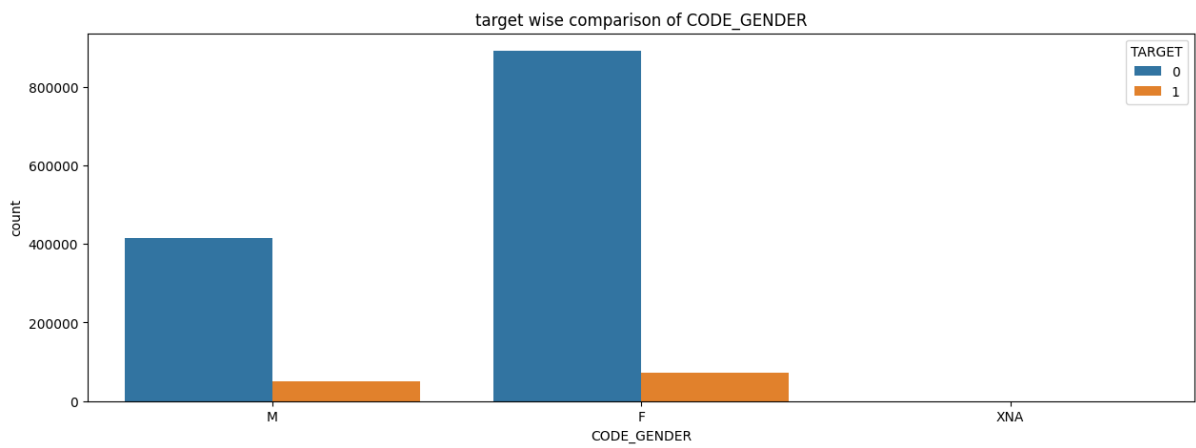
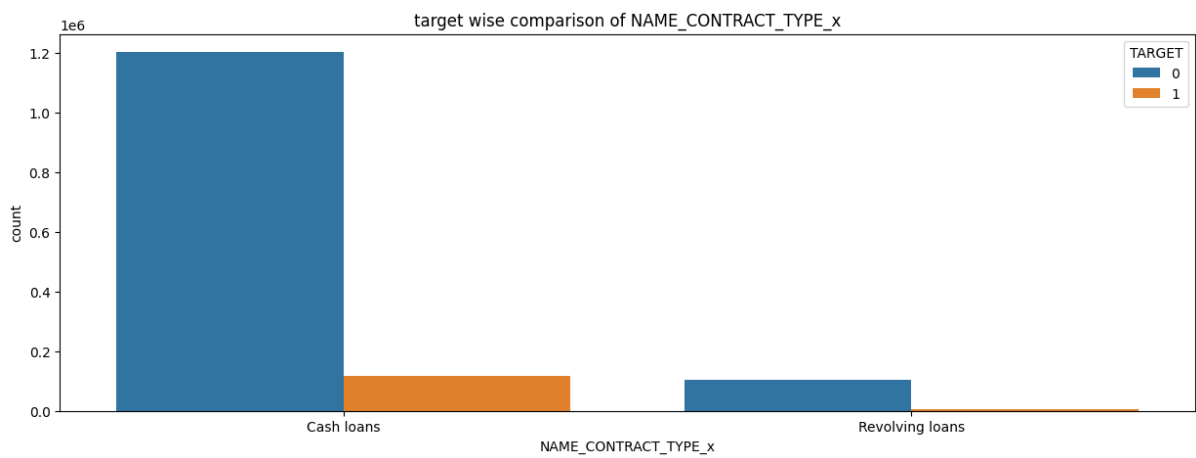
**Scatterplots :**

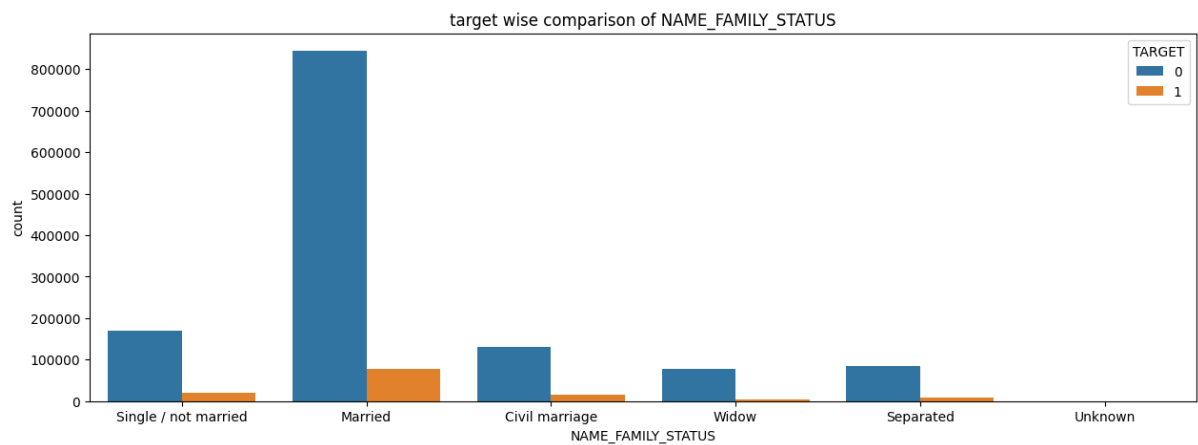
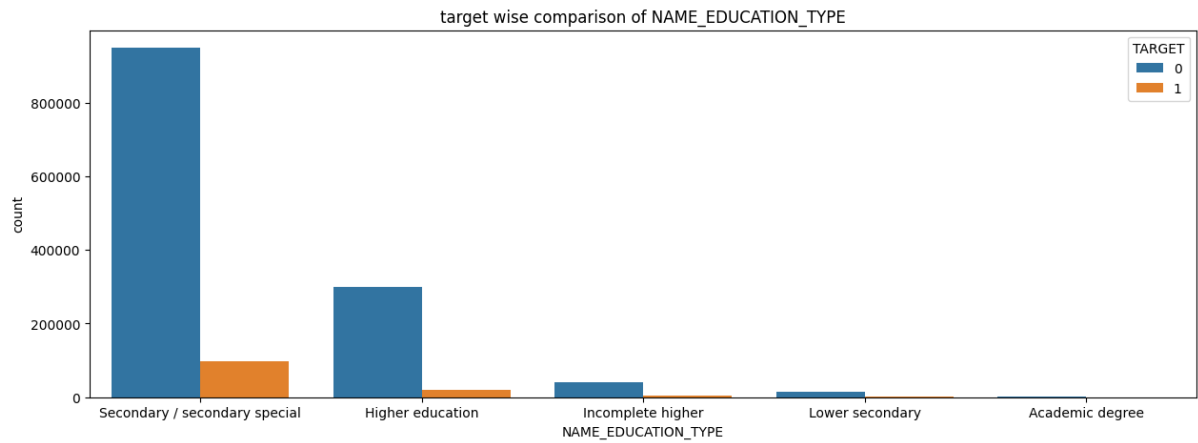
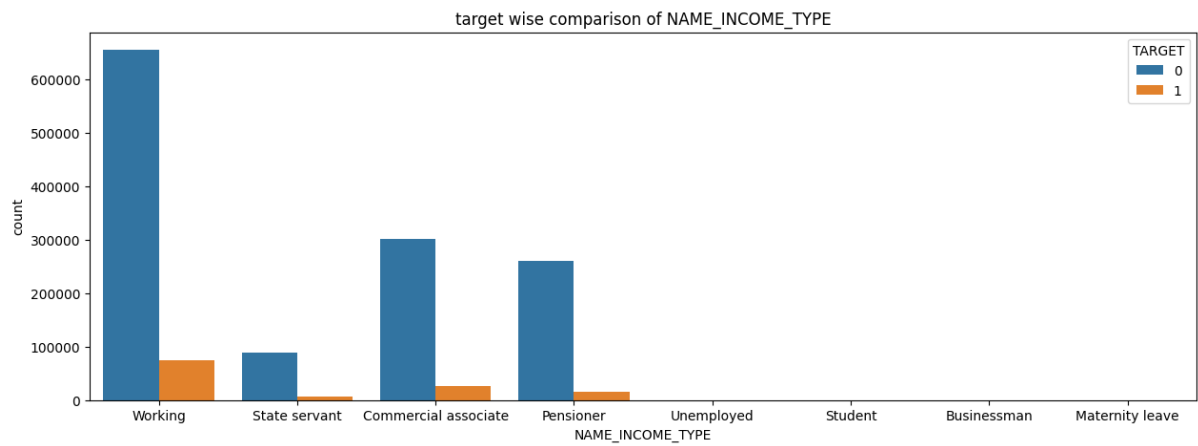
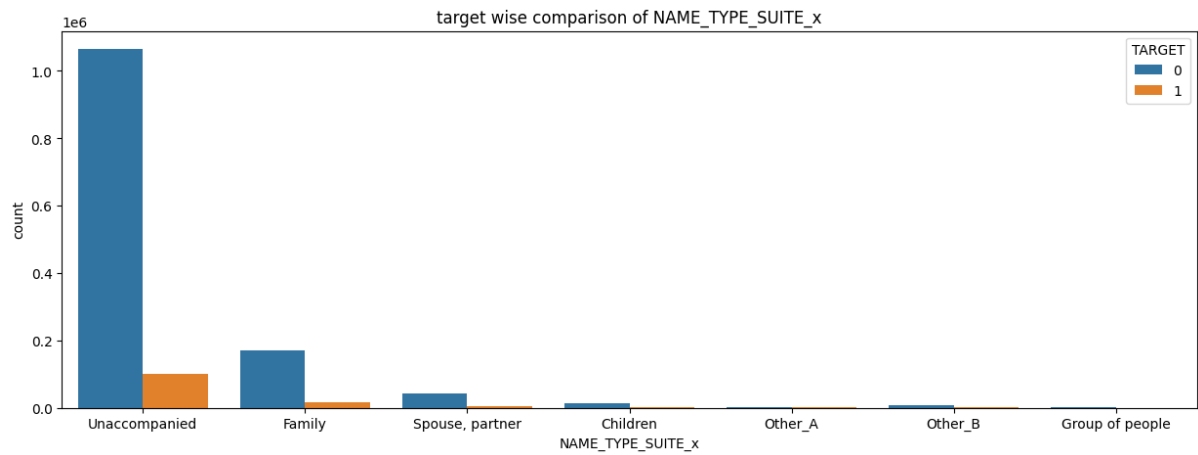




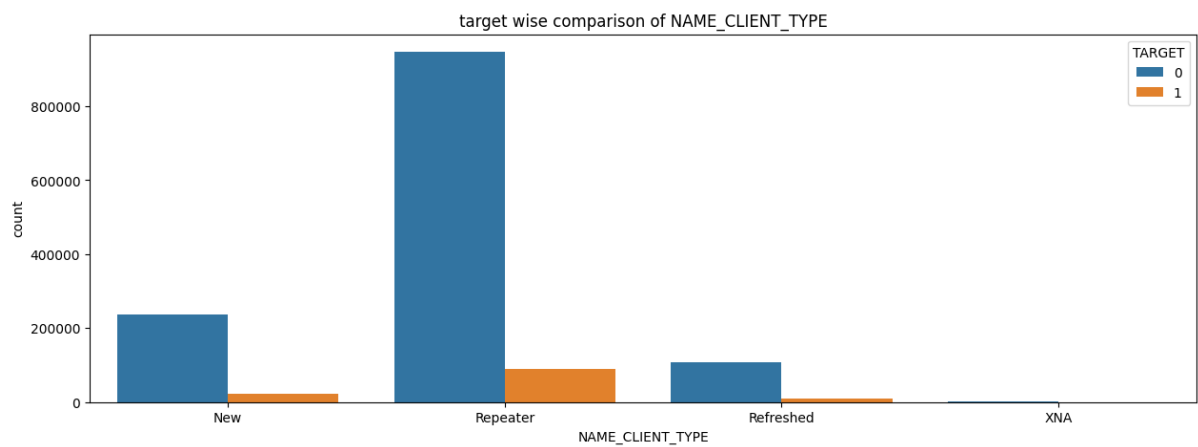
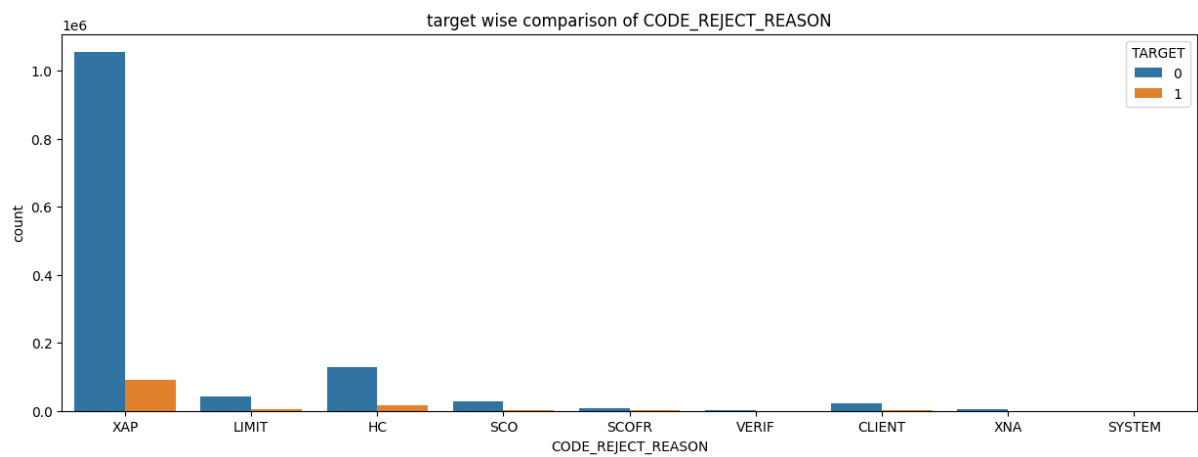
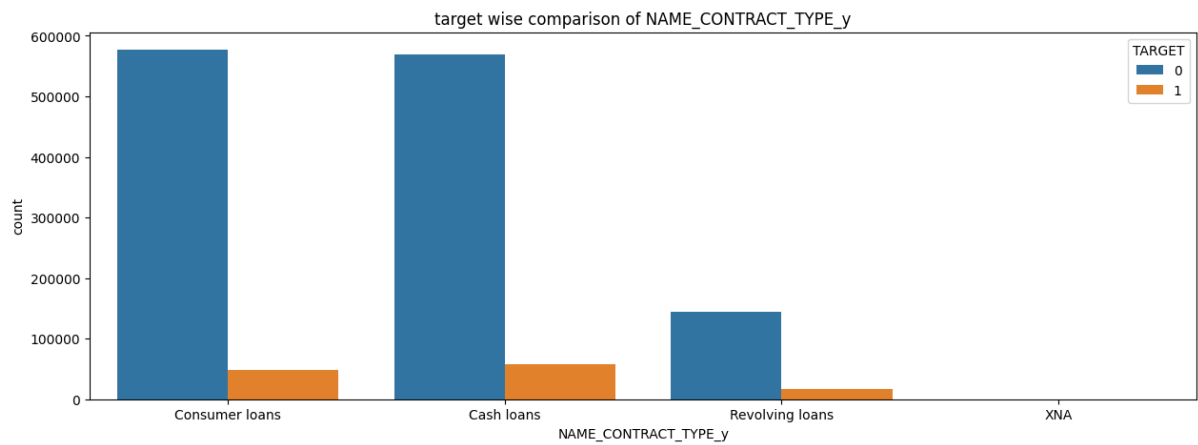
*Inference : Although there are parameters that show high correlation and some that do not - yet there is no such combination that shows a different behavior for those with Payment difficulties vs those with all other cases*

Count plots :









### **Count plot inferences:**

#### 1. Contract type wise :

very high imbalance on Cash loans.

ratio of 0 and 1 seems similar for all loan types. hence no distinguishing factor.

## 2. Gender :

Females apply far more for the loans compared to men. However, male feel find more difficulties paying.. to be proven with ratio.

## 3. Car Ownership:

very high number of people who apply loans do not have car.

no visible differentiation in behavior between those who own and those who dont own a car.

## 4. Realty ownership:

there are more applicants with realty who apply for loan than those who don't  
same like car ownership, no differentiation.

## 5. Suite type (Who accompanied customer during prev application) - (OLD and NEW)

People mostly visit unaccompanied, followed by Family.

visibly no differentiation.. ratio to be checked.

## 6. Income type:

Working people form biggest chunk of applicants.

payment difficulty ratio seems higher for commercial associates and Pensioners. ratio calculation to be done to objectify number

## 7. Education:

applicants are mostly either secondary educated and then followed by those with "higher education". all others constitute a small part.

No visible differentiating behavior

8. Family status:

Most of applicants are married, followed by Single.

those in civil marriage have higher visible ratio of applicants with payment difficulties compared to others. objectify with ratio calculation.

9. housing type :

Most of applicants have House/Apartment.

those who live with parents have higher ratio of applicants with payment difficulties

10. Application processing day (OLD and Current)

no visible difference. lesser business of credit on weekends.

11. EMERGENCYSTATE\_MODE:

this parameter is heavily biased on "NO"

no visible differentiation

12. Contract type (loan type):

Consumer and cash loans are the most in demand, though almost equal to each other. Revolving loans on third number.

Applicants are finding slightly more difficulty on cash loans compared to Consumer loans. This may be due to amount of loan. can be checked.

13. Contract status (from previous application):

those who were awarded loan last time are the highest with ease of payment in current applications as well.

However, differentiation ratio across types to be checked.. visibly not identifiable.

14. Payment type:

Most of people prefer to pay by "Cash through the bank".

XNA type applicants have higher numbers having issues in paying back.

15. Rejection\_reason:

for prev applications, XAP was the biggest reason for loan rejection. however, a very good percentage of those have not faced difficulty in paying back now.

HC (High credit ratio ? ) applicants in previous applications, have higher ratio of paying back again

16. Client type:

Repeaters are the biggest chunk of the current applicants as well.

differentiation ratio seems same across all types.

17. Portfolio :

POS is the most common Portfolio type, followed by Cash.

ratio to be calculated . "CASH" category seems to have higher ratio of payment difficulties

18. Channel\_type:

Credit and cash offices bring the maximum business for company.

County wide seems to have higher ratio. to be checked

19. yield group (Loan rate grouping):

XNA is the most common type of RoI group though not clear about its meaning.

ratio seems similar across categories.

**The analysis indicates a clear need for calculating ratio of all categories to confirm. This ratio or percentage is more easily visible in pie plots. Hence, together with inference of pie plots, final outcome can be concluded.**

## **Step 10 : CONCLUSION**

Following are the most important Parameters that affect the outcome as mentioned in the above cell

**Gender , Income type, Education , Family Status , Client type**

out of these parameters, the most effective correlations seen are :

**1. Gender**

**2. Education**

**3. Family status**

### **Project Conclusion:**

***For use case :*** important insights are derived from the project for the client if any.

***for Personal benefit :*** This project has given me enough confidence to handle big data and deal with multiple various parameters without having deep interactive understanding from the client.