

CREDIT EDA CASE STUDY

Group Name:

1. Shantam
2. Swati Hota

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.

In other words, the company wants to understand the driving factors (or driver variables) behind loan default, i.e. the variables which are strong indicators of default. The company can utilize this knowledge for its portfolio and risk assessment.

To develop your understanding of the domain, you are advised to independently research a little about risk analytics - understanding the types of variables and their significance should be enough).

The given data contains the information about the loan application at the time of applying for the loan. It contains two types of scenarios:

- **The client with payment difficulties:** he/she had late payment more than X days on at least one of the first Y instalments of the loan in our sample,
- **All other cases:** All other cases when the payment is paid on time.

When a client applies for a loan, there are four types of decisions that could be taken by the client/company):

Approved: The Company has approved loan Application

Cancelled: The client cancelled the application sometime during approval. Either the client changed her/his mind about the loan or in some cases due to a higher risk of the client he received worse pricing which he did not want.

Refused: The company had rejected the loan (because the client does not meet their requirements etc.).

Unused offer: Loan has been cancelled by the client but on different stages of the process.

In this case study, you will use EDA to understand how consumer attributes and loan attributes influence the tendency of default.

Three datasets given:

1. '*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**.

2. '*previous_application.csv*' contains information about the client's previous loan data. It contains the data whether the previous application had been **Approved, Cancelled, Refused or Unused offer**.

3. '*columns_description.csv*' is data dictionary which describes the meaning of the variables.

ANALYSIS

- Missing values ratio acknowledgement
- Data Imbalance
- Outliers
- Removing missing values
- Univariate analysis of the new and old applications
- Bivariate analysis of the new and old applications
- Correlation

Data Quality checks and handling missing values:

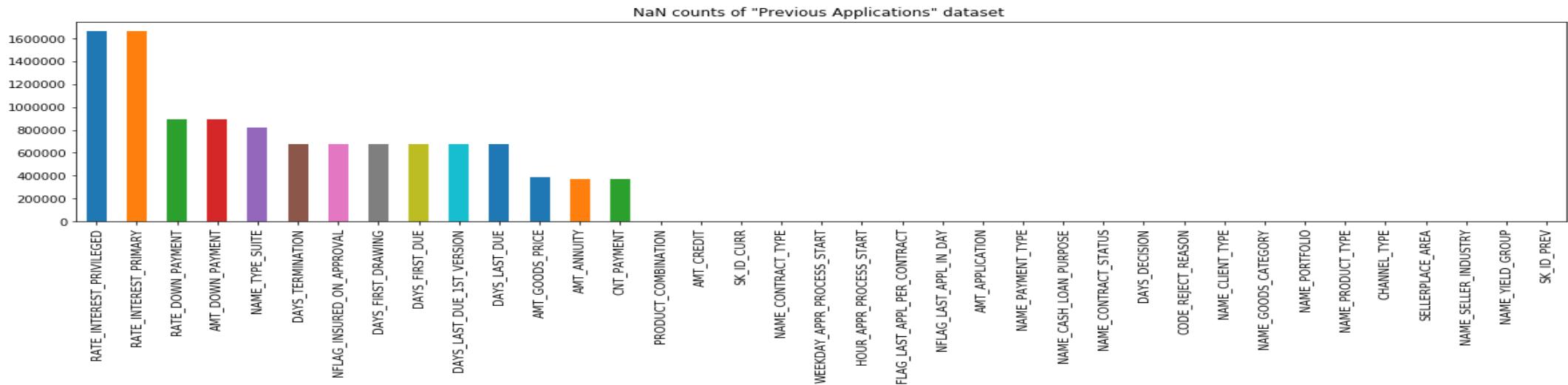
Two data sets were provided as part of this case study

- Application Data
- Previous application data

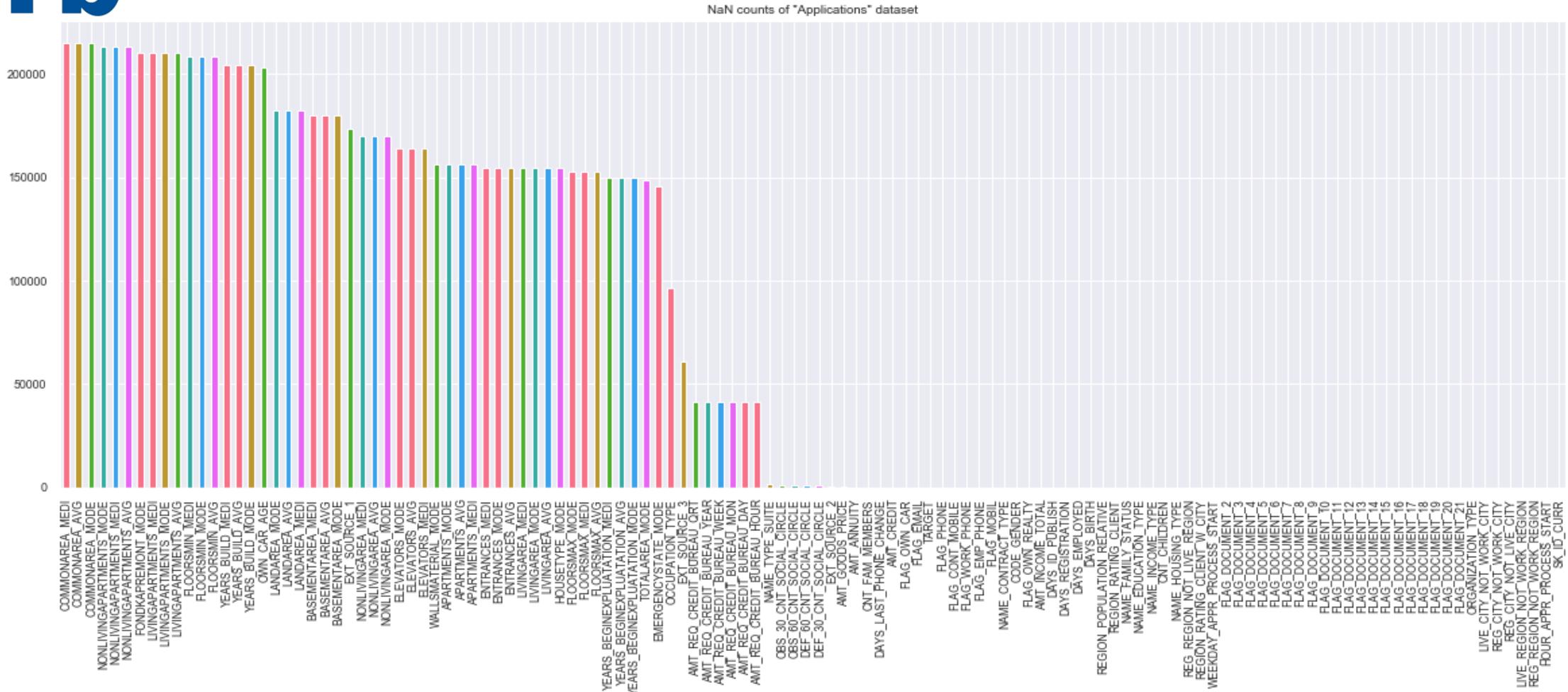
To find out the missing values and handle it :

We have observed there were many missing values in this data set, so using the indexing and count of missing values in each column, identified the % of missing values for each column.

- And dropped all columns from data frame for which missing values % is more than 50. We have also found few more columns which are having around 47% missing values. Since these are almost around 50% ,we have removed these columns as well.
- To extend the data enhancements and maintain a data frame with accurate values, we have identified columns with NULL values and even applied the same 47% logic and removed those columns from the data frame.



Insight: Near about 37% columns are having fully or partially missing values in “PREVIOUS -APPLICATION” dataset.



Insight: Near about 46% columns are having fully or partially missing values in “APPLICATION” dataset.

Data imbalance

For "TARGET" Class, we are having a huge data imbalance for Class-0 and Class-1

Class-0 : 92% data

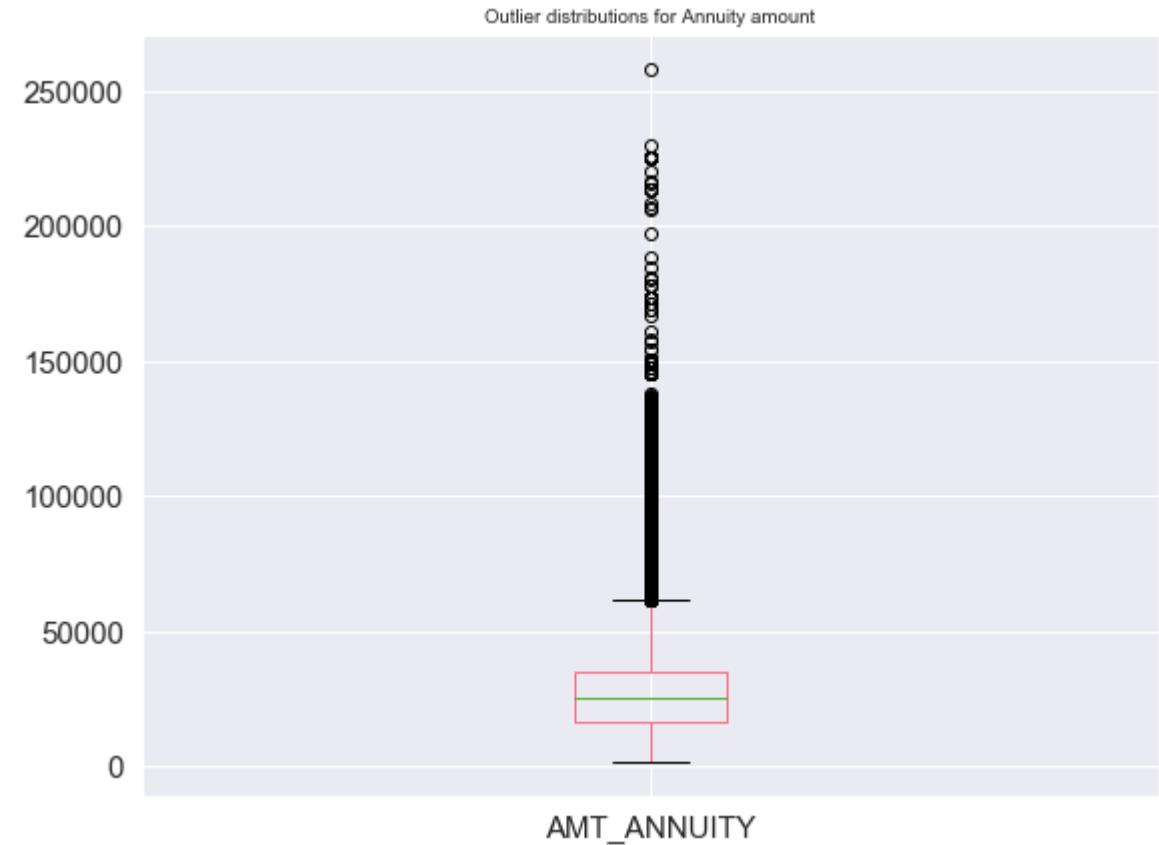
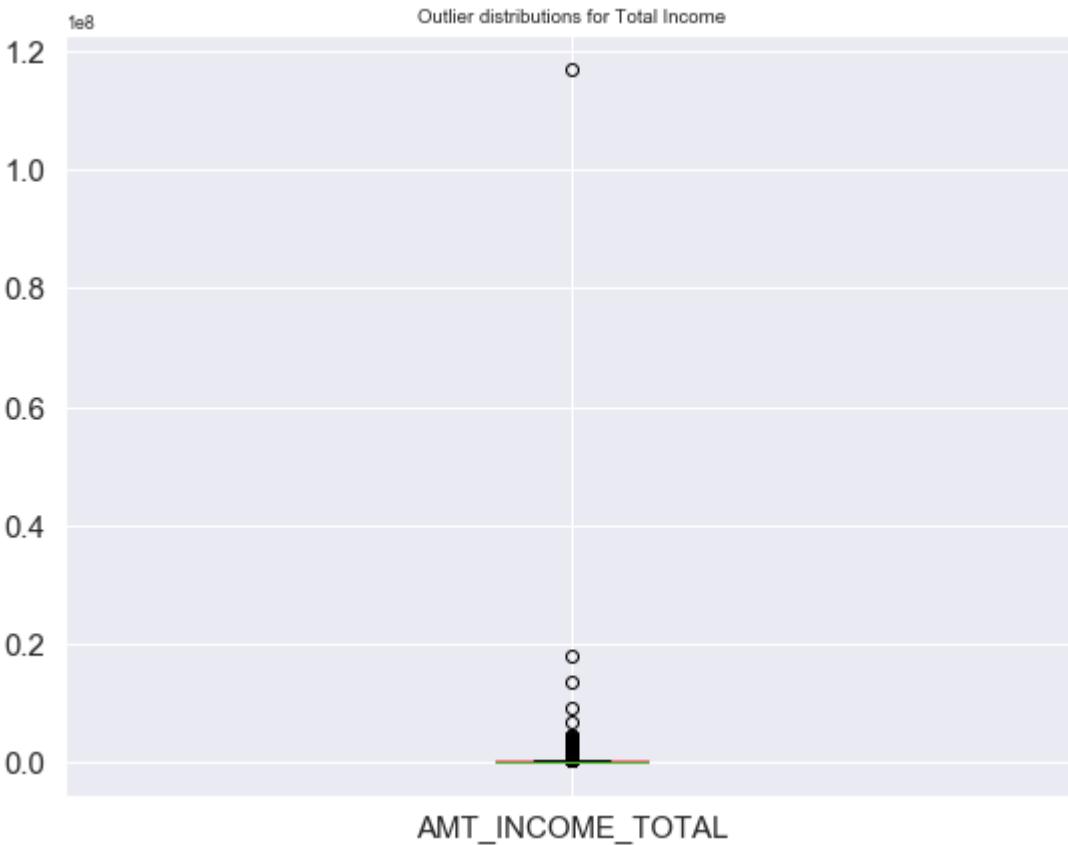
Class-1 : 8% data

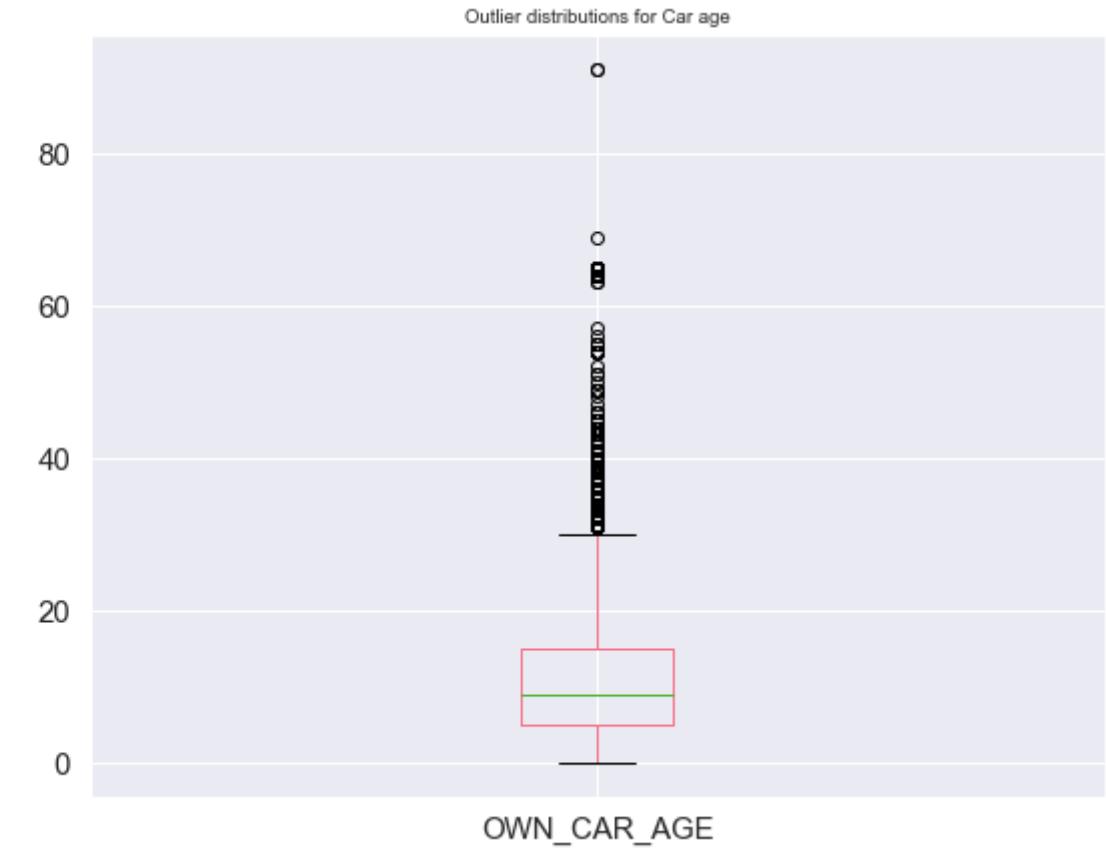
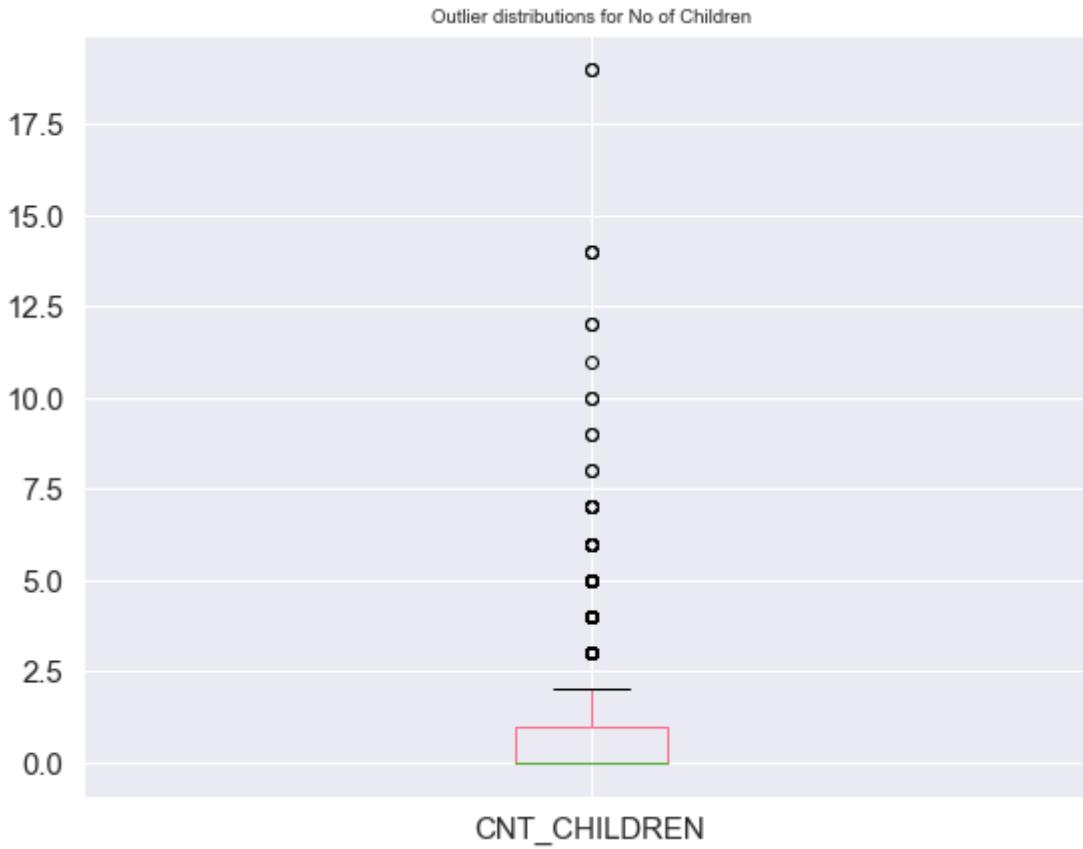
Checking the Data frame for outlier values and analyze them

- Checking the outliers for columns and understanding the reason to mention that as an outlier.
- Here in our analysis to find out the outliers, we have considered few numerical columns and analyzed the statistics of them.

Followings are few variables that have the outlier value(s)

- AMT_INCOME_TOTAL
- AMT_ANNUITY
- CNT_CHILDREN
- OWN_CAR AGE

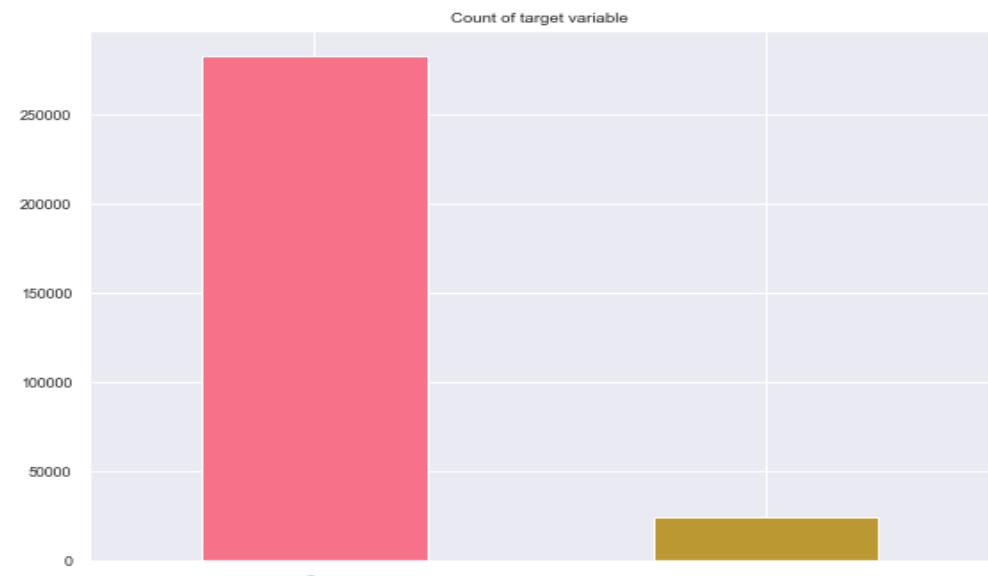
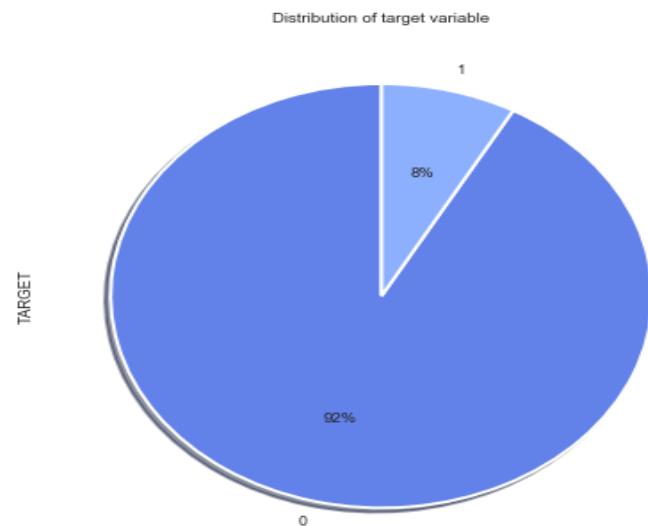




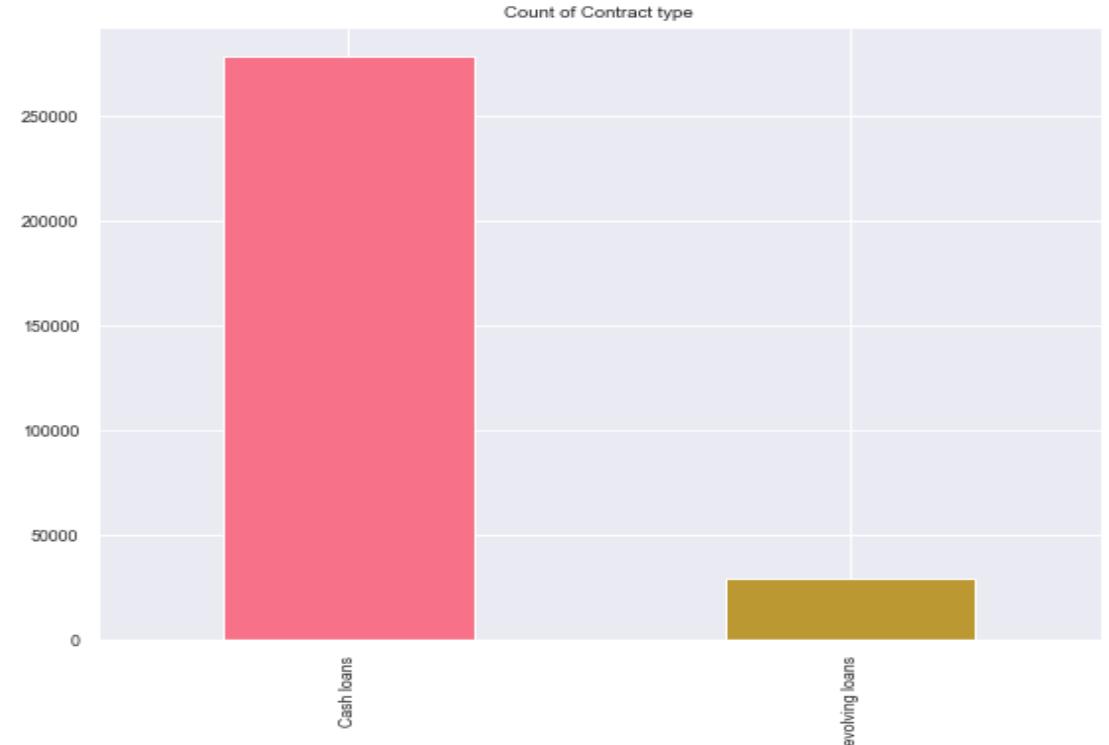
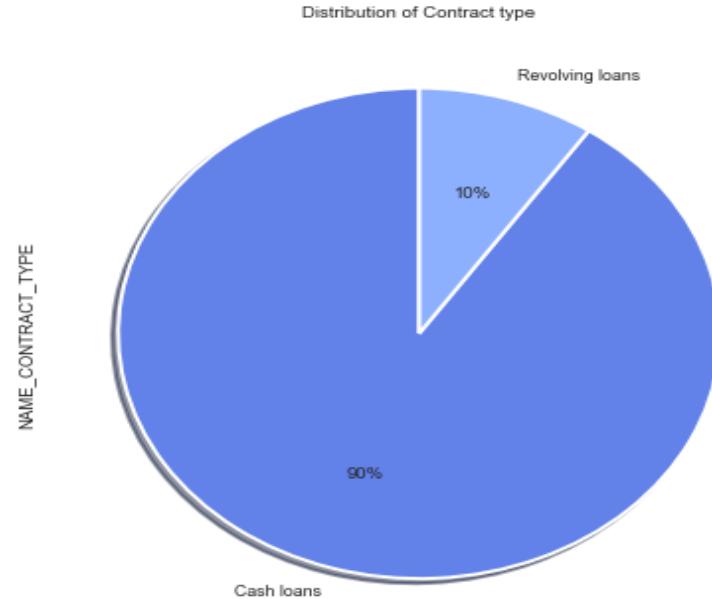
UNIVARIATE ANALYSIS:

Univariate analysis is perhaps the simplest form of statistical analysis. Like other forms of statistics, it can be inferential or descriptive. The key fact is that only one variable is involved. Univariate analysis can yield misleading results in cases in which multivariate analysis is more appropriate.

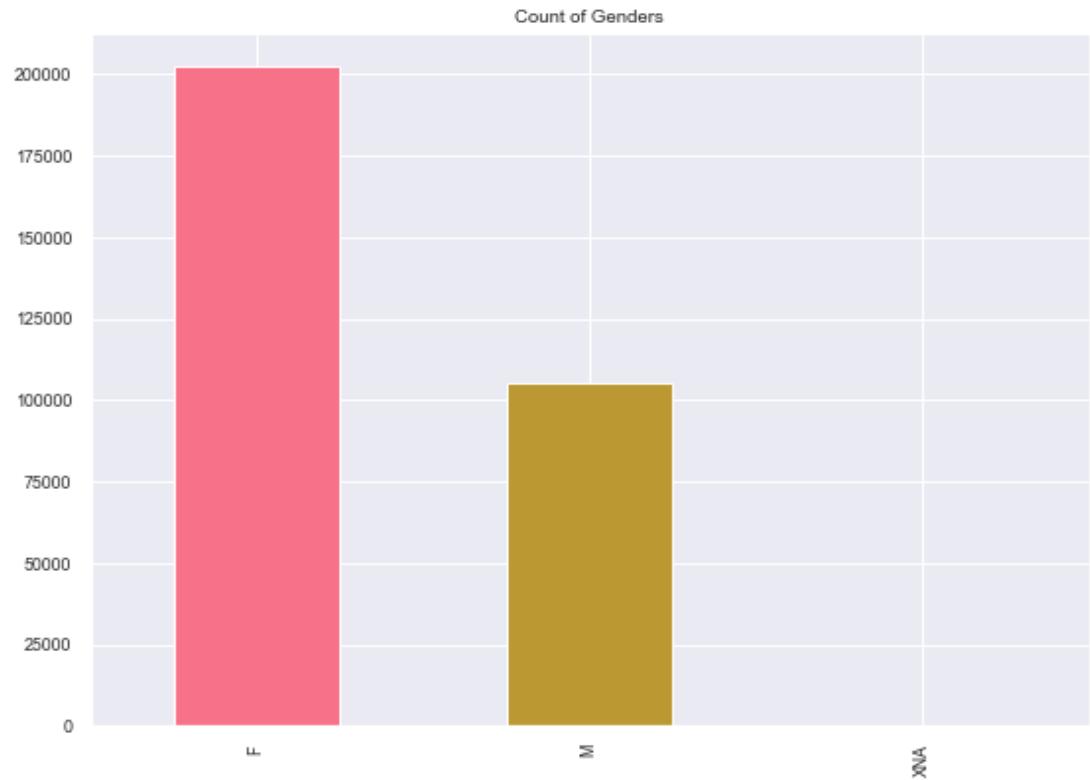
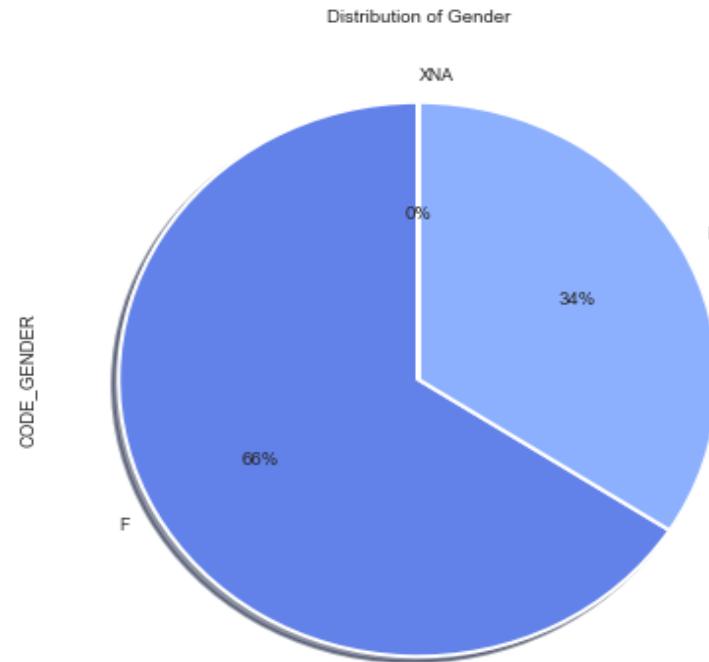
- While working with the data frame based on Target variable, **As shown in the below screenshot, we can clearly see the imbalance between Target type 1 and 0.**
- **Ratio is of 92 : 8**



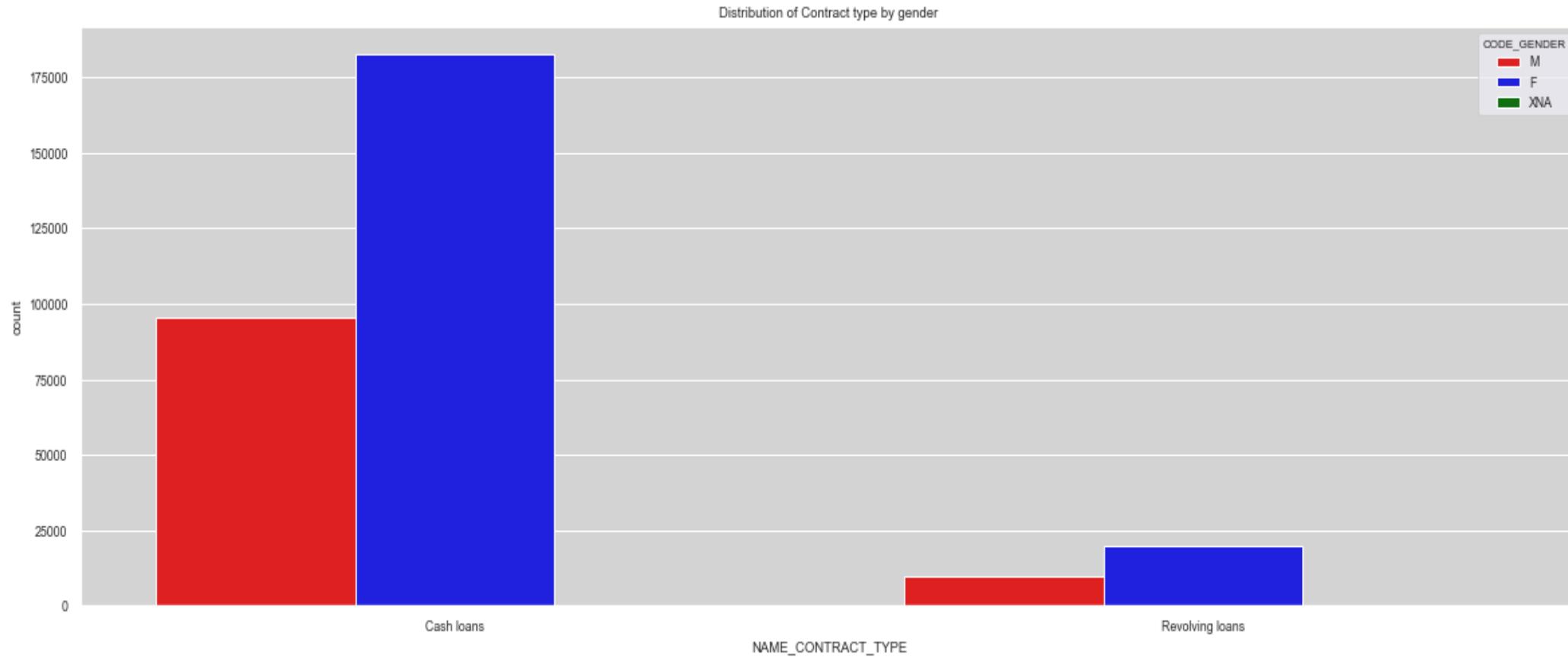
Insight: Out of total new loan applications about 8% are defaulters.



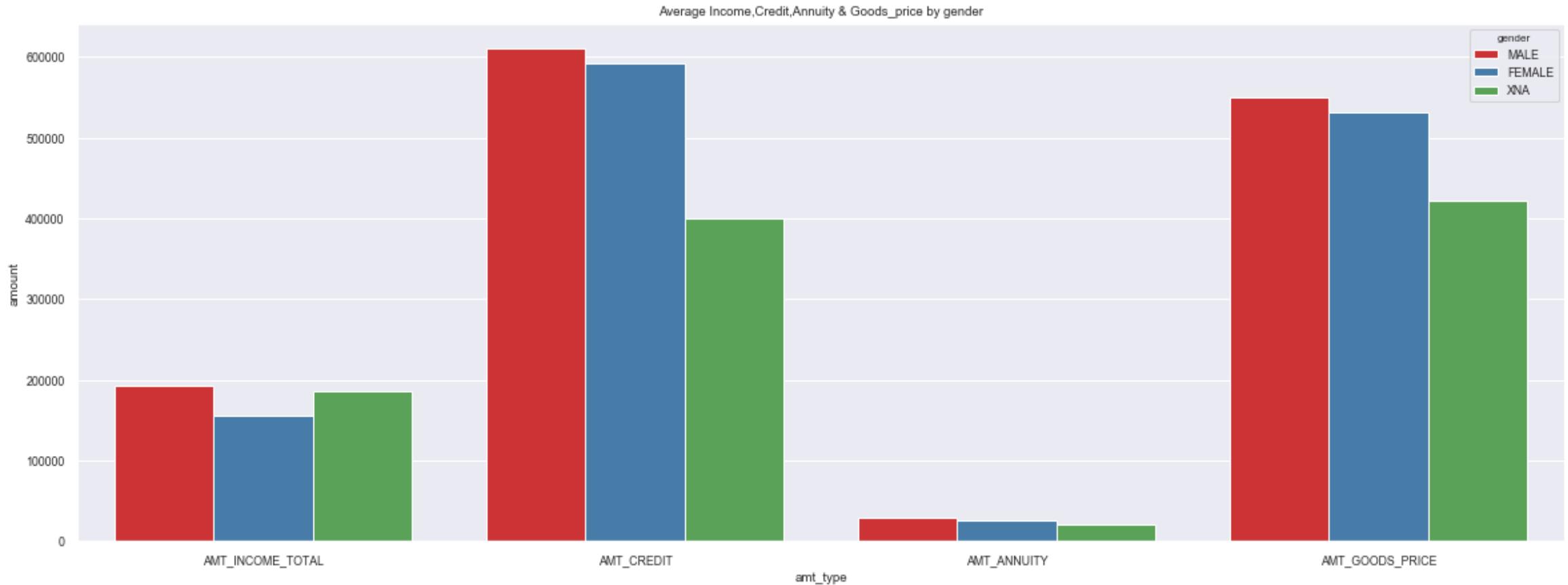
Insight: 90% loans are Cash loan, so most people prefer cash loans.



Insight: Out of total new loan applications 66% are female clients, Female clients are taking more loan than Male clients.



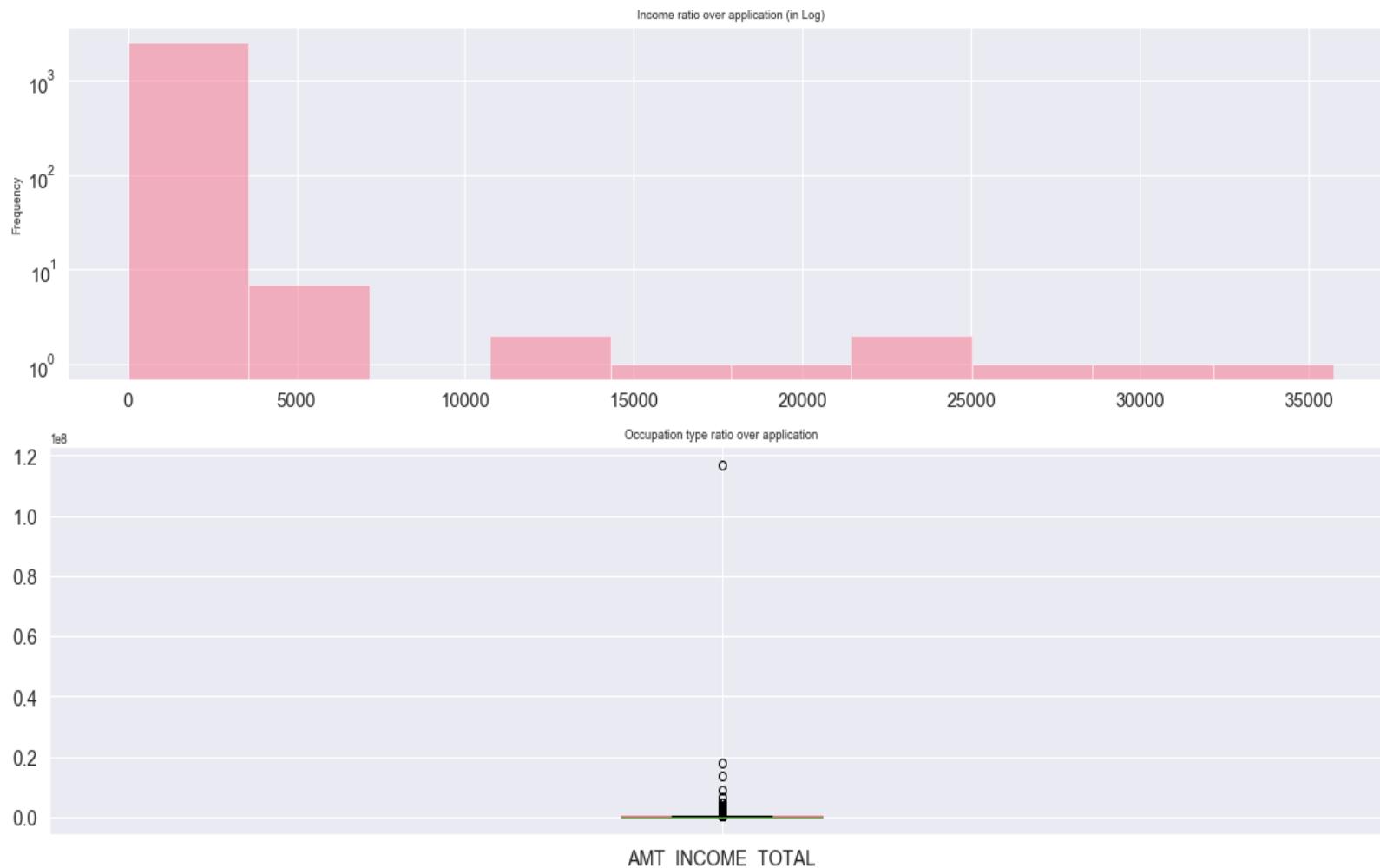
Insight: Here in both the scenario Female clients are taking more loans than Male clients whether its cash loan or revolving loan.



Insight: Here in above case, we observe that Male clients tops in almost every category with respect to the average amount.

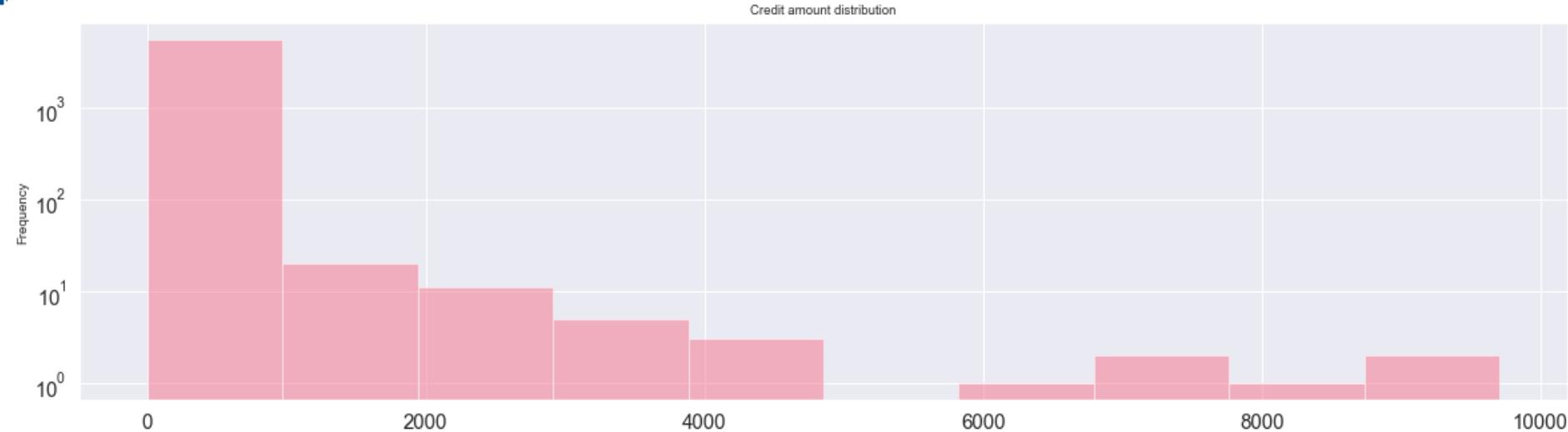
Income Distribution

AMT_INCOME_TOTAL Income Distribution

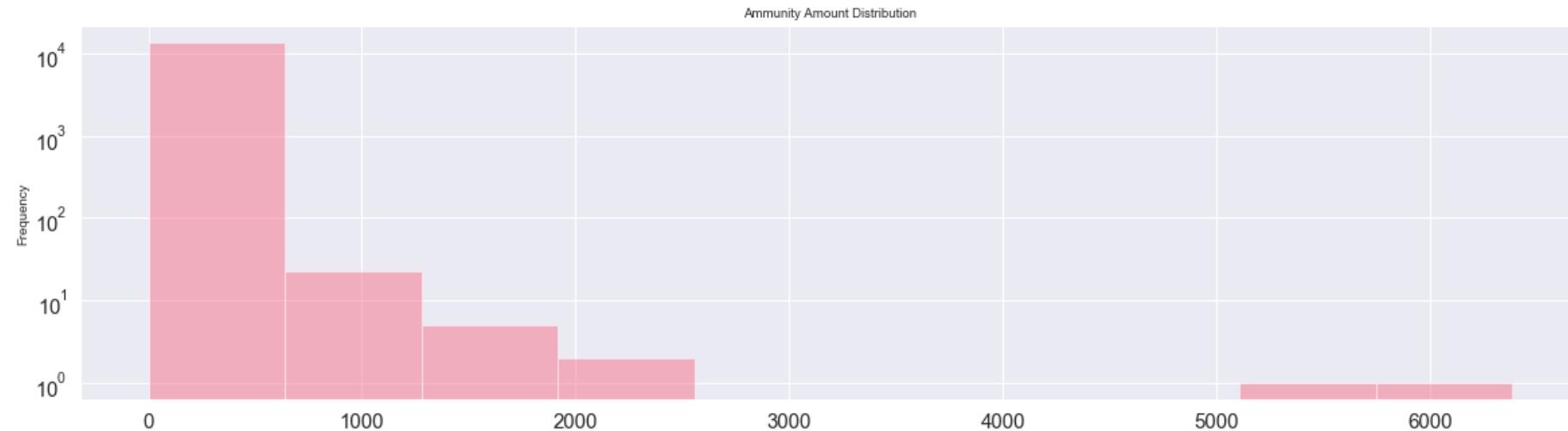


Insight: Mostly data is normally distributed with an exception outlier value in the "AMT_INCOME_TOTAL" column

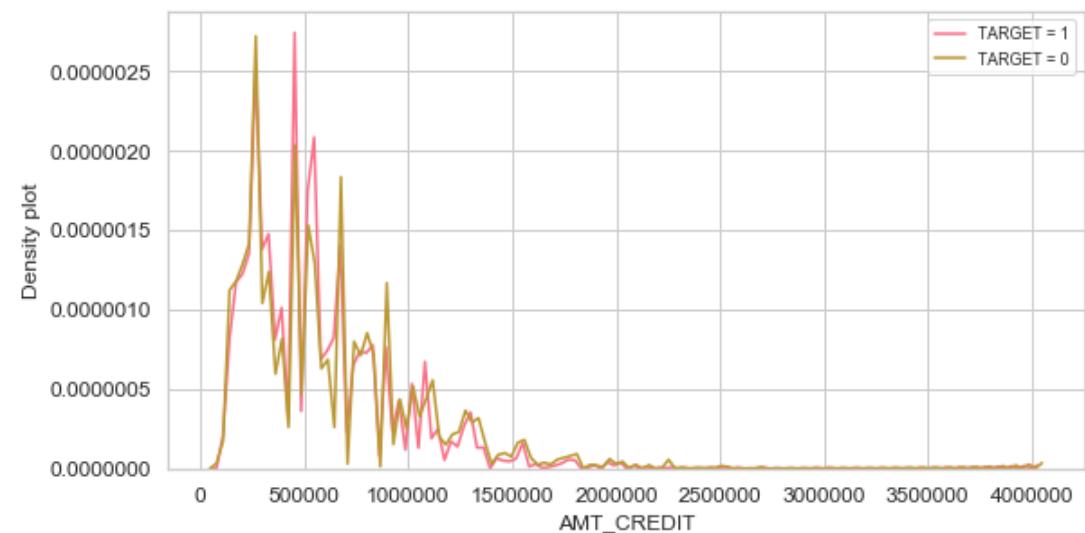
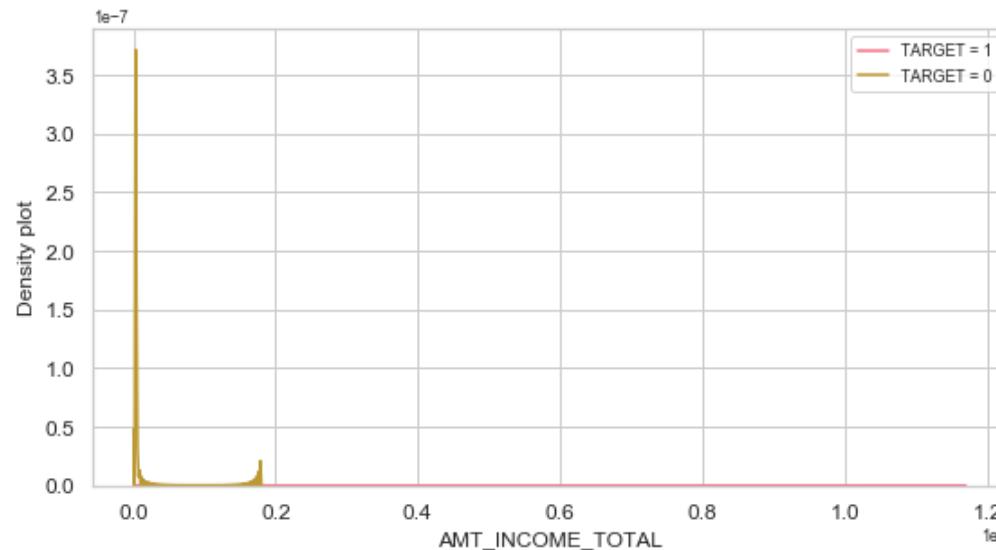
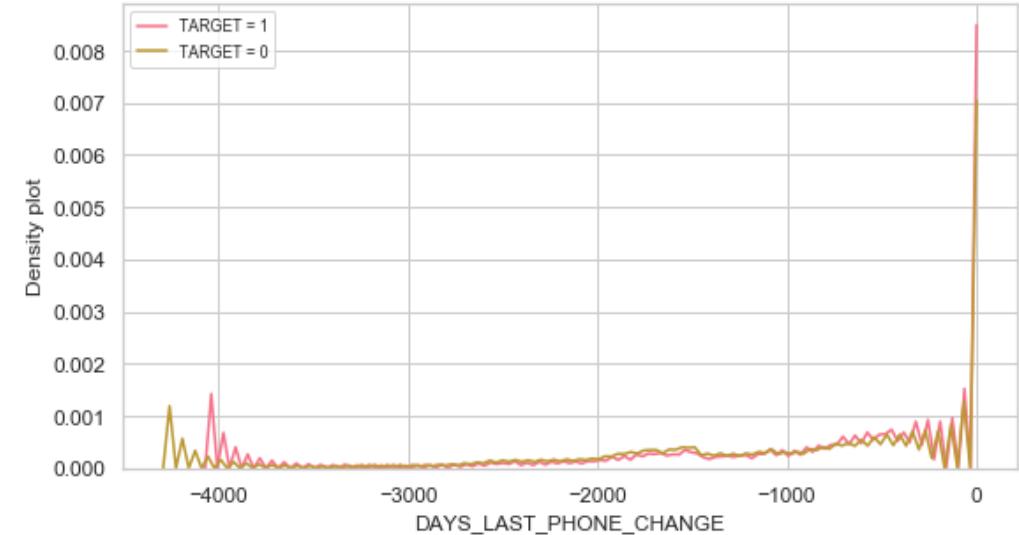
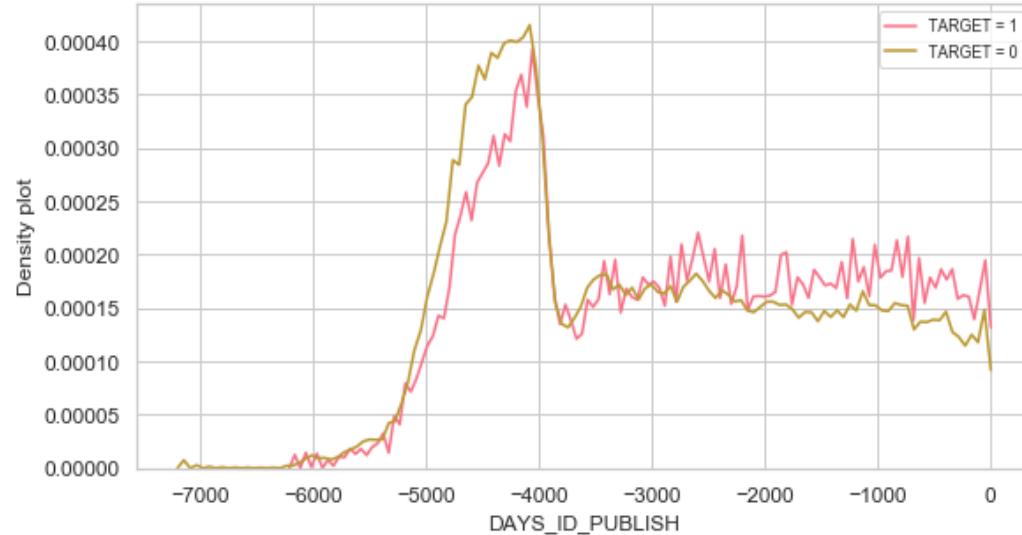
Amount Credit Distribution (In Log values)



Annuity Distribution (In Log values)



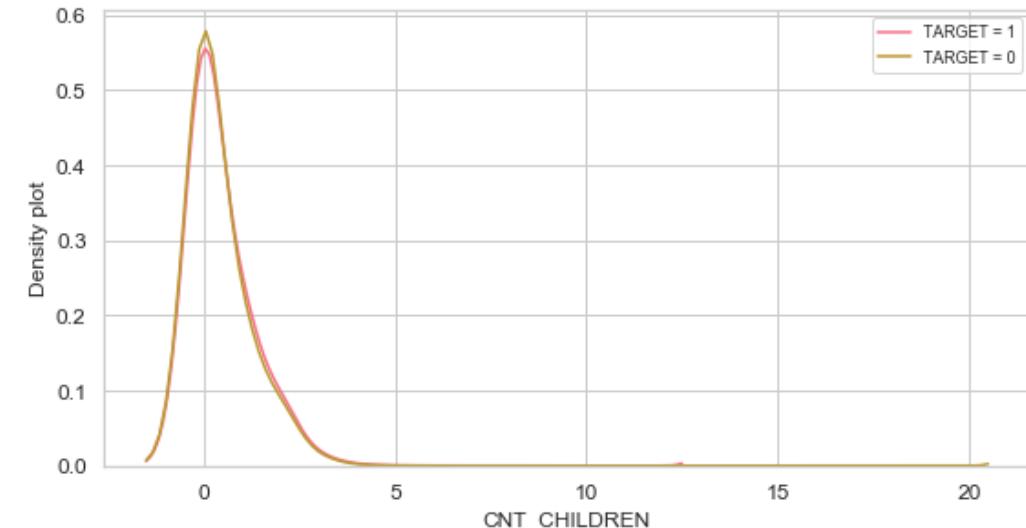
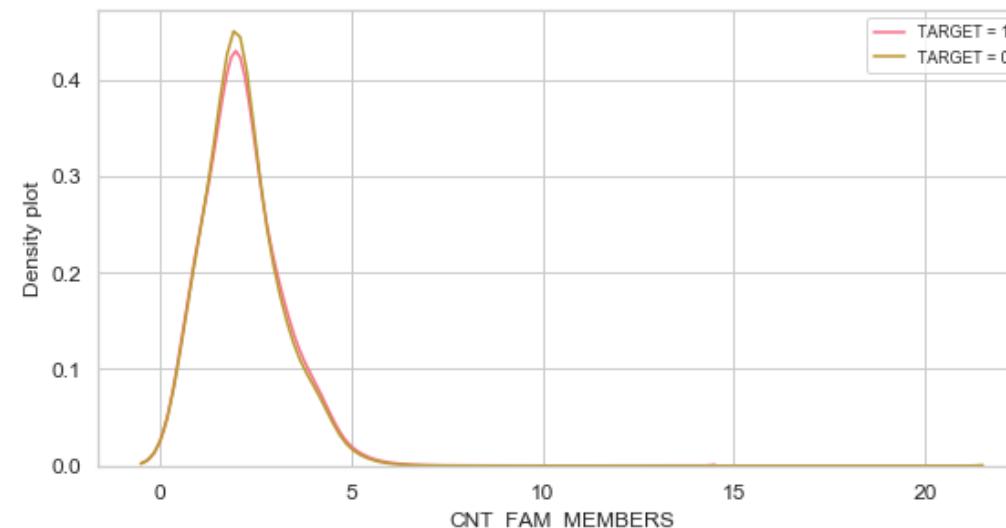
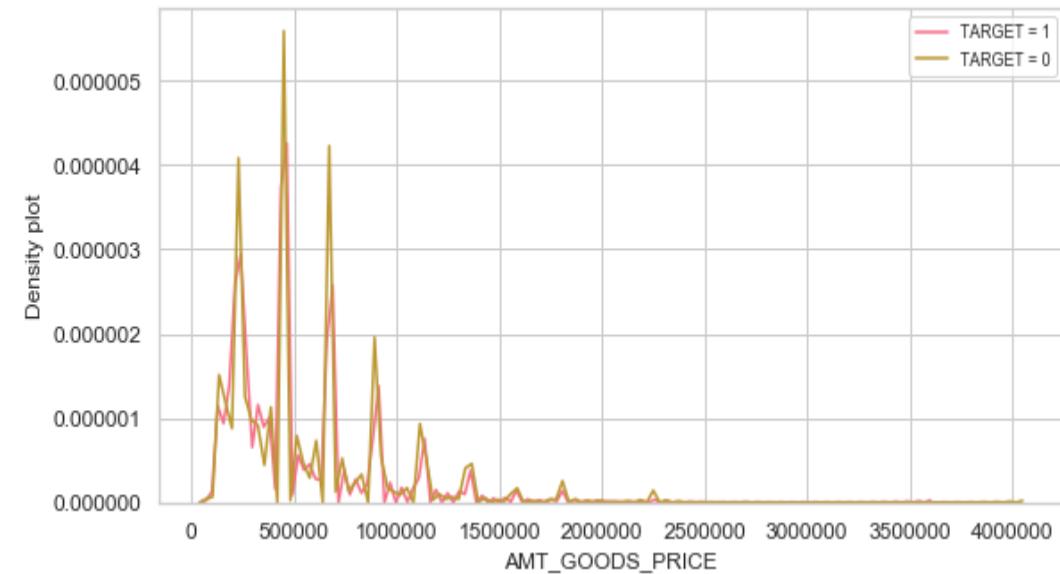
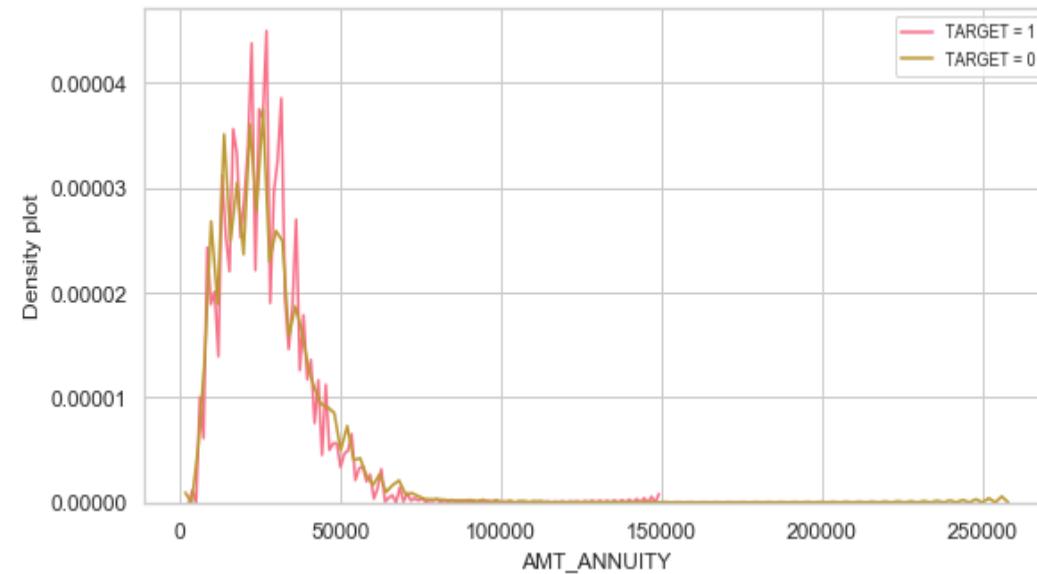
Comparative distribution of various variables for Target=0 and Target=1

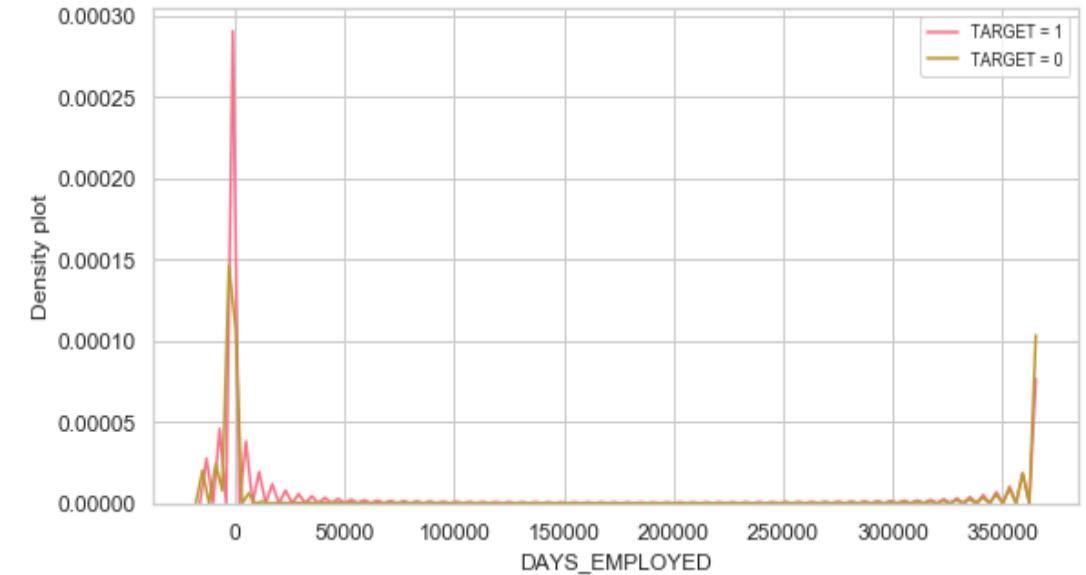
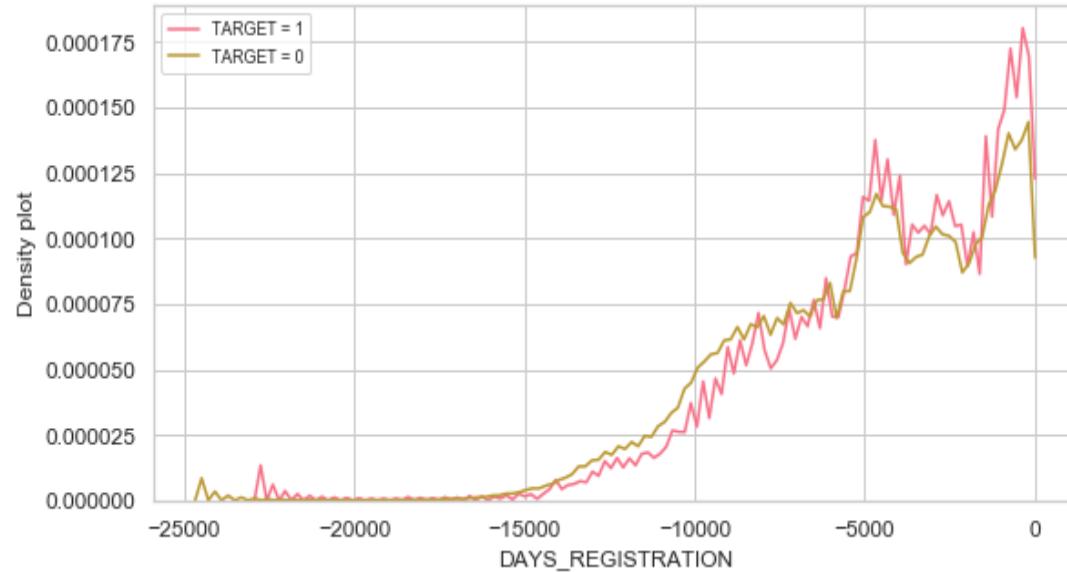




iit.b

UpGrad

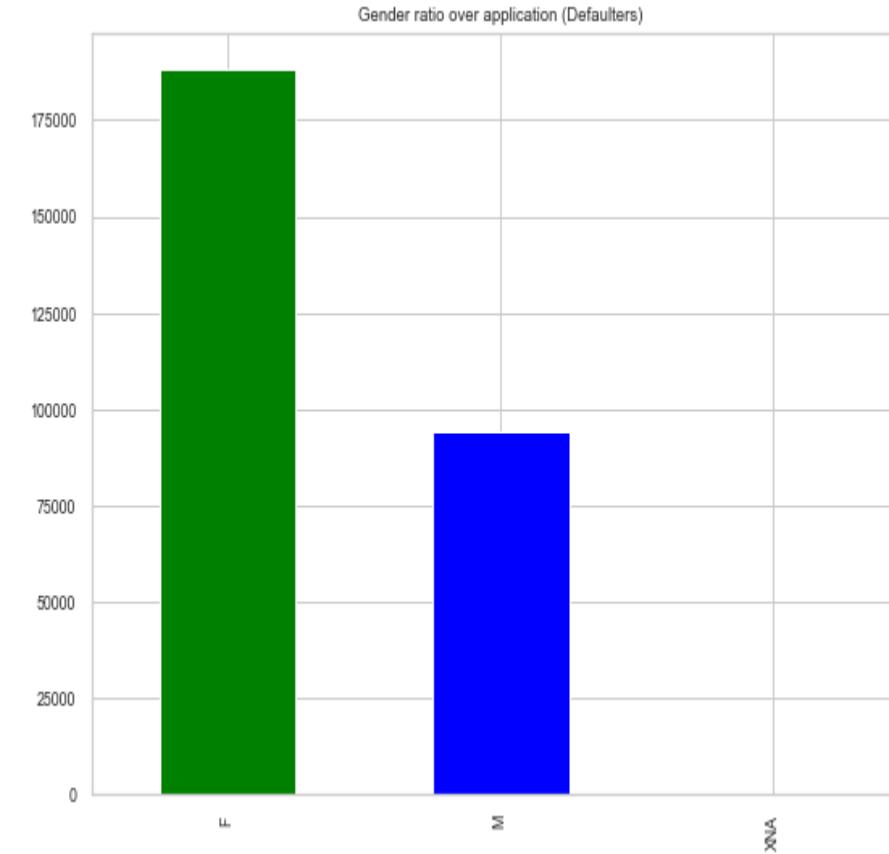
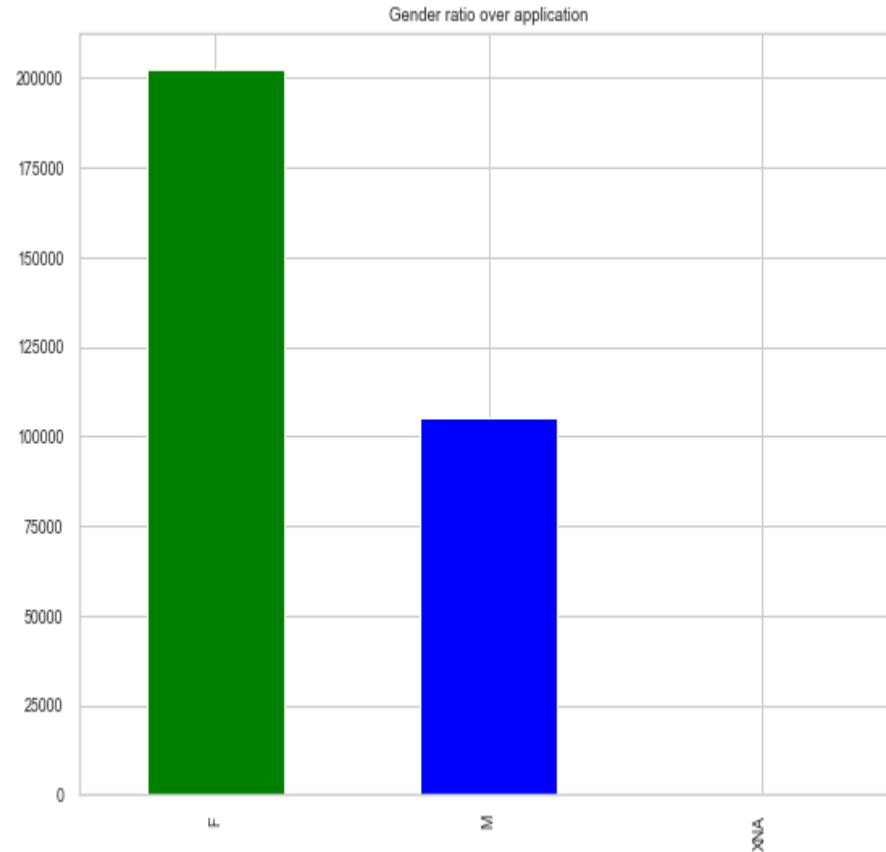




Insight: In every graph the defaulters are fluctuating mostly same as other variable fluctuate.

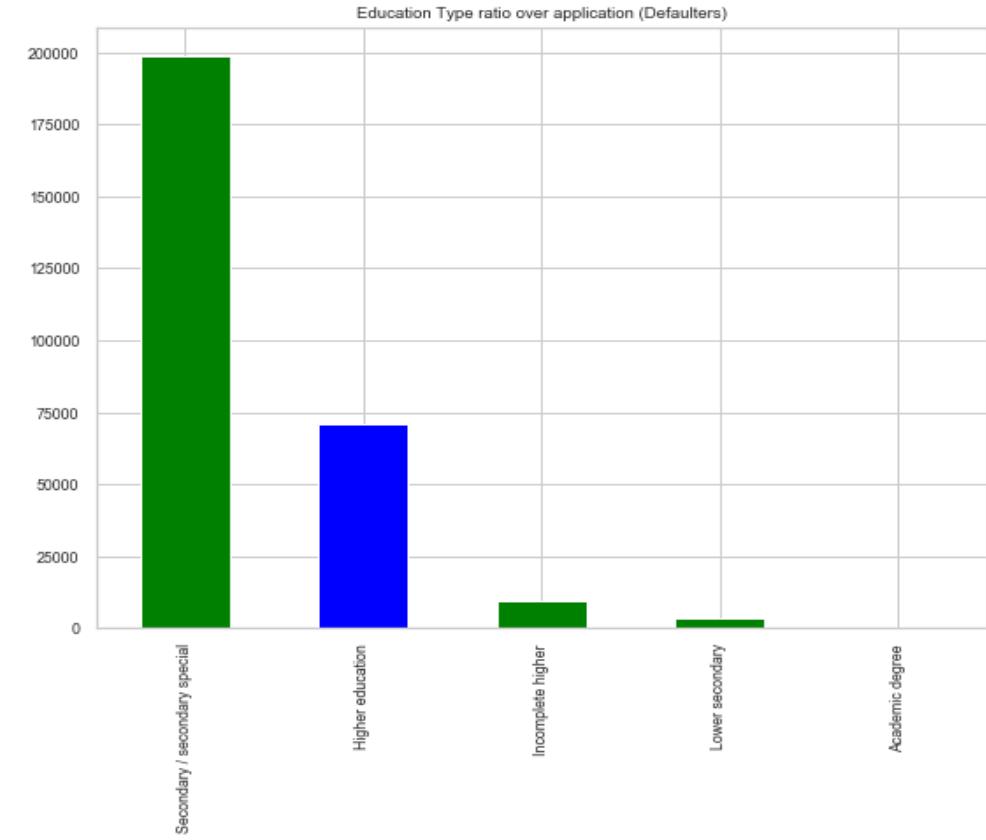
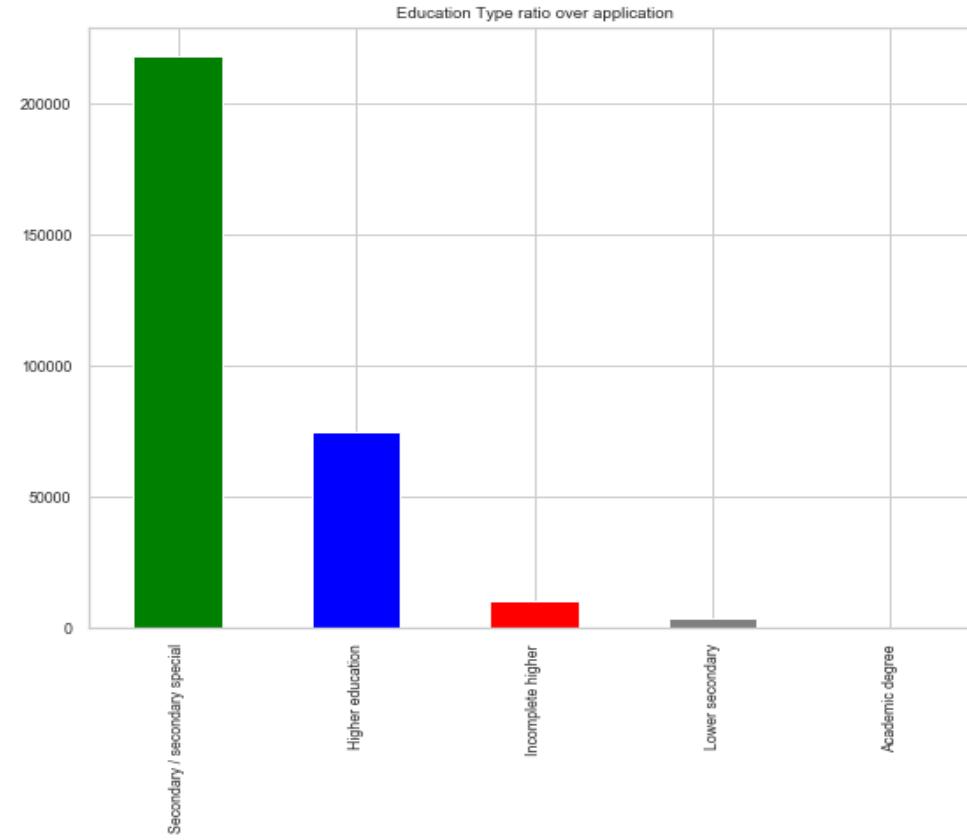
Ratio on the various categorical dimensions

Gender Ratio over application



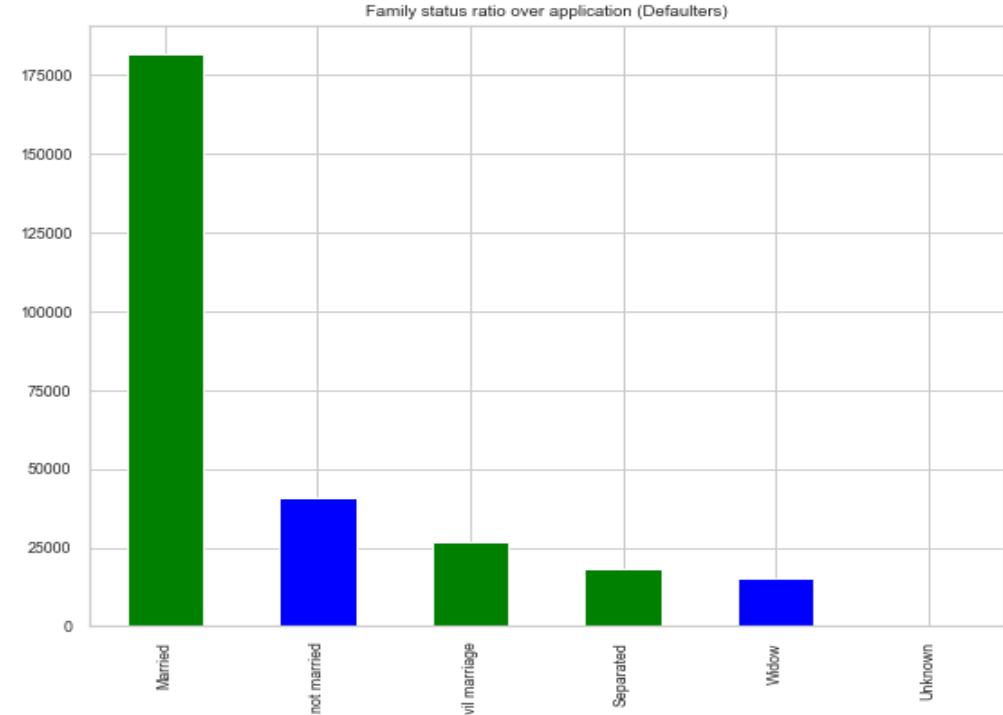
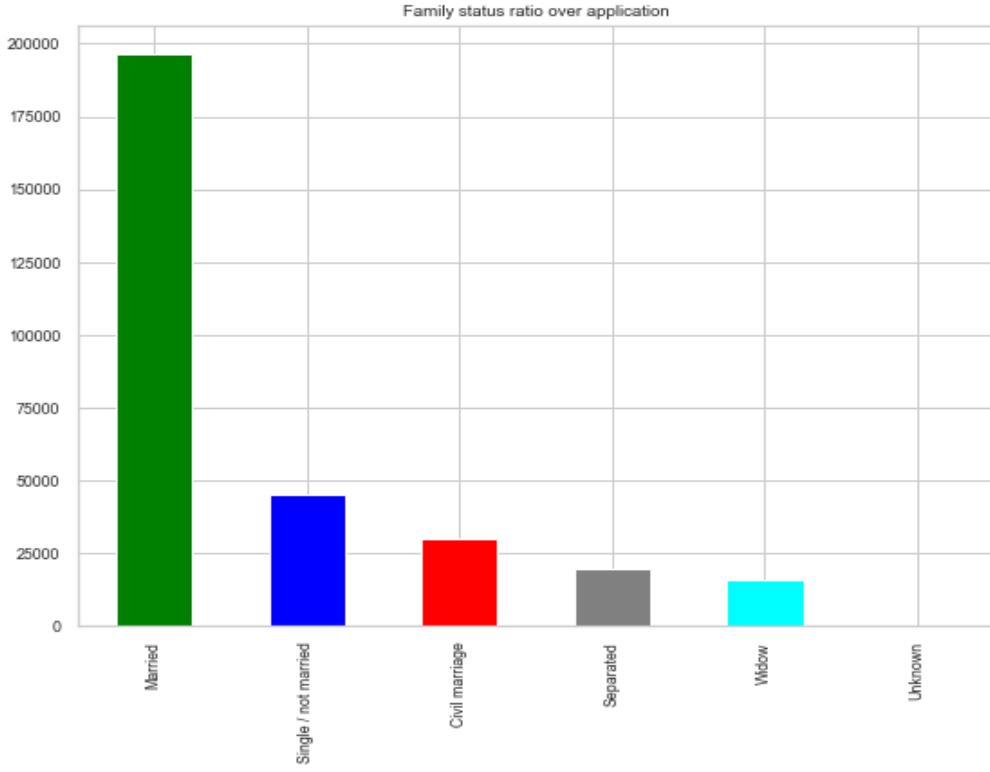
Insight: Female clients are most likely to be defaulters than the male clients

Education Ratio over application



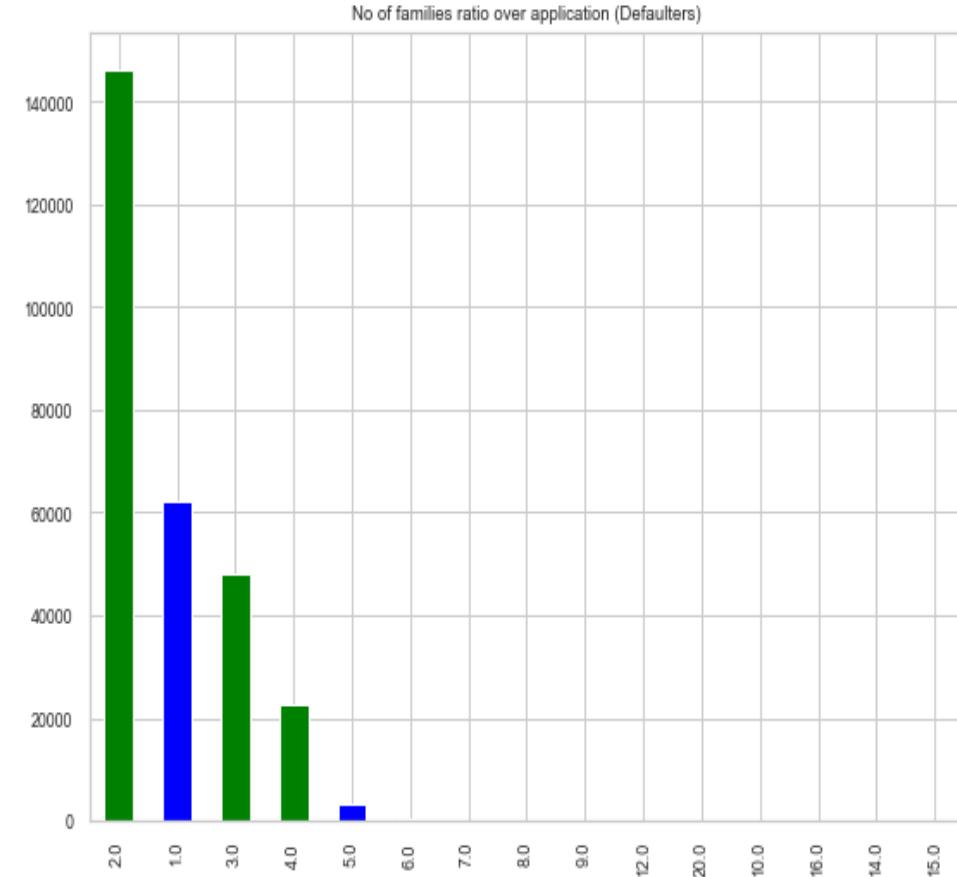
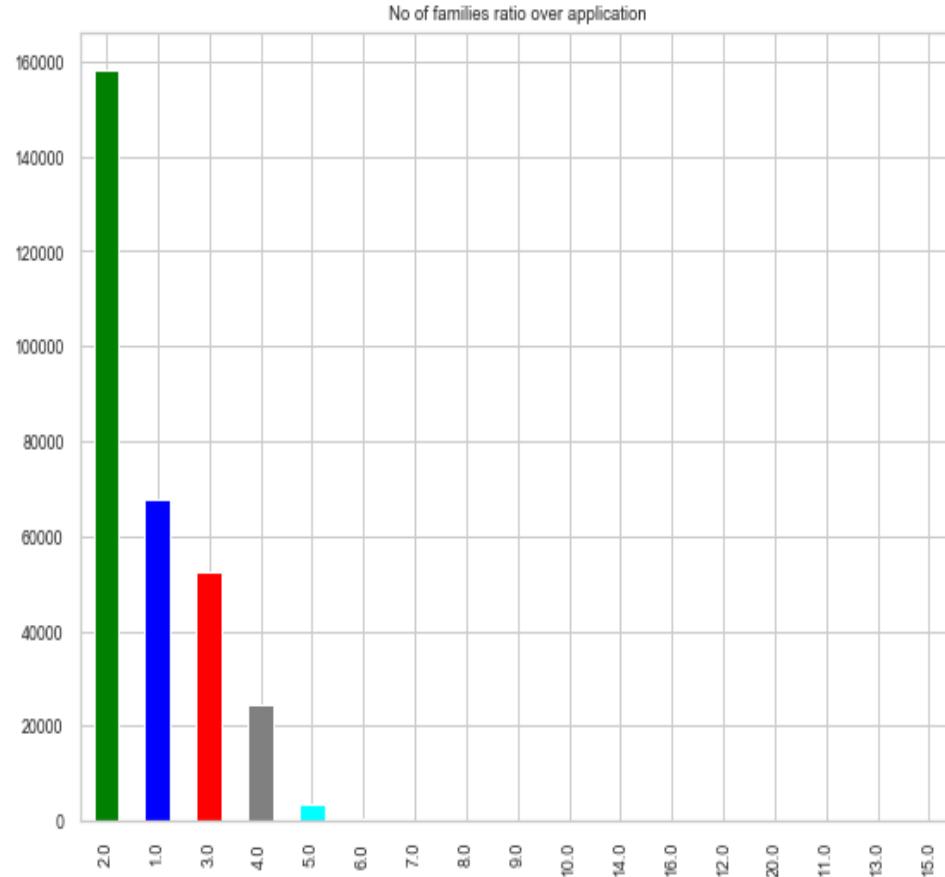
Insight: "Secondary/Secondary Special" Education is more likely to be defaulters.

Family Status Ratio over application



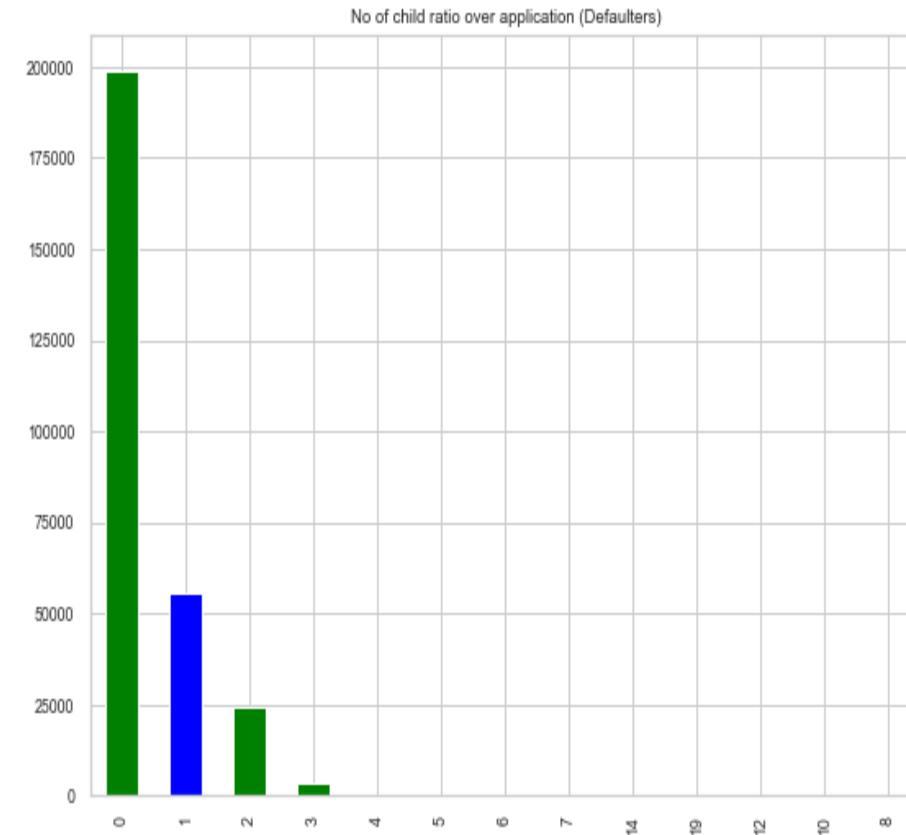
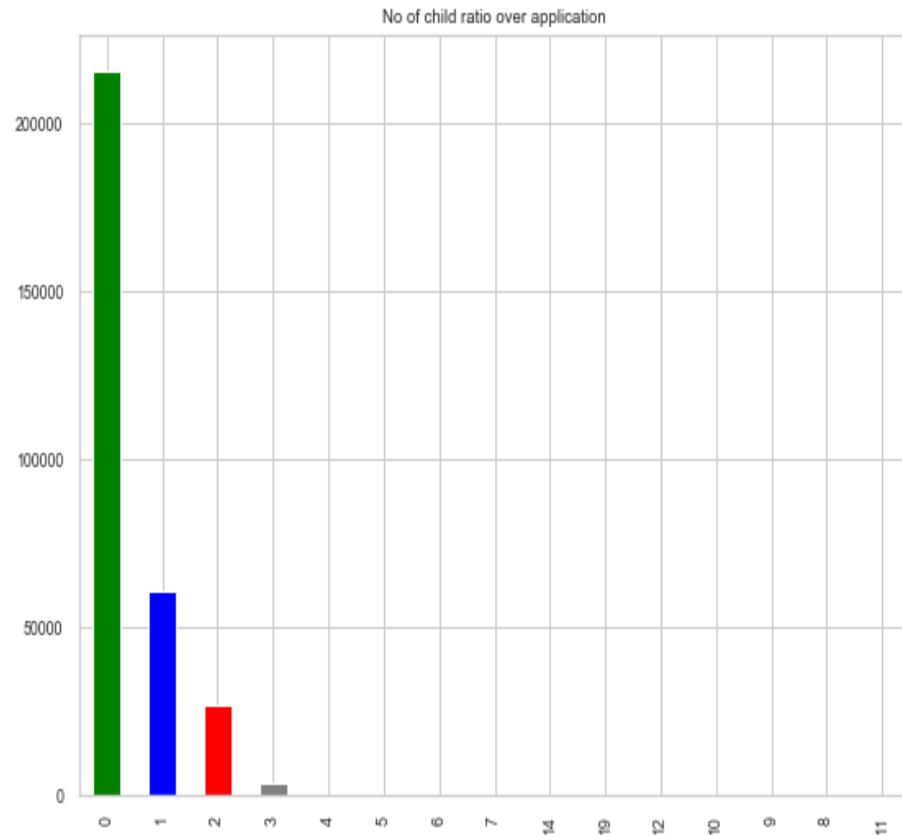
Insight: Married people are most likely to be defaulters

No of families ratio over application : CNT_FAM_MEMBERS



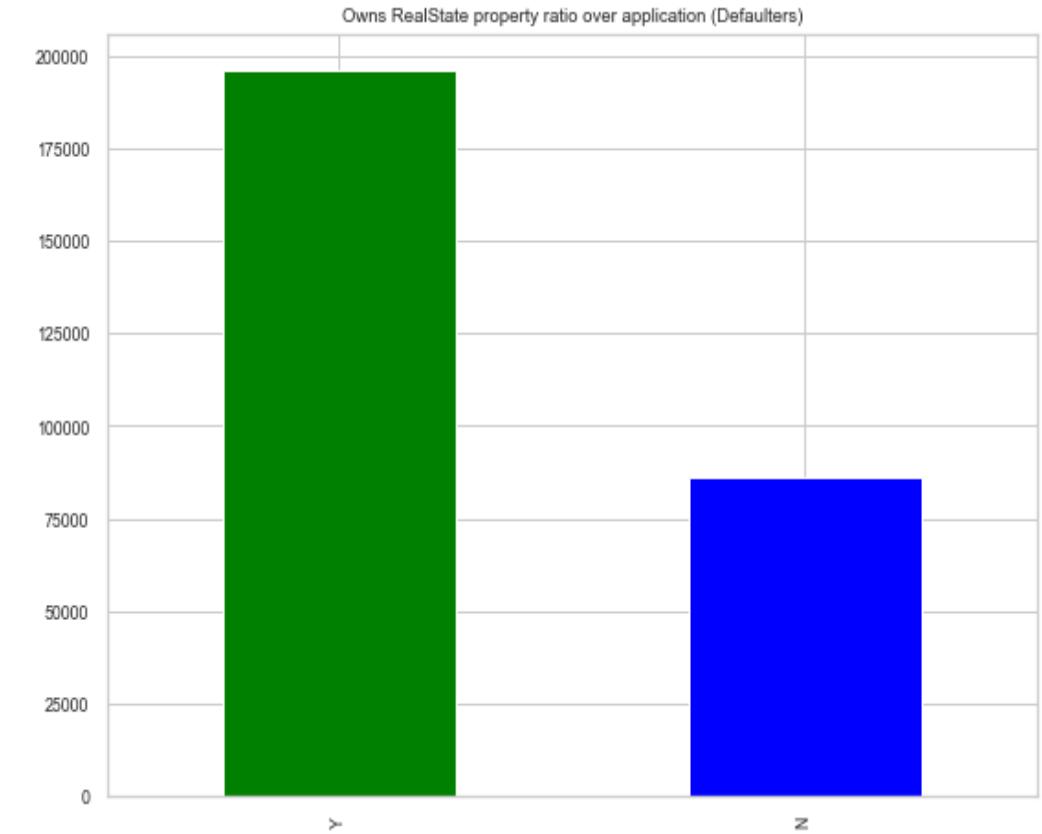
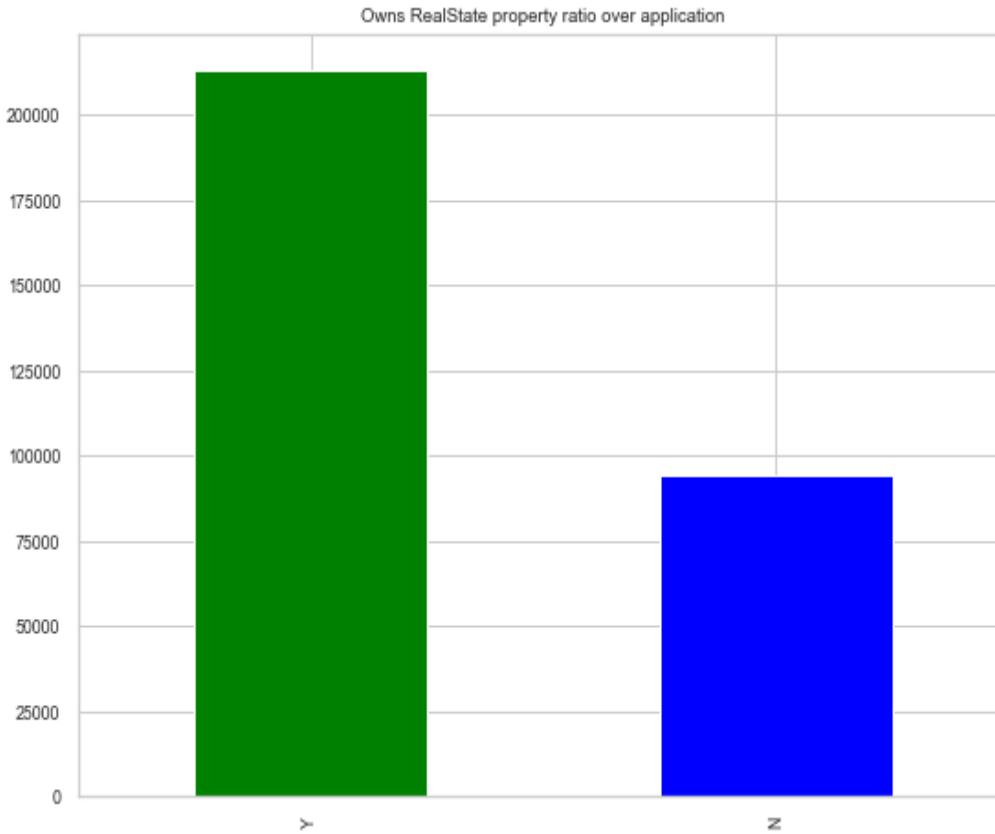
Insight: People who has two family members are most likely to be defaulters.

CNT_CHILDREN ratio over application



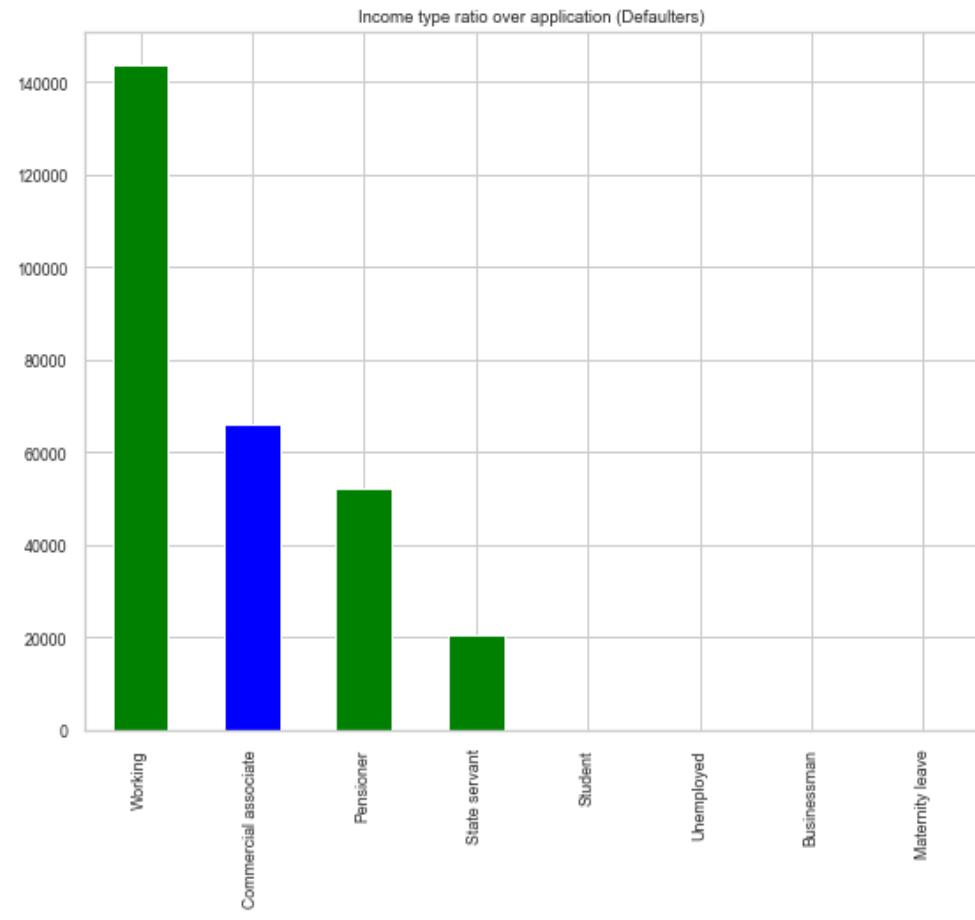
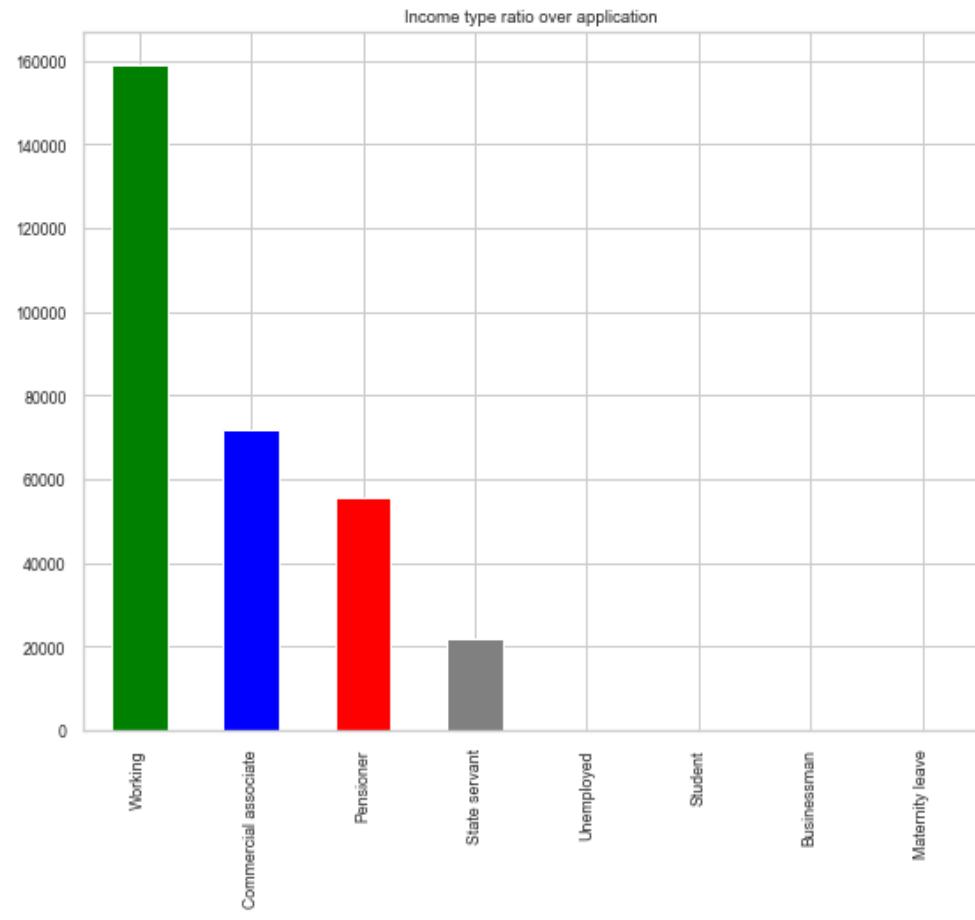
Insight: People who have no children in their family are most likely to be defaulters.

FLAG_own_realty



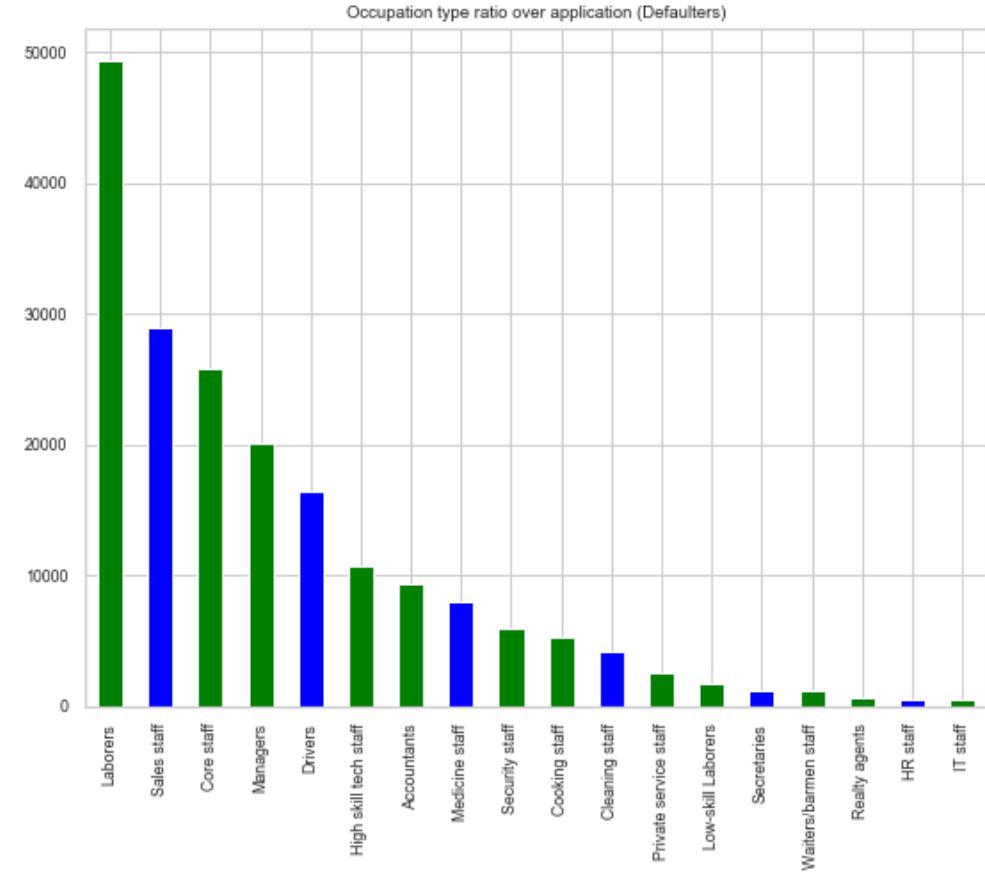
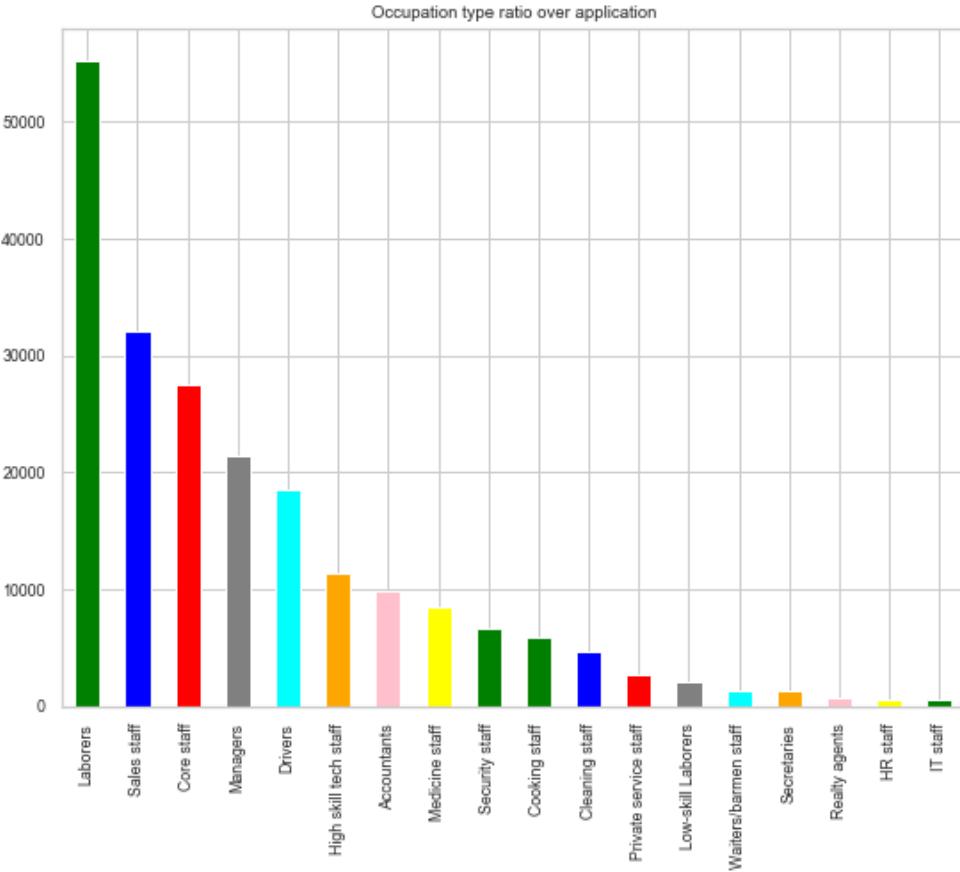
Insight: People who have Real State properties are most likely to be defaulters.

Income type ratio over application



Insight: People who are Working, are most likely to be defaulters.

OCCUPATION_TYPE ratio over application



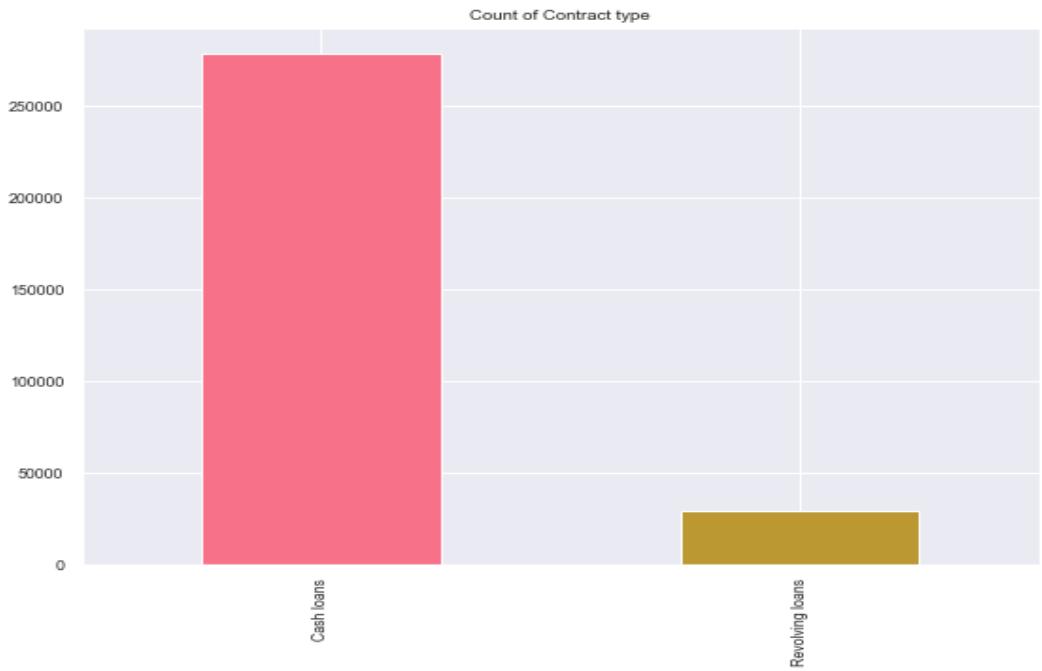
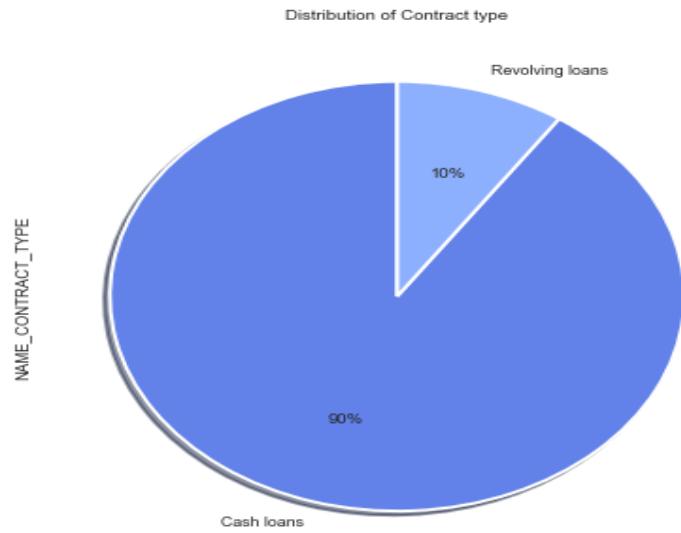
Insight: People who are Laborers, are most likely to be defaulters.

SEGMENTED UNIVARIATE ANALYSIS:

Univariate analysis is perhaps the simplest form of statistical analysis. Like other forms of statistics, it can be inferential or descriptive. The key fact is that only one variable is involved. Univariate analysis can yield misleading results in cases in which multivariate analysis is more appropriate

Insight: 10% out of total client population have difficulties in revolving loans.

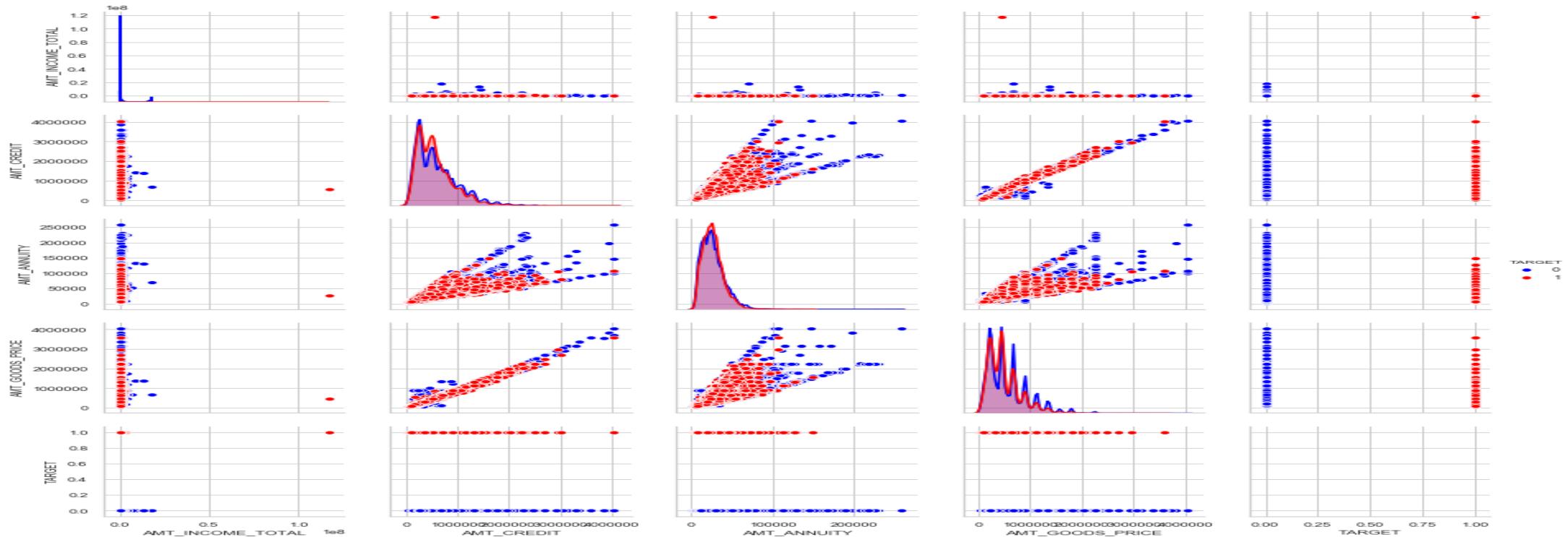
Insight: 90% are of cash loans.

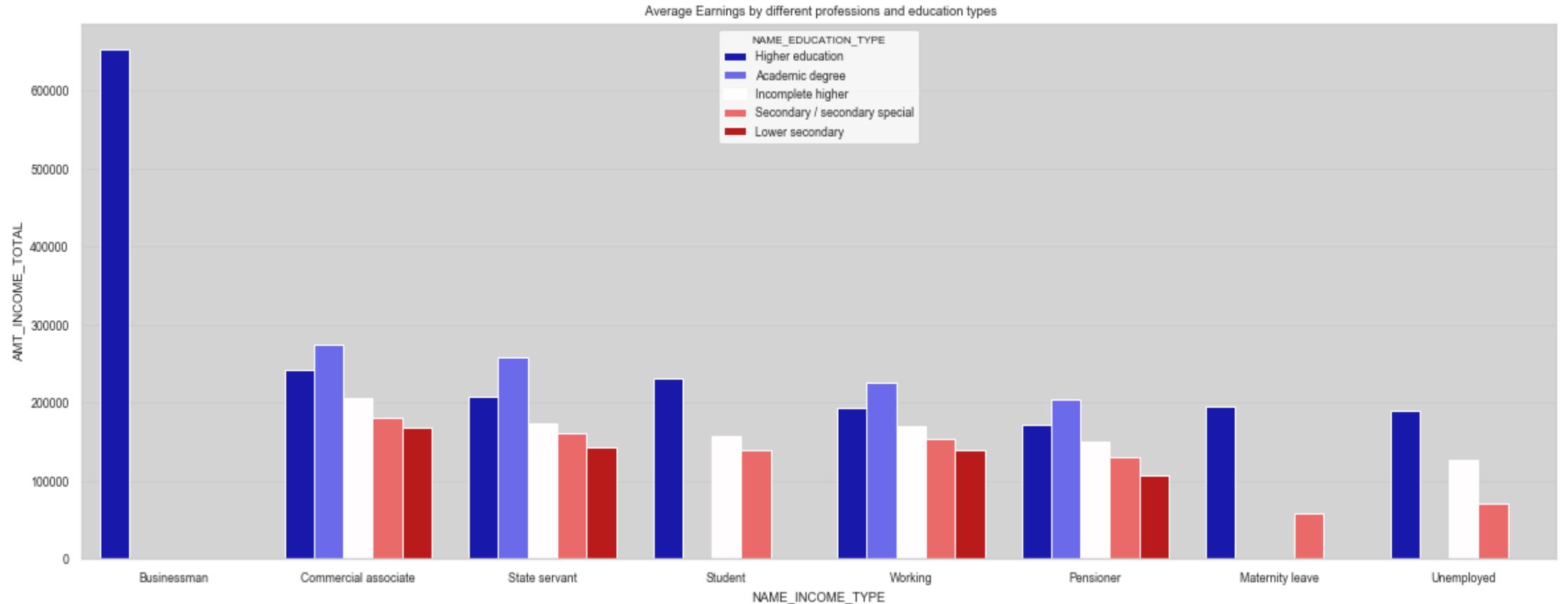


BIVARIATE ANALYSIS:

Bivariate analysis is the simultaneous analysis of two variables (attributes). It explores the concept of relationship between two variables, whether there exists an association and the strength of this association, or whether there are differences between two variables and the significance of these differences.

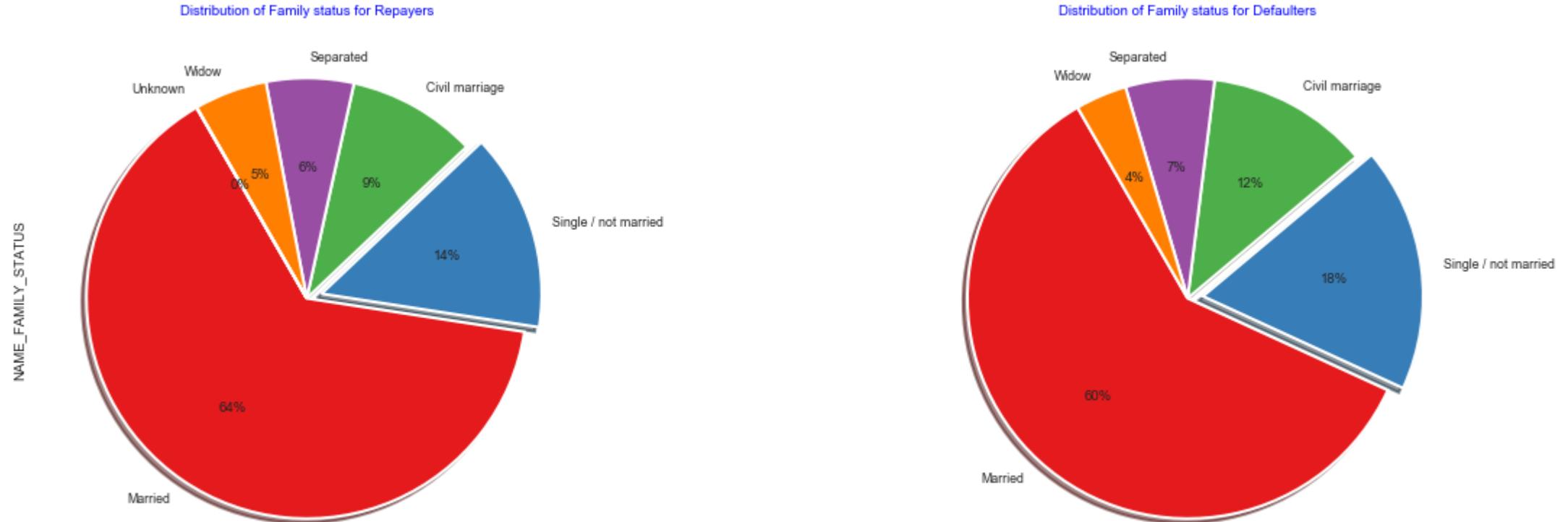
Bivariate Analysis For Amount





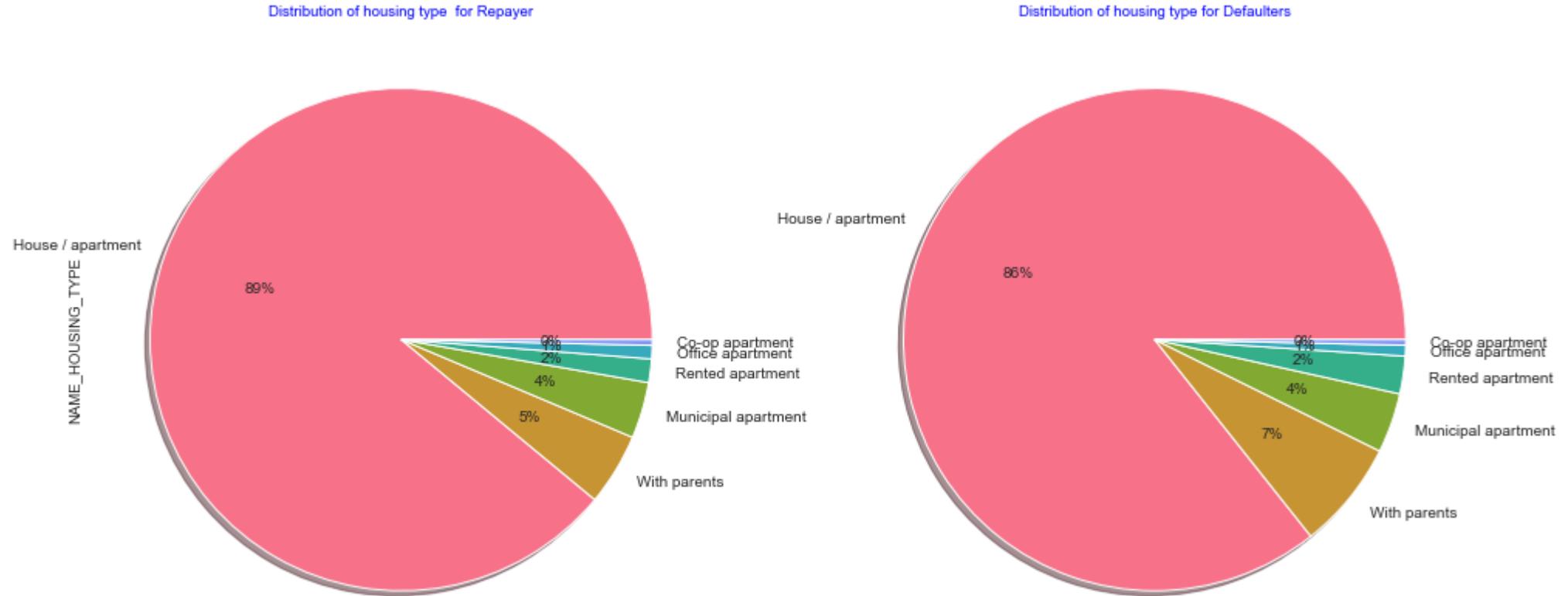
Insight: Businessman and Students who are having Higher education, are having more Income. For the rest, who are having Academic degree are having more income.

Distribution of Education type by loan repayment status



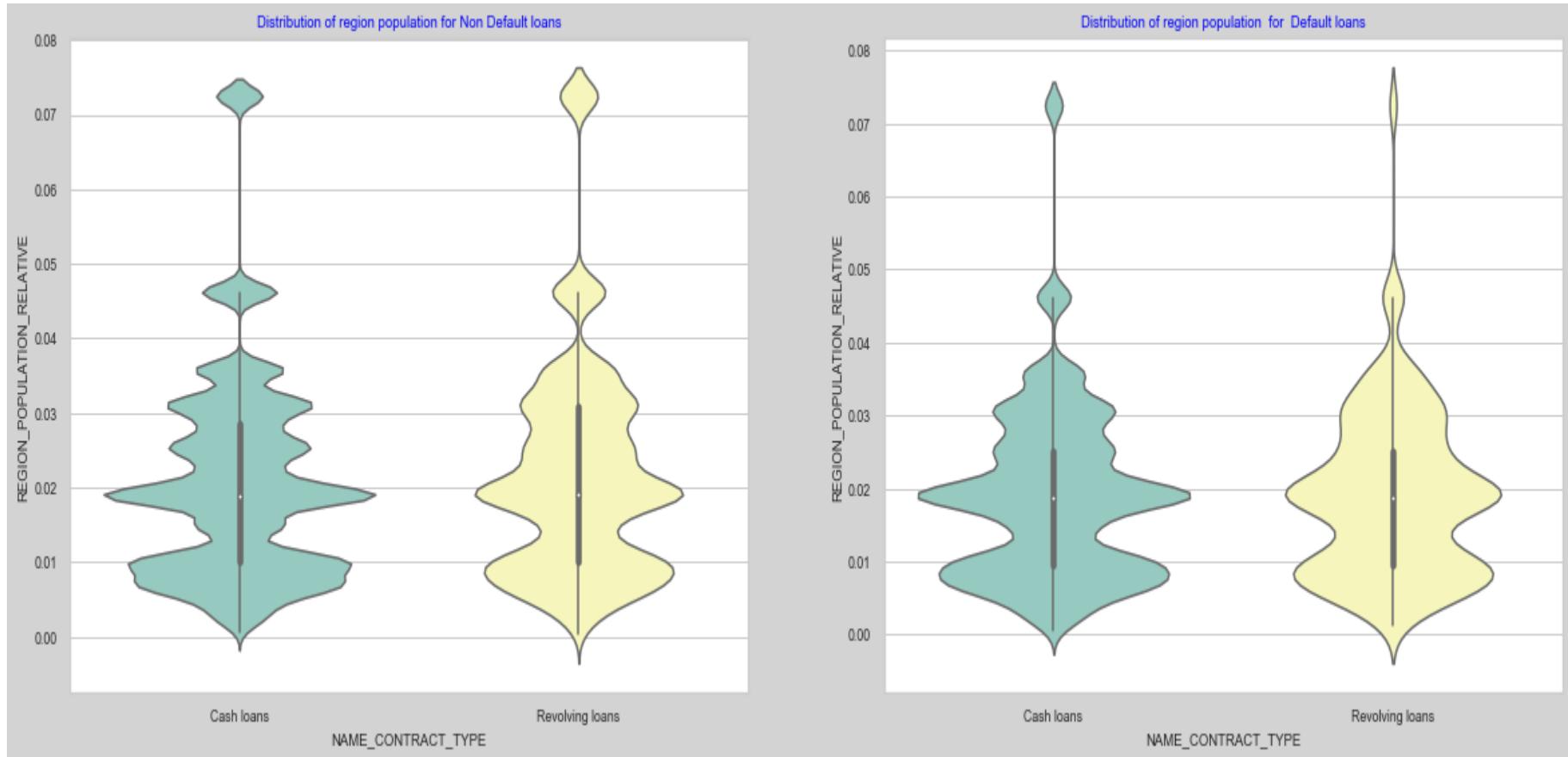
Insight: People with 60% who are married, are most likely to be defaulters.

Distribution of housing type for Defaulters



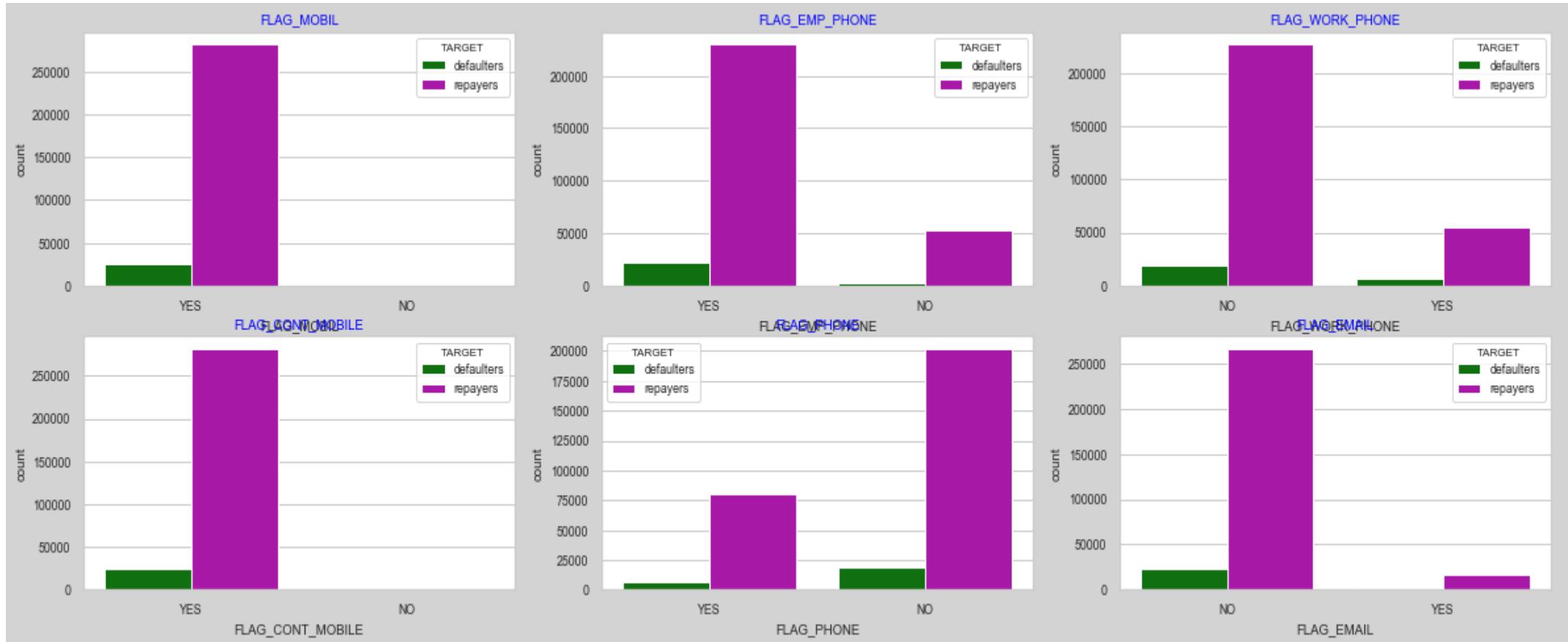
Insight: People(mostly 86%) who lives in House/Apartment, are most likely to be defaulters.

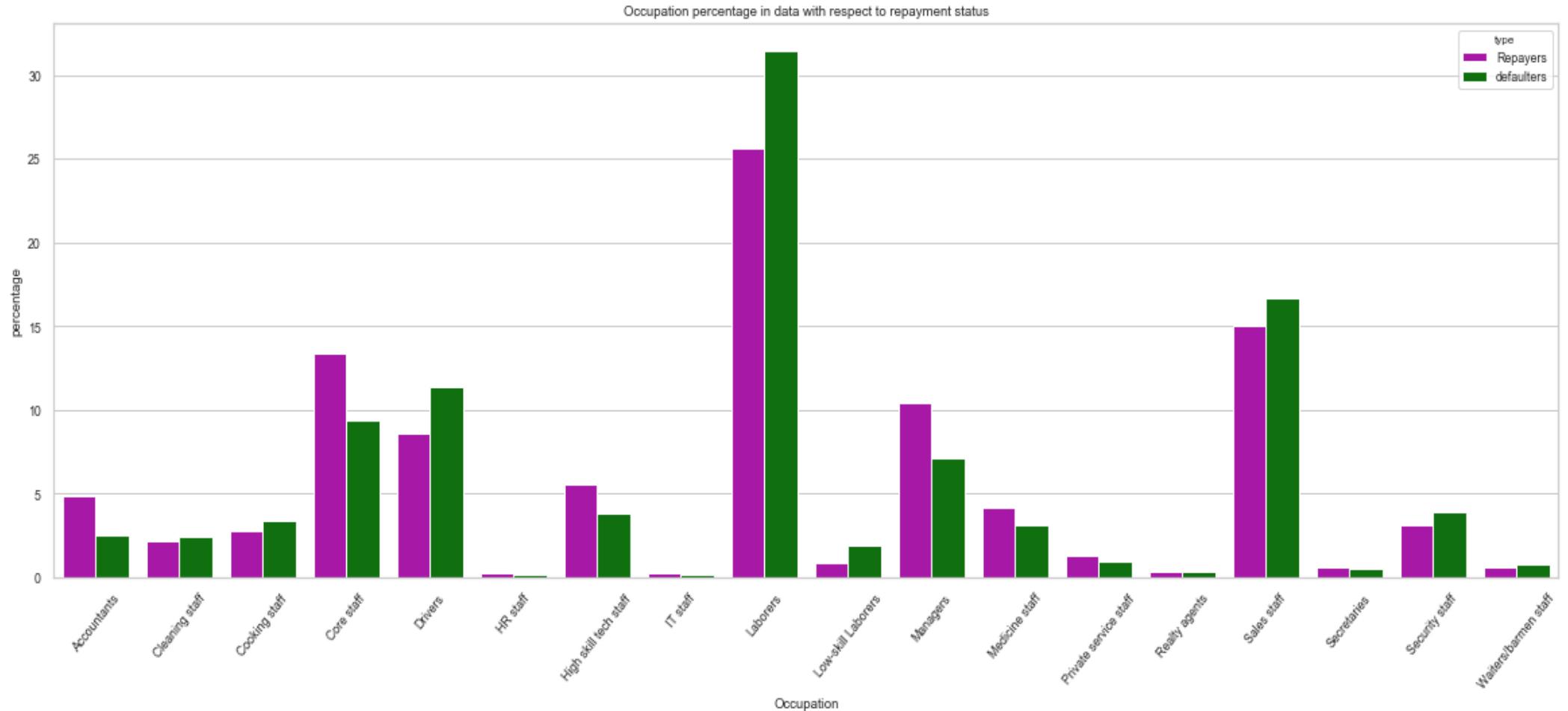
Distribution normalized population of region where client lives by loan repayment status



Insight: In High population density regions people are less likely to default on loans.

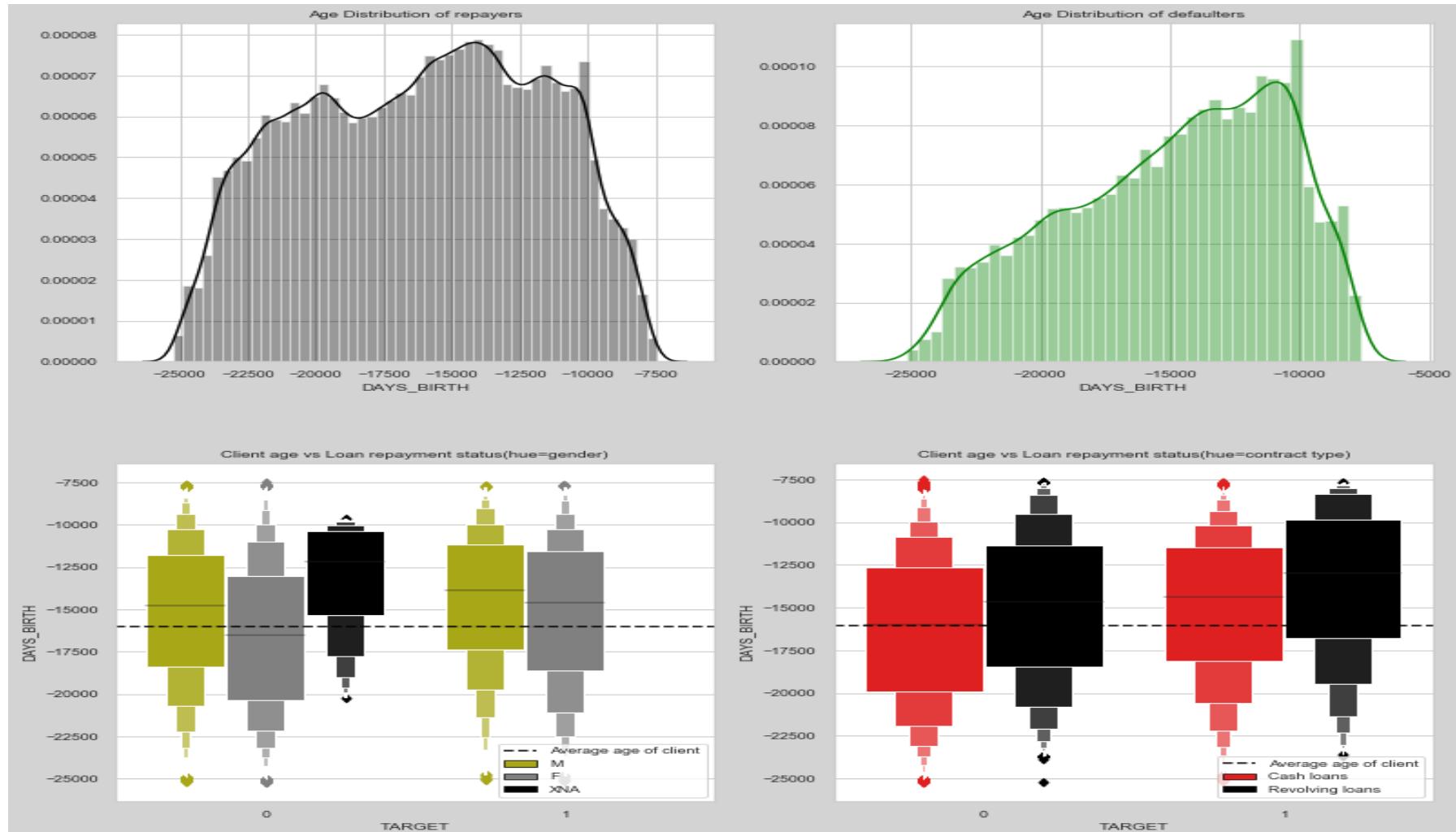
Distribution of contact information provided by client





Insight: Occupations like Cleaning staff ,Cooking staff, Drivers ,Laborers , Low-skill Laborers ,Sales staff, Security staff are more likely to default in loans.

Distribution of Client's age



Insight: As per the graph we observe that people with middle and old ages are repayers and mostly old age people are defaulters

CORRELATION

In statistics, dependence or association is any statistical relationship, whether causal or not, between two random variables or bivariate data. In the broadest sense correlation is any statistical association, though it commonly refers to the degree to which a pair of variables are linearly related.

Correlation between variables (Categorized by "TARGET" variable)

Getting Correlation from *New Loan Application*

| | SK_ID_CURR | TARGET | CNT_CHILDREN | AMT_INCOME_TOTAL | AMT_CREDIT | AMT_ANNUITY | AMT_GOODS_PRICE |
|------------------|------------|-----------|--------------|------------------|------------|-------------|-----------------|
| SK_ID_CURR | 1.000000 | -0.002108 | -0.001129 | -0.001820 | -0.000343 | -0.000433 | -0.000232 |
| TARGET | -0.002108 | 1.000000 | 0.019187 | -0.003982 | -0.030369 | -0.012817 | -0.039645 |
| CNT_CHILDREN | -0.001129 | 0.019187 | 1.000000 | 0.012882 | 0.002145 | 0.021374 | -0.001827 |
| AMT_INCOME_TOTAL | -0.001820 | -0.003982 | 0.012882 | 1.000000 | 0.156870 | 0.191657 | 0.159610 |
| AMT_CREDIT | -0.000343 | -0.030369 | 0.002145 | 0.156870 | 1.000000 | 0.770138 | 0.986968 |

5 rows × 8 columns

Getting Correlation from *Previous Loan Application*

| | SK_ID_PREV | SK_ID_CURR | AMT_ANNUITY | AMT_APPLICATION | AMT_CREDIT | AMT_DOWN_PAYMENT | AMT_GOODS_PRICE |
|-----------------|------------|------------|-------------|-----------------|------------|------------------|-----------------|
| SK_ID_PREV | 1.000000 | -0.000321 | 0.011459 | 0.003302 | 0.003659 | -0.001313 | 0.015293 |
| SK_ID_CURR | -0.000321 | 1.000000 | 0.000577 | 0.000280 | 0.000195 | -0.000063 | 0.000369 |
| AMT_ANNUITY | 0.011459 | 0.000577 | 1.000000 | 0.808872 | 0.816429 | 0.267694 | 0.820895 |
| AMT_APPLICATION | 0.003302 | 0.000280 | 0.808872 | 1.000000 | 0.975824 | 0.482776 | 0.999884 |
| AMT_CREDIT | 0.003659 | 0.000195 | 0.816429 | 0.975824 | 1.000000 | 0.301284 | 0.993087 |

5 rows × 8 columns



HEAT MAP FOR CORRELATION DATA ON NEW APPLICATION

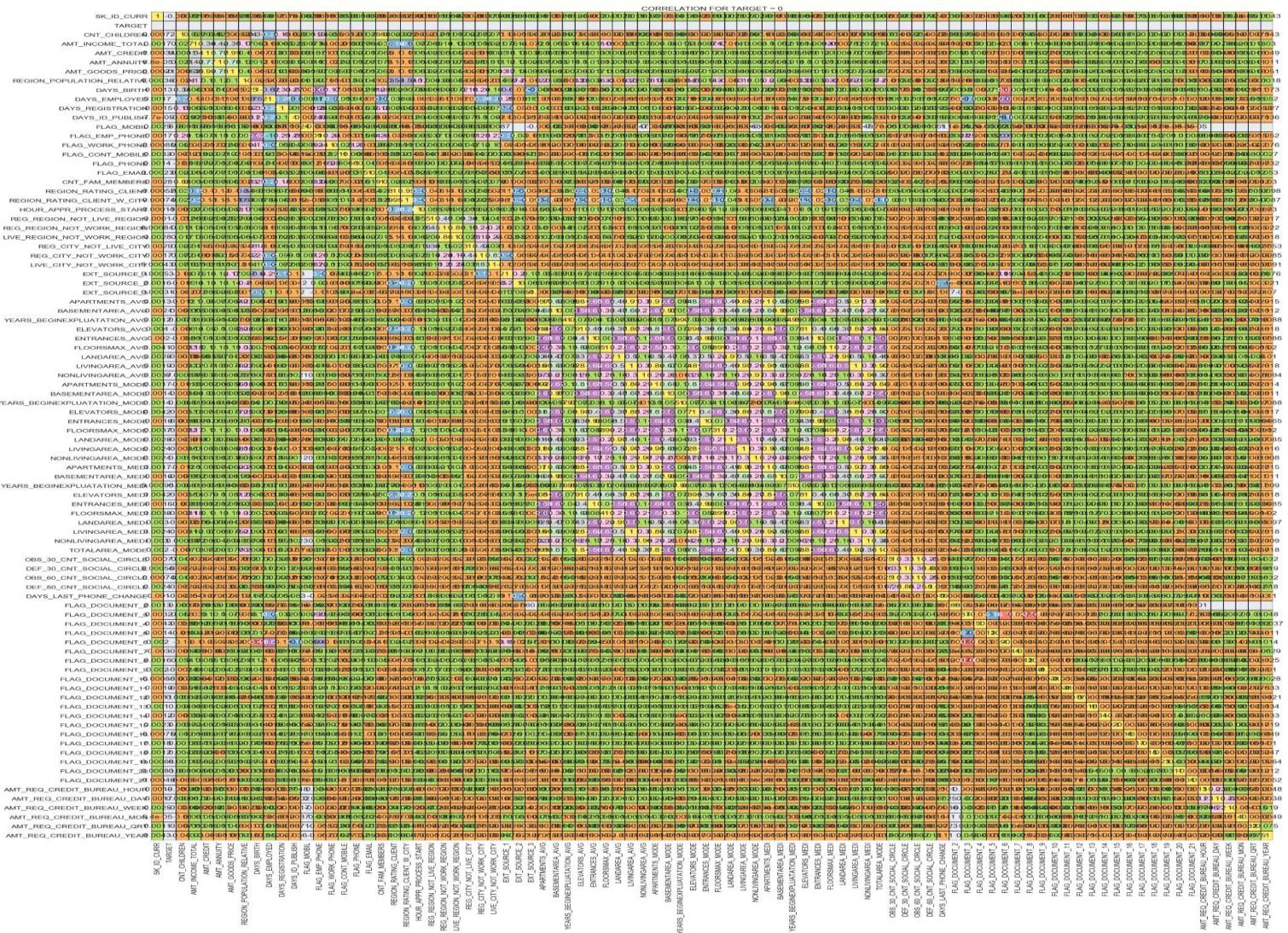
UpGrad





Categorizing data by "TARGET" variable

HEAT MAP FOR CORRELATION DATA ON TARGET=0

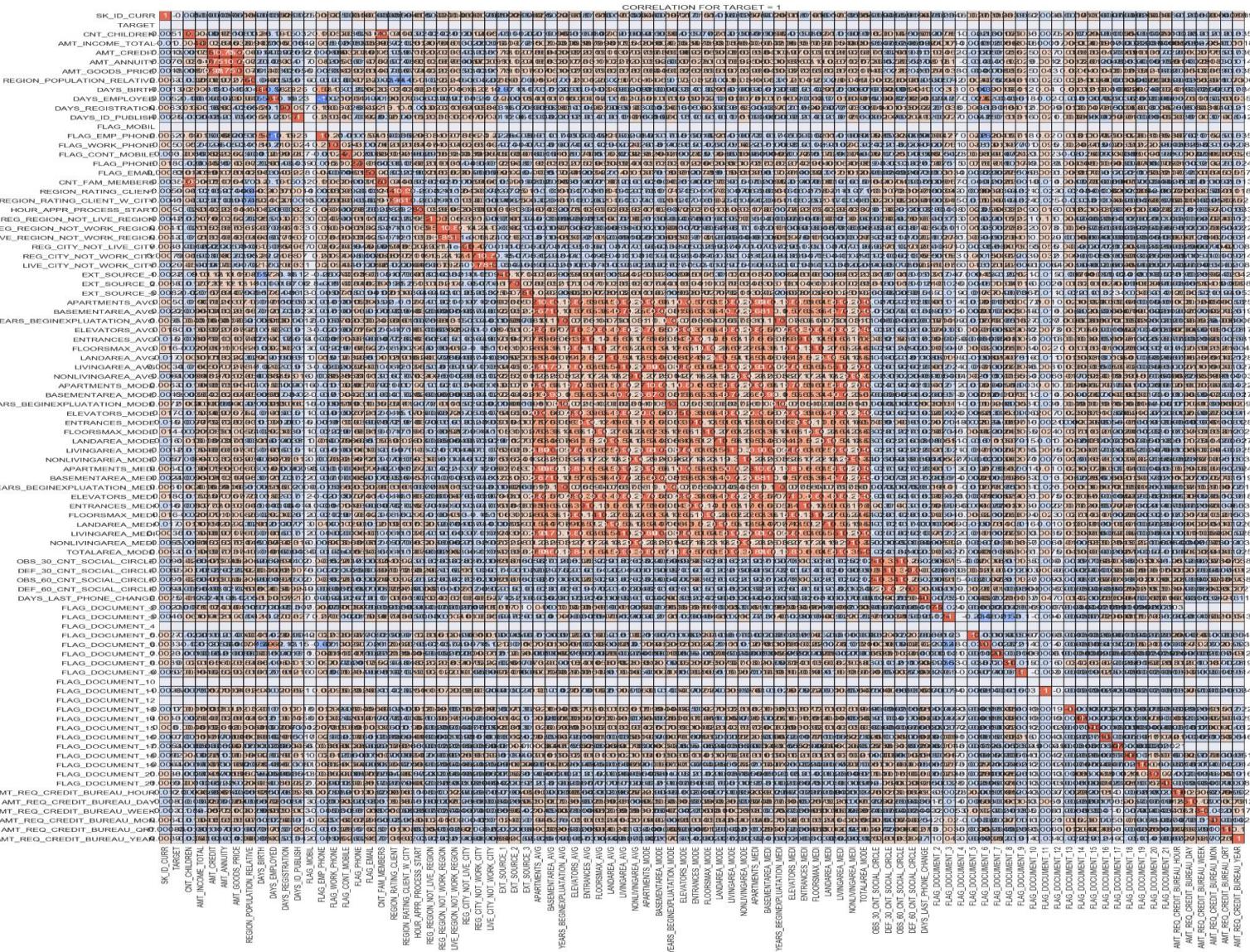


UpGrad



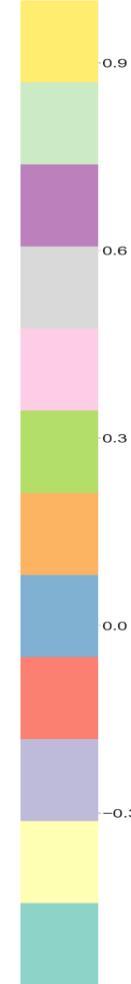
HEAT MAP FOR CORRELATION DATA ON TARGET=1

UpGrad



CORREALTION FOR SPECIFIC VARIABLES

| CORREALTION FOR SPECIFIC VARIABLES | | | | | | | | | | |
|------------------------------------|------------------|------------|-------------|-----------------|-----------------|--------------|-------------------|---------------|------------|--------------|
| | AMT_INCOME_TOTAL | AMT_CREDIT | AMT_ANNUITY | AMT_GOODS_PRICE | CNT_FAM_MEMBERS | CNT_CHILDREN | DAYS_REGISTRATION | DAYS_EMPLOYED | DAYS_BIRTH | EXT_SOURCE_1 |
| AMT_INCOME_TOTAL | 1 | 0.038 | 0.046 | 0.038 | 0.0067 | 0.0048 | 0.00016 | -0.015 | 0.0031 | 0.01 |
| AMT_CREDIT | 0.038 | 1 | 0.75 | 0.98 | 0.051 | -0.0017 | -0.026 | -0.00097 | -0.14 | 0.18 |
| AMT_ANNUITY | 0.046 | 0.75 | 1 | 0.75 | 0.076 | 0.031 | 0.034 | -0.083 | -0.014 | 0.11 |
| AMT_GOODS_PRICE | 0.038 | 0.98 | 0.75 | 1 | 0.047 | -0.0081 | -0.026 | 0.0036 | -0.14 | 0.19 |
| CNT_FAM_MEMBERS | 0.0067 | 0.051 | 0.076 | 0.047 | 1 | 0.89 | 0.15 | -0.19 | 0.2 | -0.051 |
| CNT_CHILDREN | 0.0048 | -0.0017 | 0.031 | -0.0081 | 0.89 | 1 | 0.15 | -0.19 | 0.26 | -0.091 |
| DAYS_REGISTRATION | 0.00016 | -0.026 | 0.034 | -0.026 | 0.15 | 0.15 | 1 | -0.19 | 0.29 | -0.16 |
| DAYS_EMPLOYED | -0.015 | -0.00097 | -0.083 | 0.0036 | -0.19 | -0.19 | -0.19 | 1 | -0.58 | 0.28 |
| DAYS_BIRTH | 0.0031 | -0.14 | -0.014 | -0.14 | 0.2 | 0.26 | 0.29 | -0.58 | 1 | -0.57 |
| EXT_SOURCE_1 | 0.01 | 0.18 | 0.11 | 0.19 | -0.051 | -0.091 | -0.16 | 0.28 | -0.57 | 1 |



Highly Correlated Variables

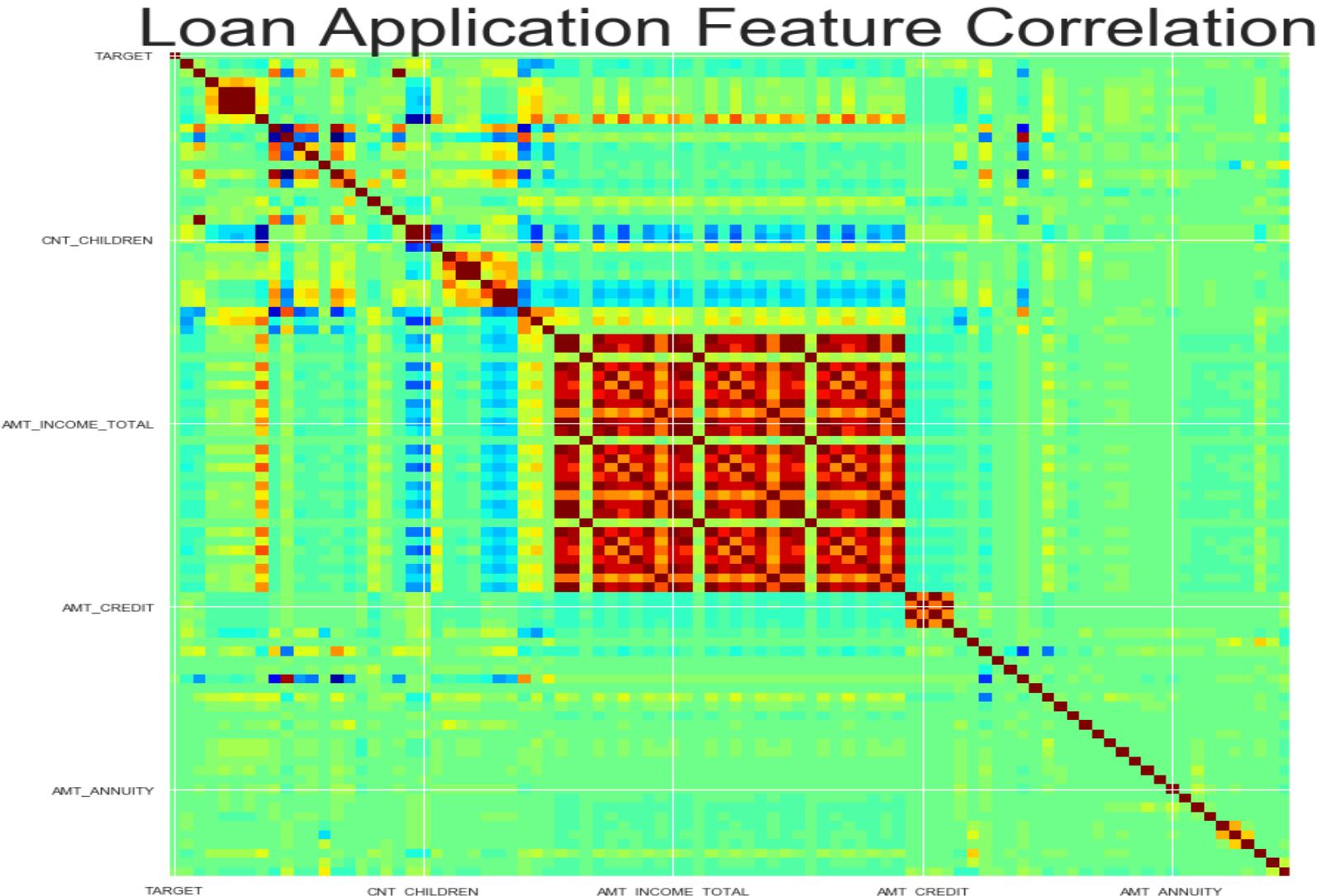
TOP 4 CORRELATED VARIABLES

| | |
|-----------------|-----------------|
| AMT_CREDIT | AMT_ANNUITY |
| AMT_ANNUITY | AMT_GOODS_PRICE |
| AMT_GOODS_PRICE | AMT_CREDIT |
| CNT_FAM_MEMBERS | CNT_CHILDREN |

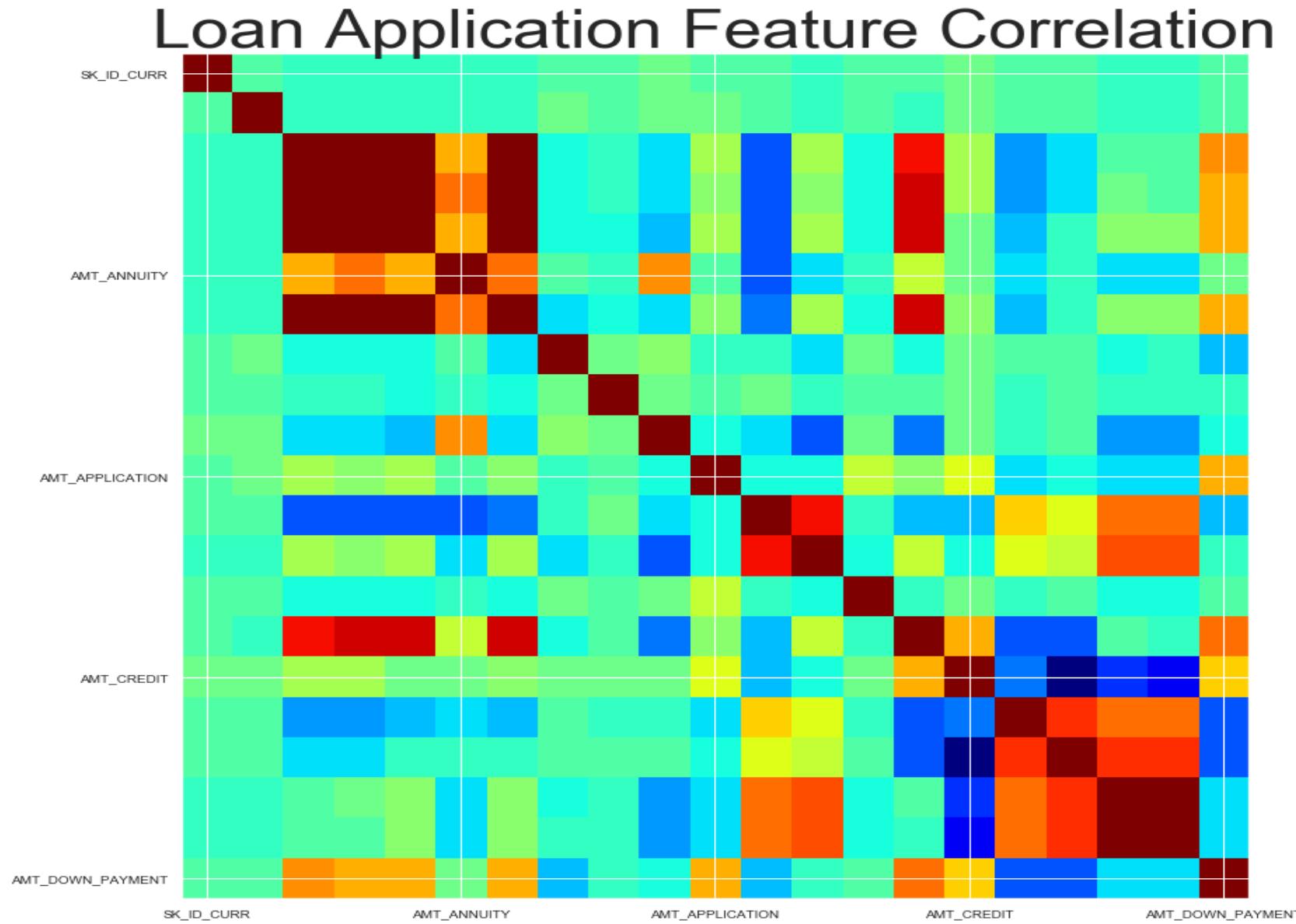
TOP 7 Correlated Variables:

1. AMT_CREDIT
2. AMT_ANNUITY
3. AMT_GOODS_PRICE
4. CNT_FAM_MEMBERS
5. DAYS_EMPLOYED
6. DAYS_BIRTH
7. EXT_SOURCE_1

```
# Plot the Heat Map For New Loan Application
```



```
# Plot the correlation pattern with New Loan Application
```

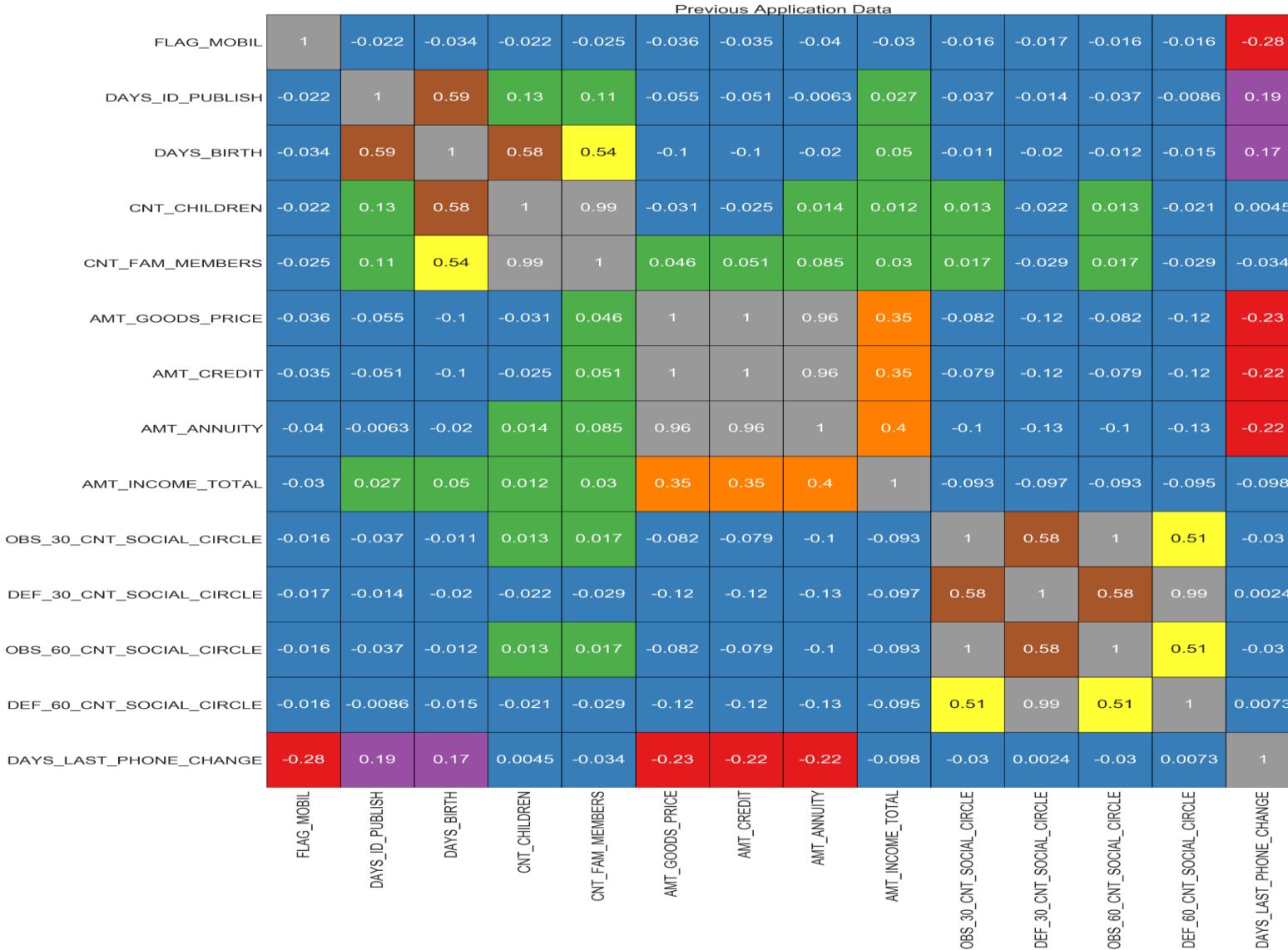




UpGrad

Extracting Highly correlated variables from New Loan Applications

Plotting the correlation pattern with Specific variables of Previous Application





UpGrad