

EDA CREDIT CASE STUDAY

BY

GAURAV SAINI



Purpose

There are two types of risks associated with the bank's decision:

- If the applicant is likely to repay the loan, then not approving the loan results in a loss of business to the company
- If the applicant is not likely to repay the loan, i.e. he/she is likely to default, then approving the loan may lead to a financial loss for the company.

Analysis of the data set has been done in python on a Jupiter Notebook

Approach

Steps :

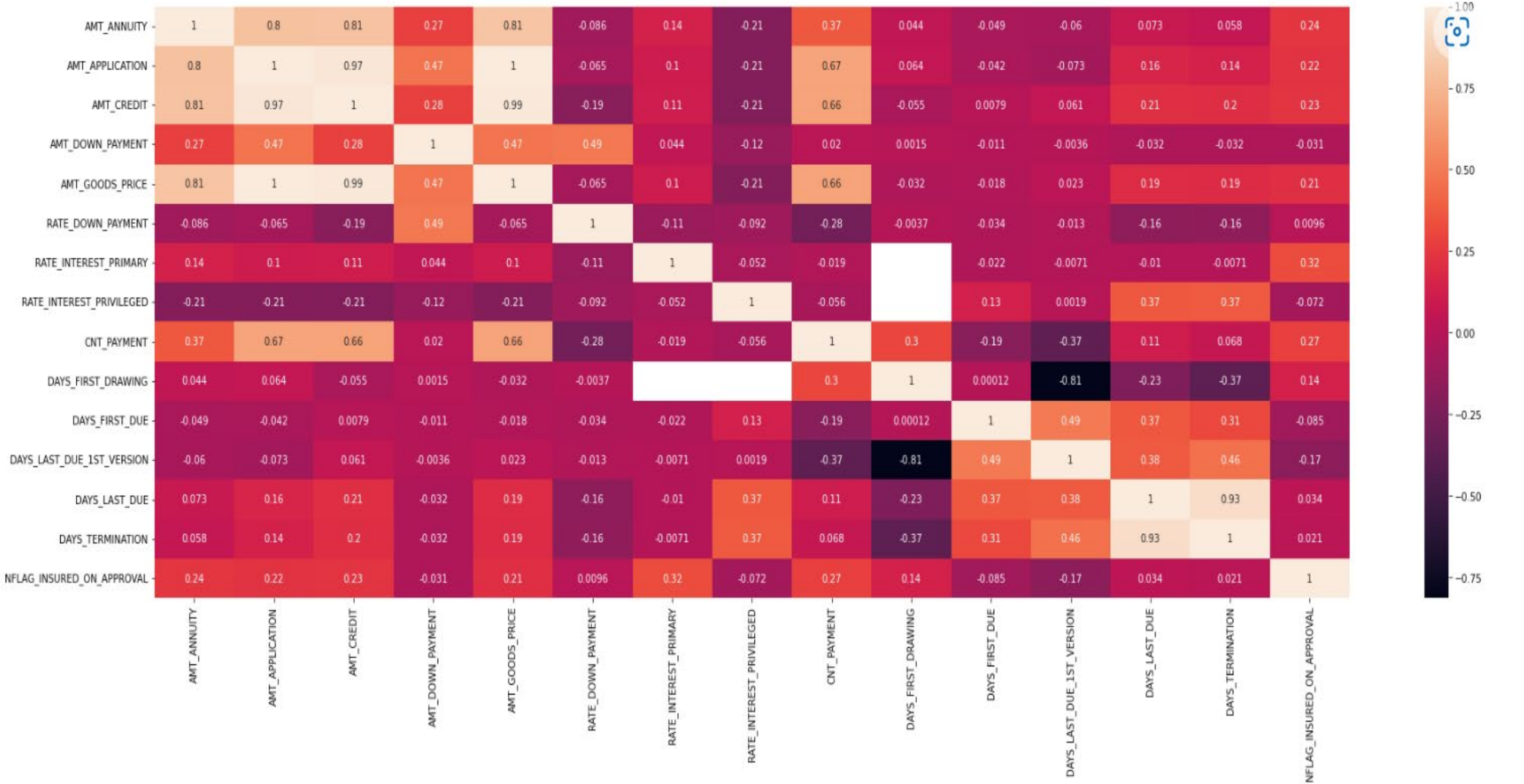
- *Importing Module***
- *Reading the Dataset into Pandas DataFrame***
- *Checking missing values and handling them***
- *Checking outliers, data imbalance, ration of imbalance***
- *We have divided the features into small segments and analyzed segment wise using a smaller dataframe containing only relevant categories.***
- *Plots and percentage wise Defaulters calculation are done segment wise as well.***

Correlation Between Numeric Feature of Previous Application Data

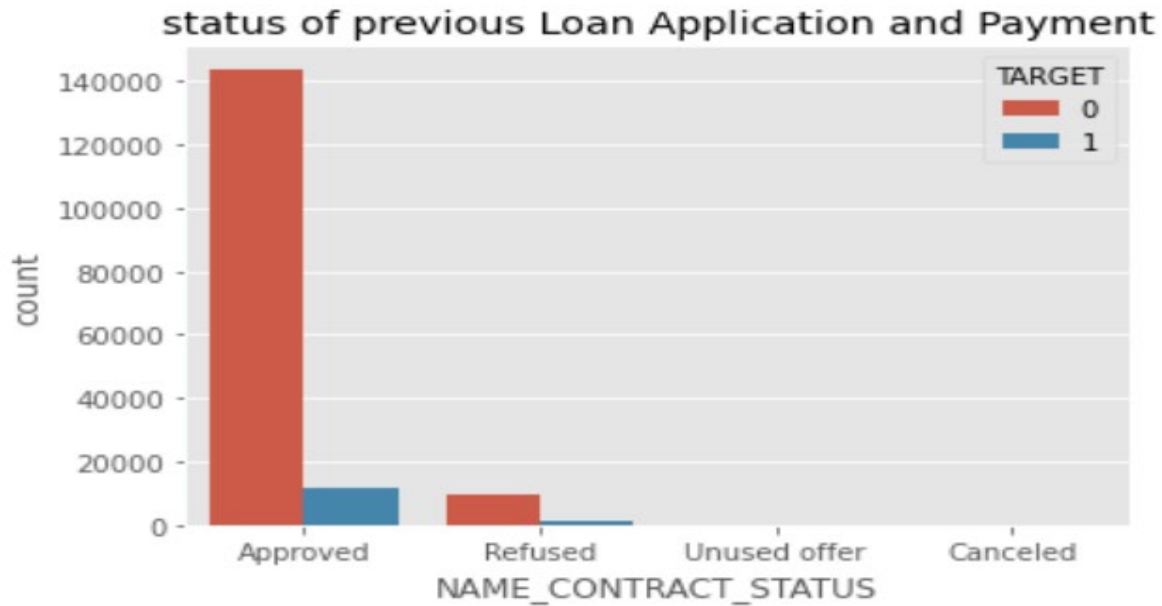
→ As you can see below(next slide) the lighter shade means very highly correlated, therefore DAYS_LAST_DUE and DAYS_TERMINATION are highly correlated

→ similarly AMT_ANNUTY, AMT_APPLICATION, AMT_CREDIT, AMT_GGODS_PRICE are highly correlated.

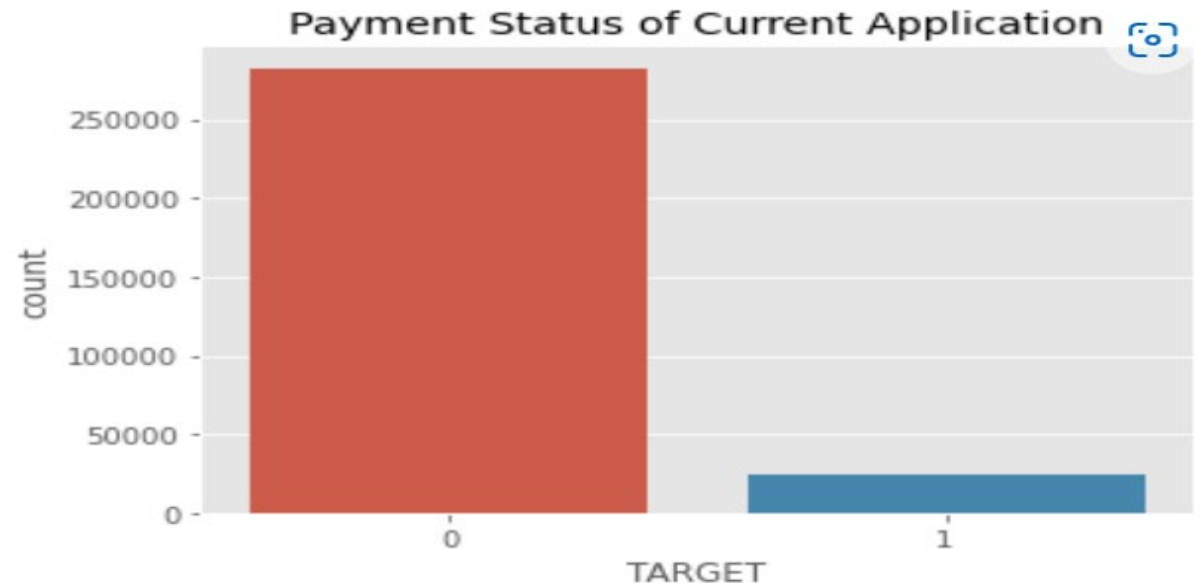
This kind of high correlation value is not accepted in linear model and this will cause multi collinearity and it will like increase number of feature we don't need this many feature so we safely drop it.



Imbalance Data



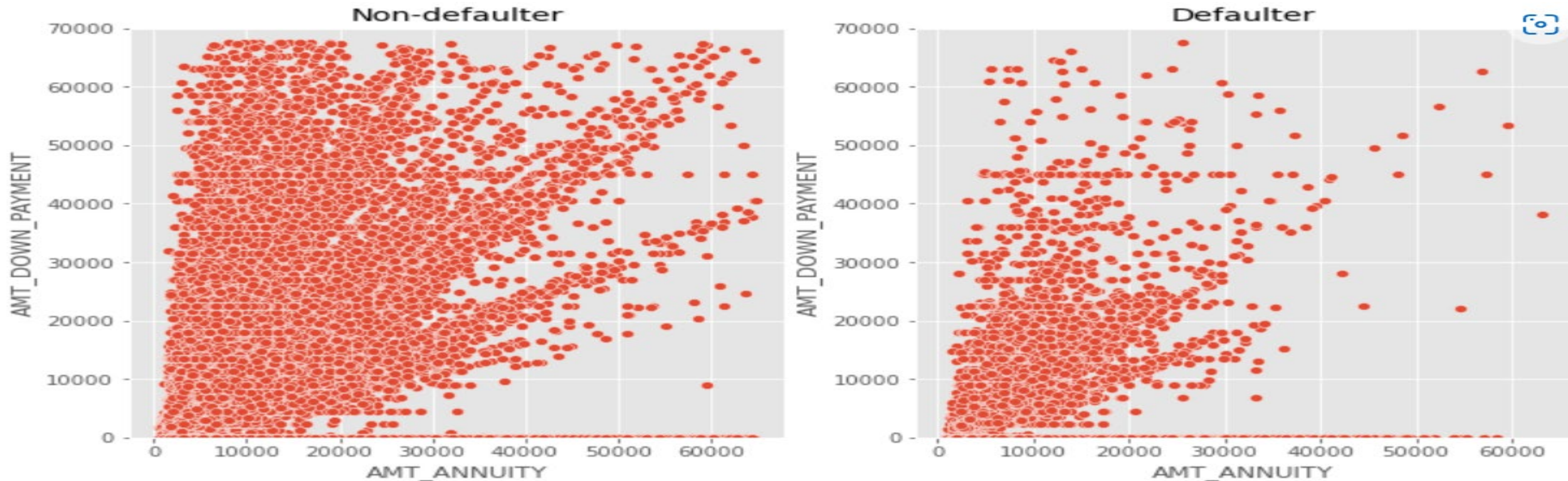
- The data is very highly imbalanced.
- Majority of loans approved some are refused.



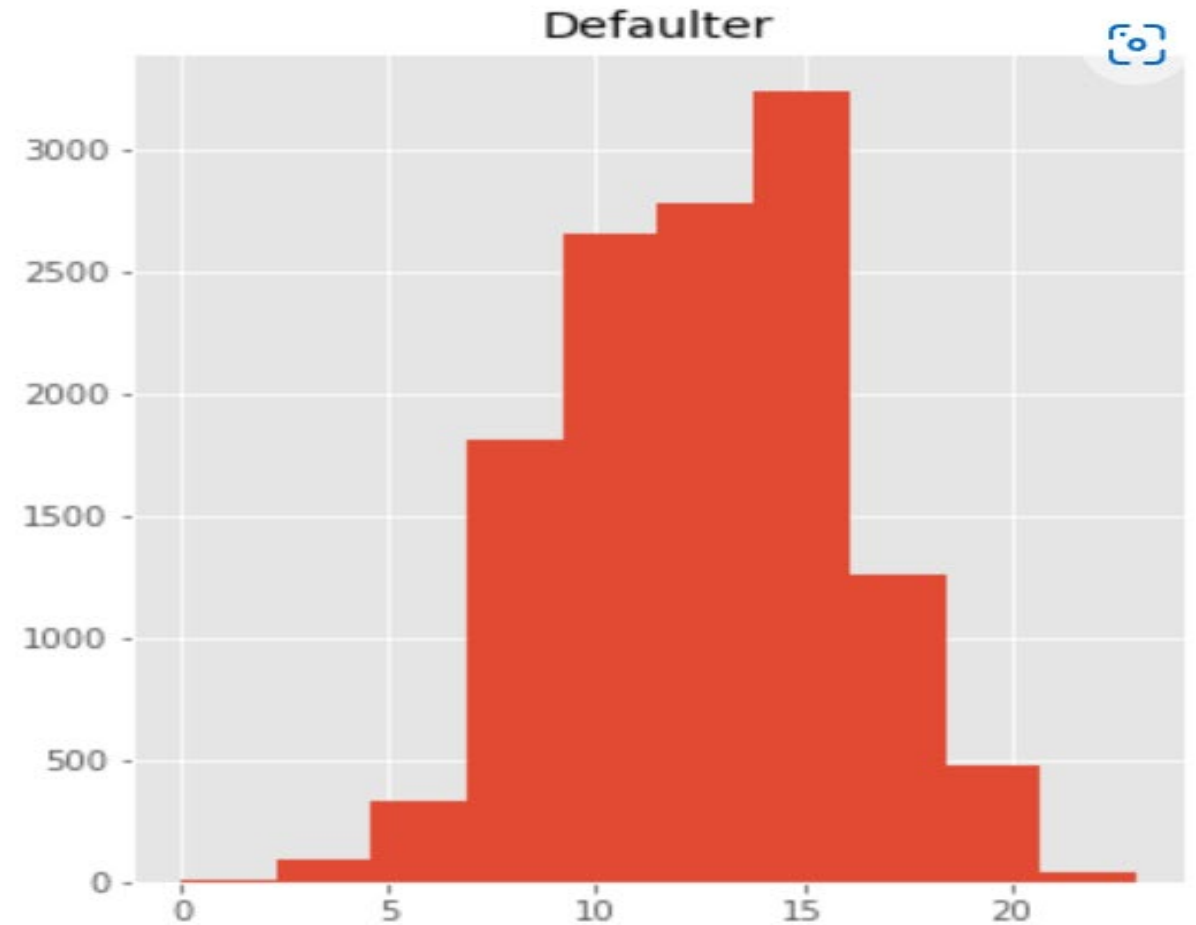
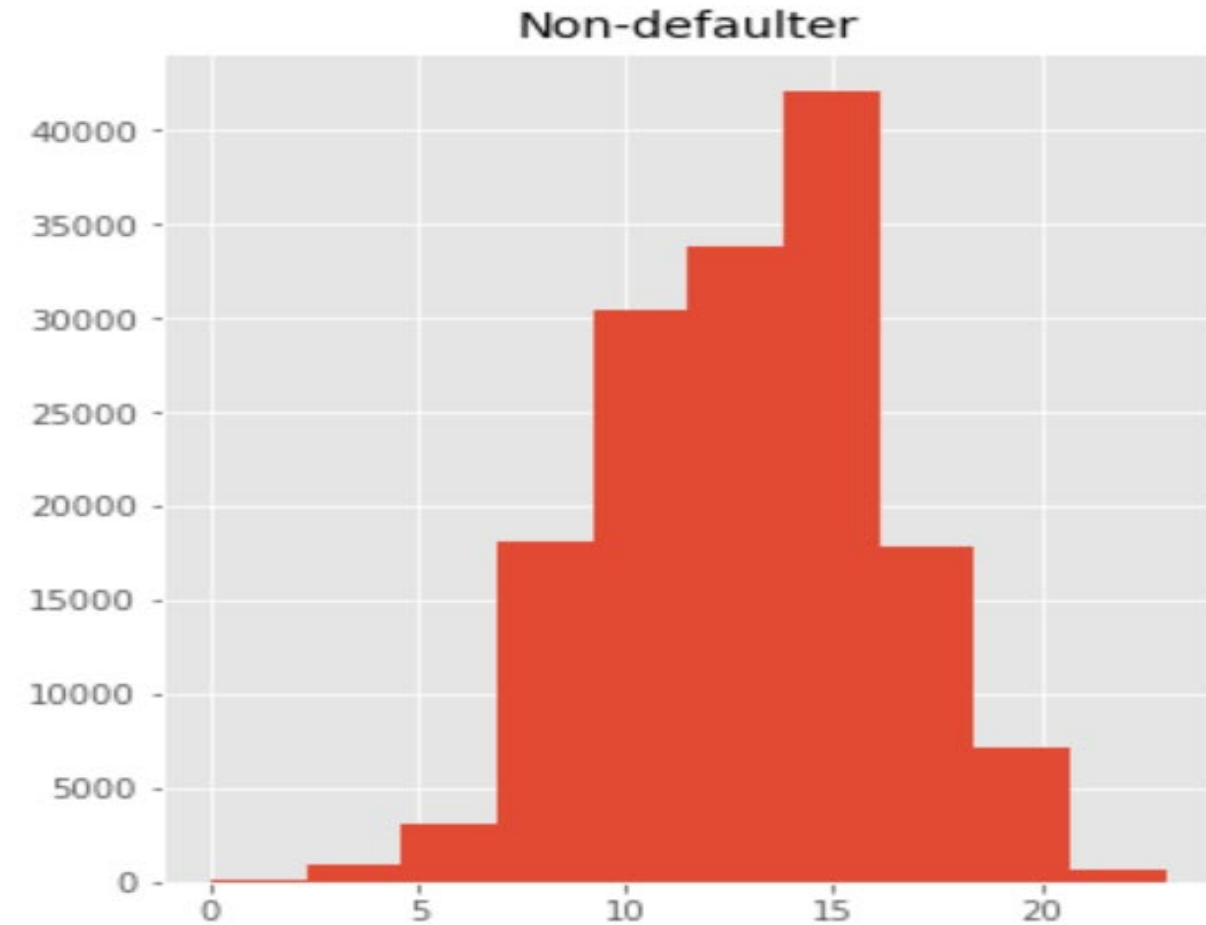
- As well as data is very highly imbalanced.
- Number of defaulter is very less in total application.

Numeric Feature of Previous Application Data

- NUMBER OF DEFAULTERS ARE LESS FOR LARGER AMOUNT OF ANNUITY OF PREVIOUS APPLICATION
- FOR HIGHER DOWN PAYMENT DEFAULTER CASES ARE LESS

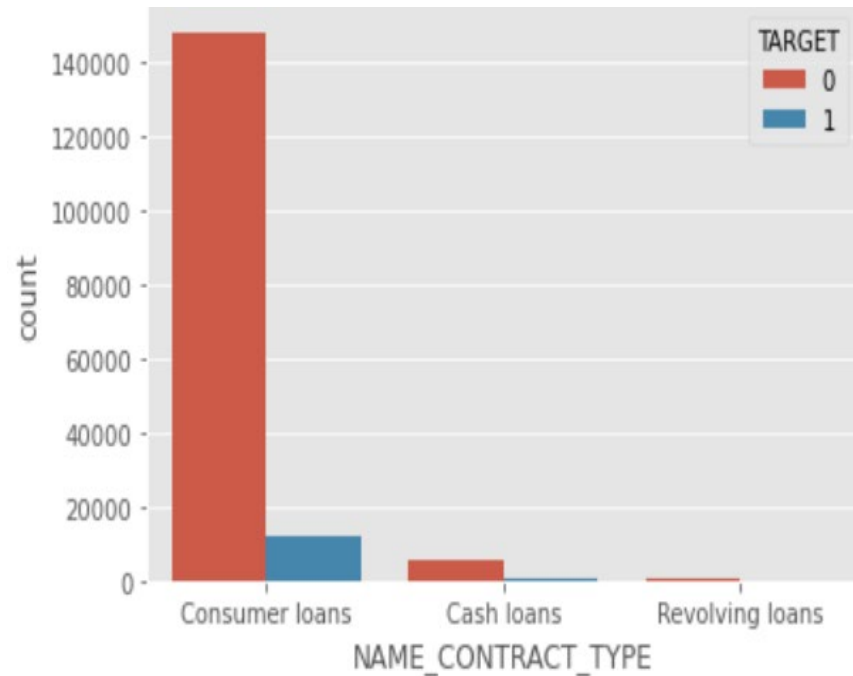


Visualization



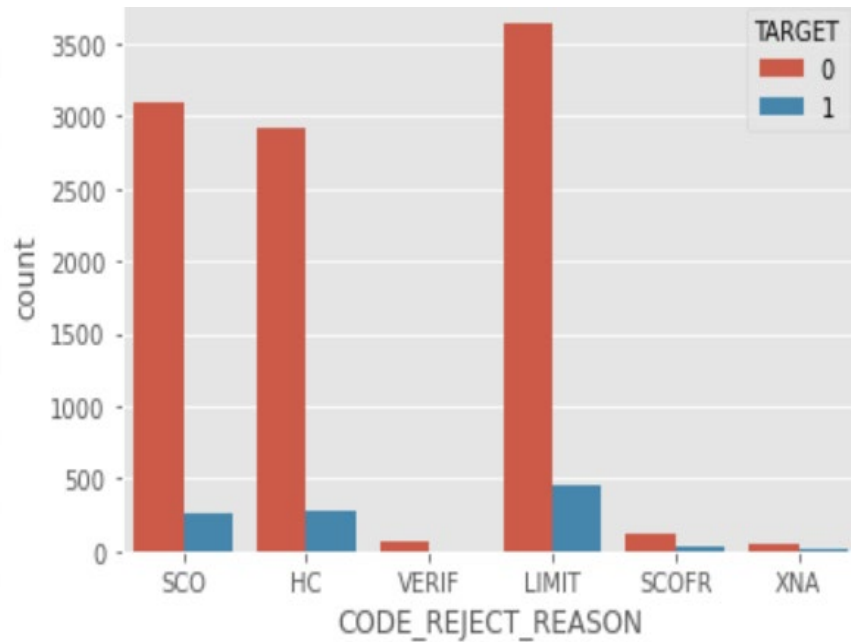
- Most of the loans are applied around 15:00 hours . This feature is does not have visible impact on TARGET variable

Applied, Rejection, Approval



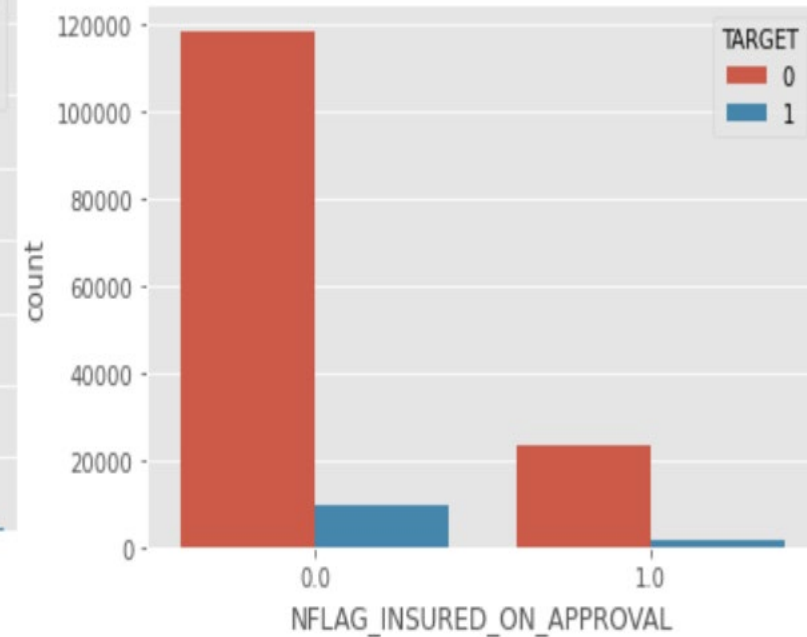
→ We can understand what kind of loans applied.

→ Highest number of loans are applied for Consumer loans.



→ This is the reason why . loans are rejected

→ As seen in the above plot . 'SCO', 'LIMIT', "HC are most common reason of rejection



→ Most of the people did not required insurance during previous loan application.

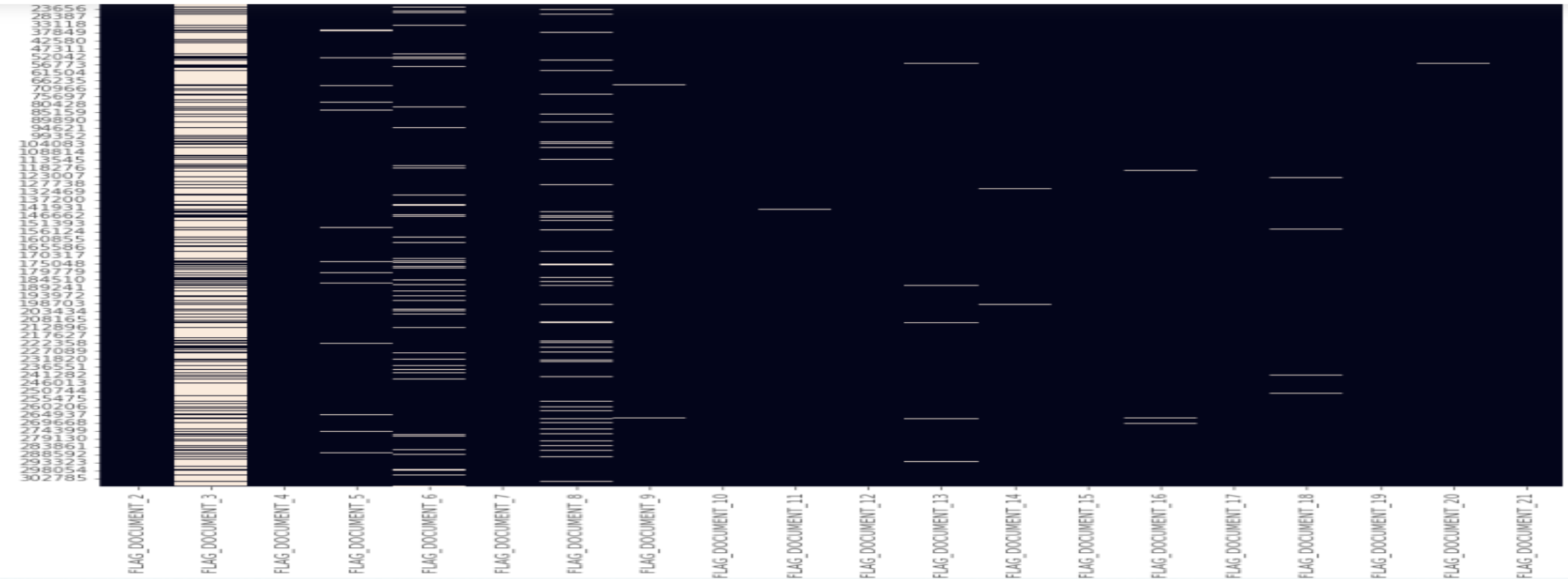
→ People are just braved they don't want insurance.

Visualization



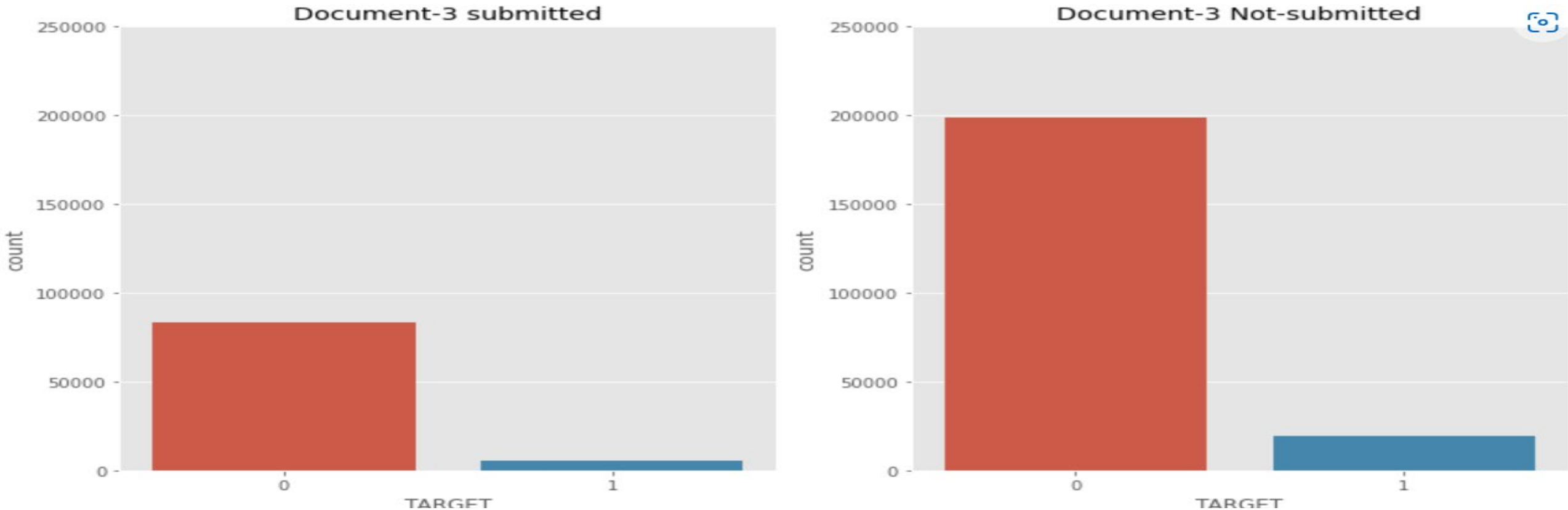
1. Most of the clients are repeaters
2. "Cash through the bank" is the most frequently use payment method

Documents Submitted By Clients



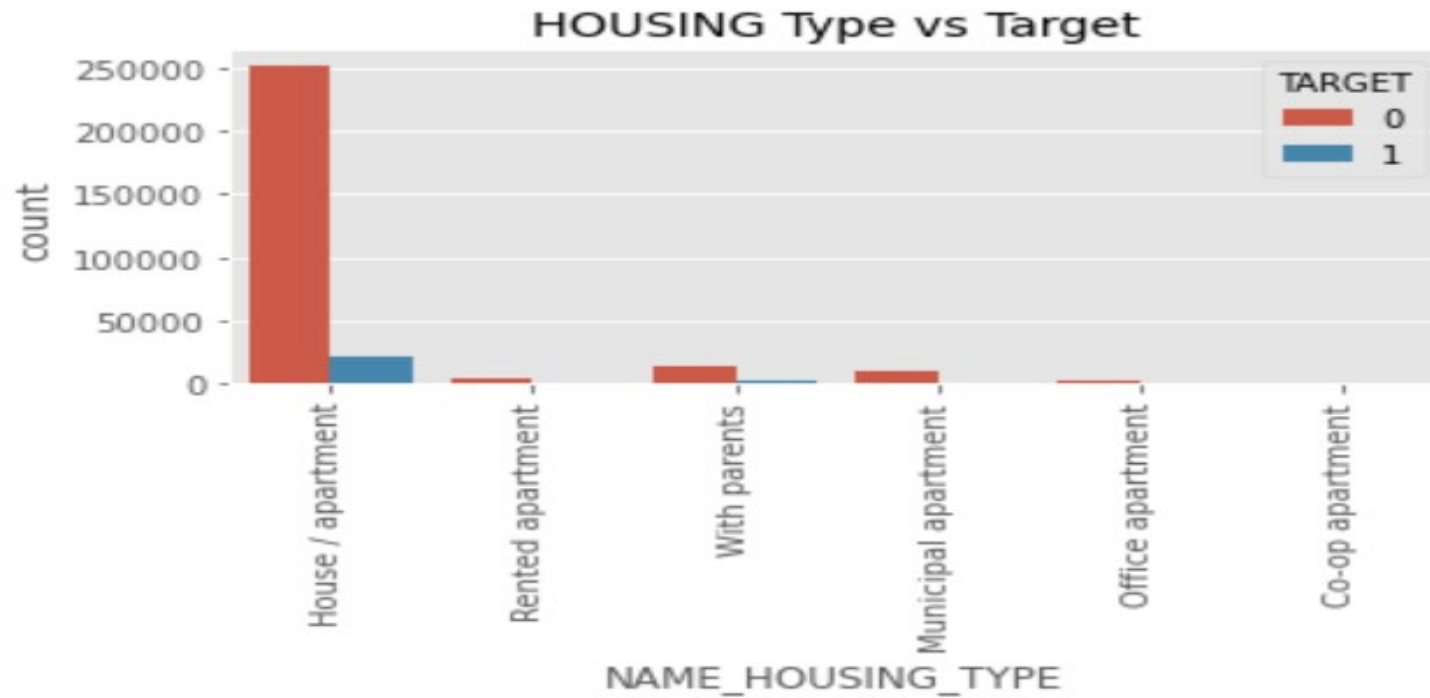
1. The heatmap suggests that all of the documents except Document 3 were not provided by clients in majority of the cases.
2. Hence we can assume all the document(except document 3) will not contribute towards analyzing the data. Hence columns can be dropped

Default and Non-Default



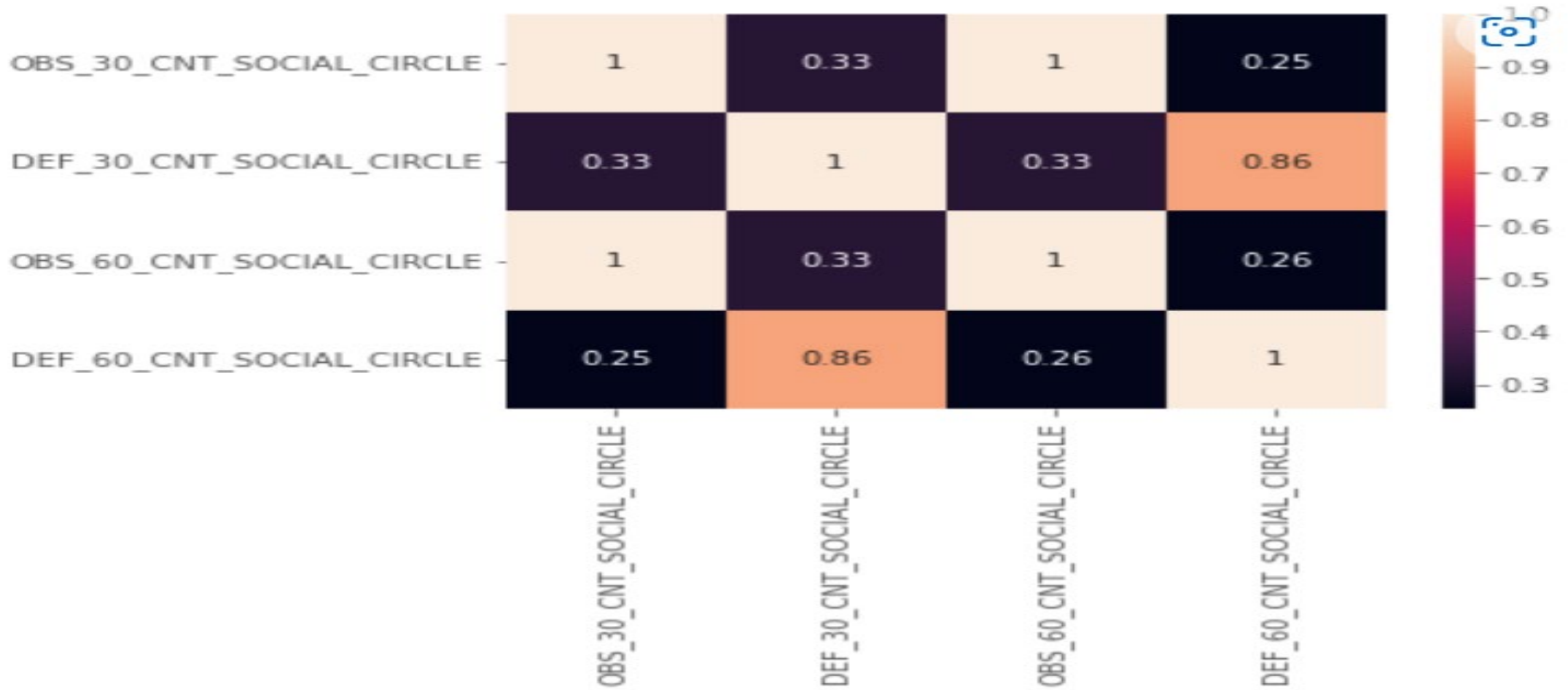
- FLAG_DOCUMENT_3 is showing similar trend for both non_defaulter and defaulter.
- Hence this column can be dropped.

Housing Information of Clients



- Most of the clients live in House/Apartment
- Clients living with their parents or in rented apartment have higher rate of default.

Social Circle Info



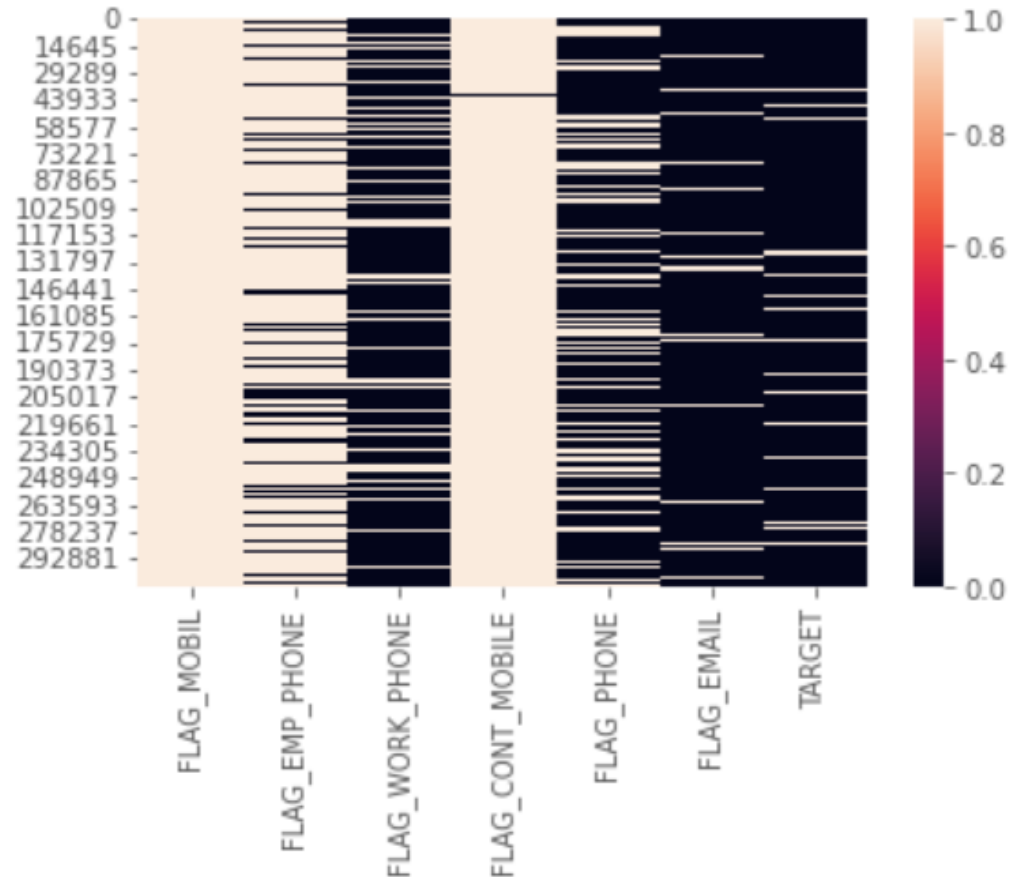
- DEF_30_CNT_SOCIAL_CIRCLE and DEF_60_CNT_SOCIAL_CIRCLE are highly correlated
- OBS_30_CNT_SOCIAL_CIRCLE and OBS_60_CNT_SOCIAL_CIRCLE are identical columns

Region Related Data



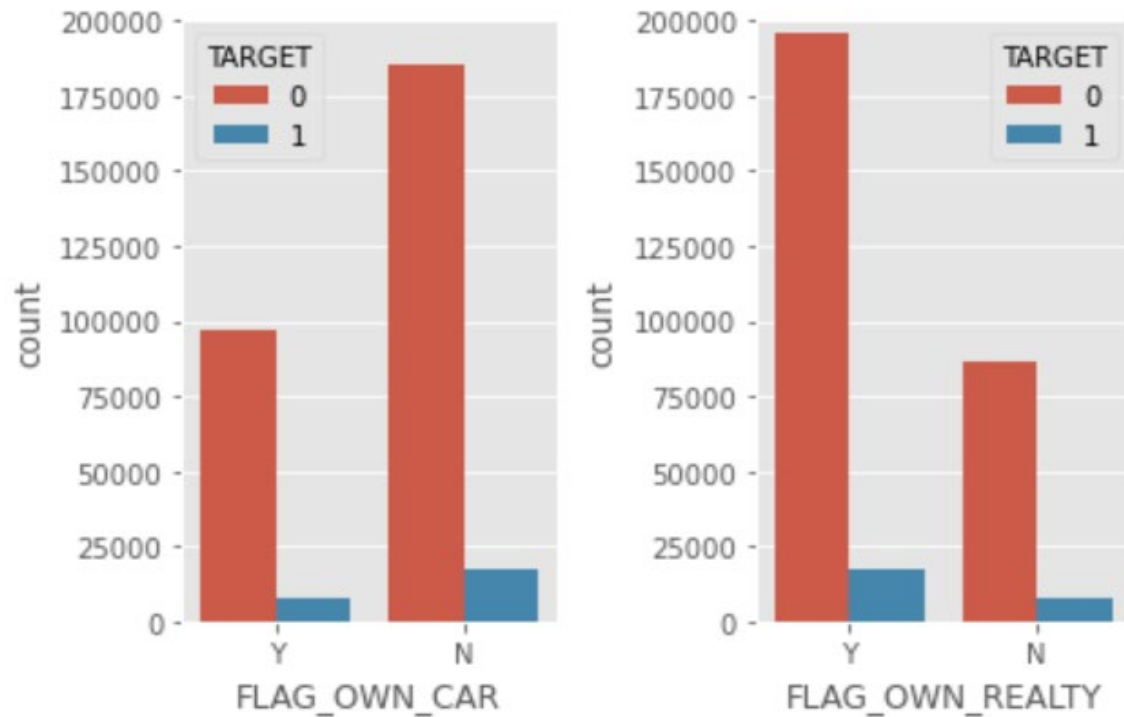
- Defaulter rate is highest when REG_REGION_NOT_WORK_REGION=0 i.e permanent address and working address is same
- Highest Application have Region rating of 2

Contact Related Info



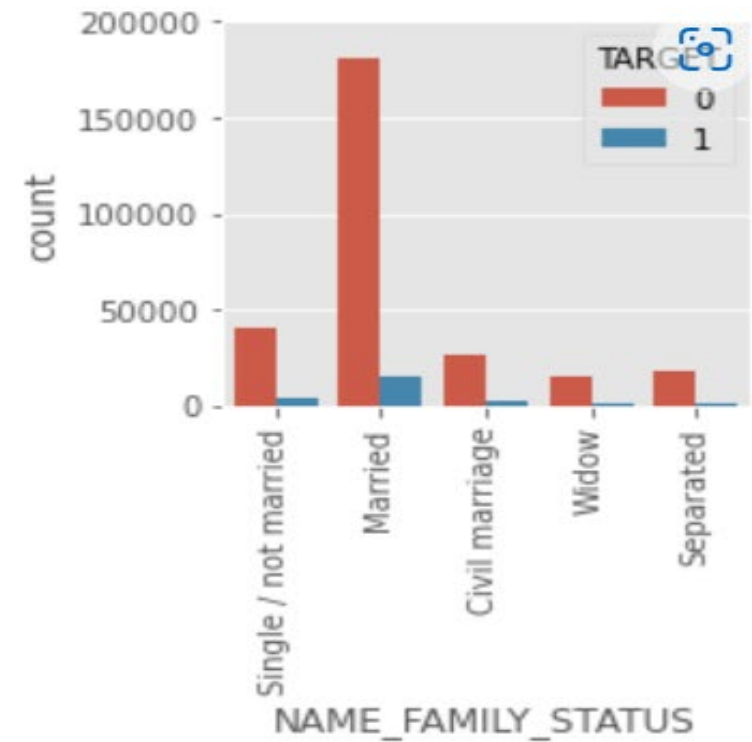
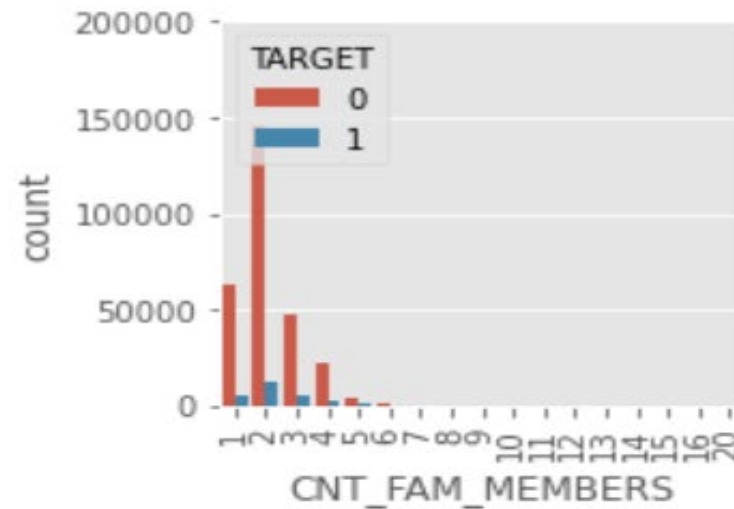
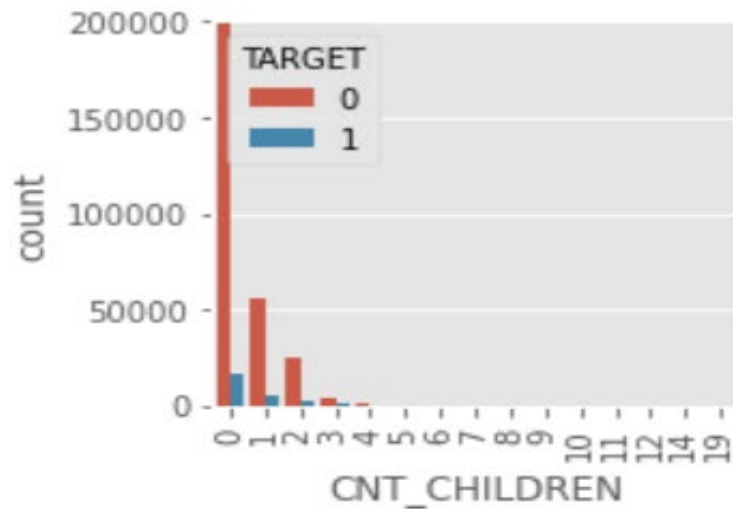
- All the features in contact_df are categorical (0 and 1)
- As there no similarity of patterns of TARGET value with the features, we are assuming the feature are not useful for analysis.
- Hence all of the features can be removed

Assets Details



- Most of the clients own realty
- Most of the clients do not own car
- People not owning realty and car and have a slightly higher default tare than the people who own realty and car

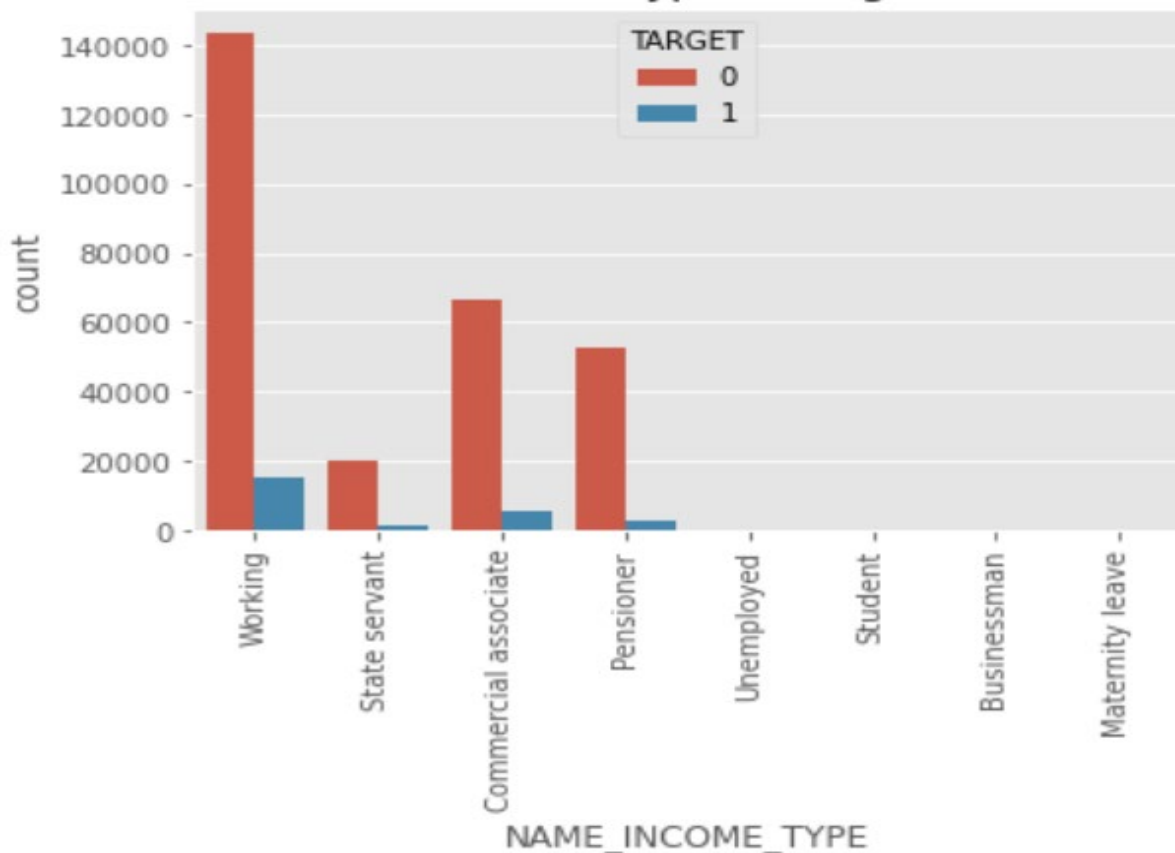
Family Related Info



- Default rate is highest for Civil Manager and Single clients
- Most of the clients are married (and/or) no children (and/or) 2 family members
- Clients with relatively more number of children (and/or) family members have higher default percentage
- For some of the cases where count children/family members high, and the default rate very high or very low. This cases cannot be taken as a conclusion as number of clients having a large family is very low.

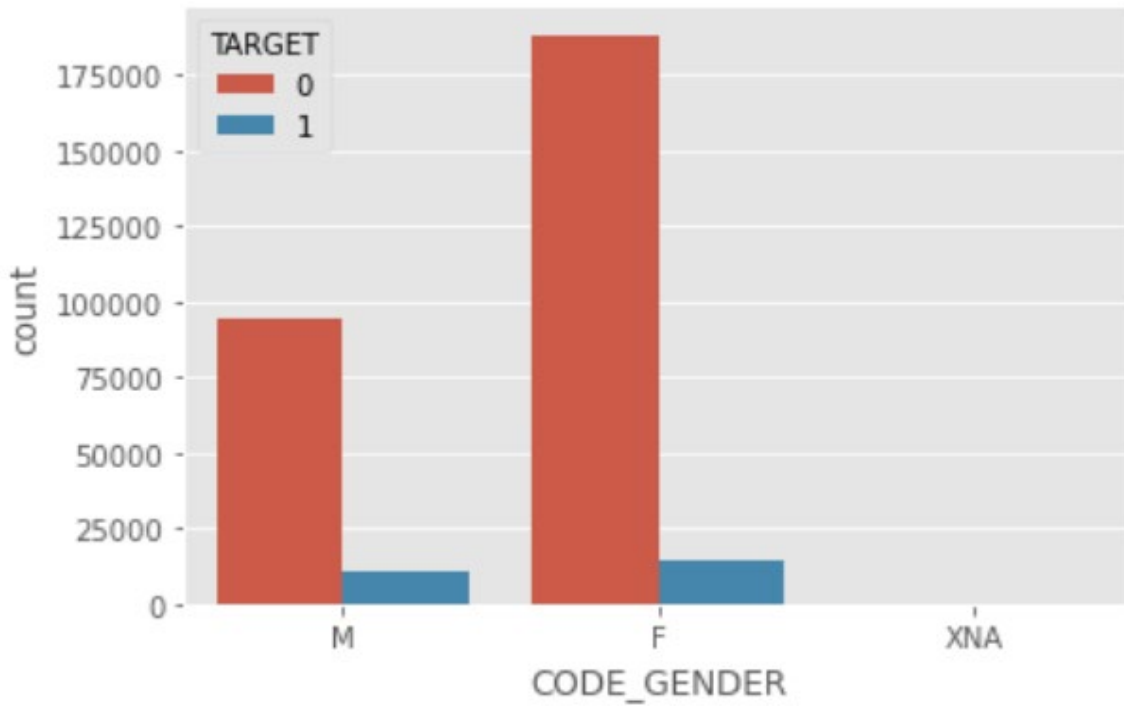
Education and Occupation Info

INCOME Type vs Target



- Most of the clients are working.
- Clients on Maternity Leave and unemployed has highest percentage of Defaulter
- Businessman have lowest(0) percentage of Defaulter However clients of income type("Unemployed", "Student", "Businessman", "Maternity leave") are very few in the dataset to contribute in the analysis.

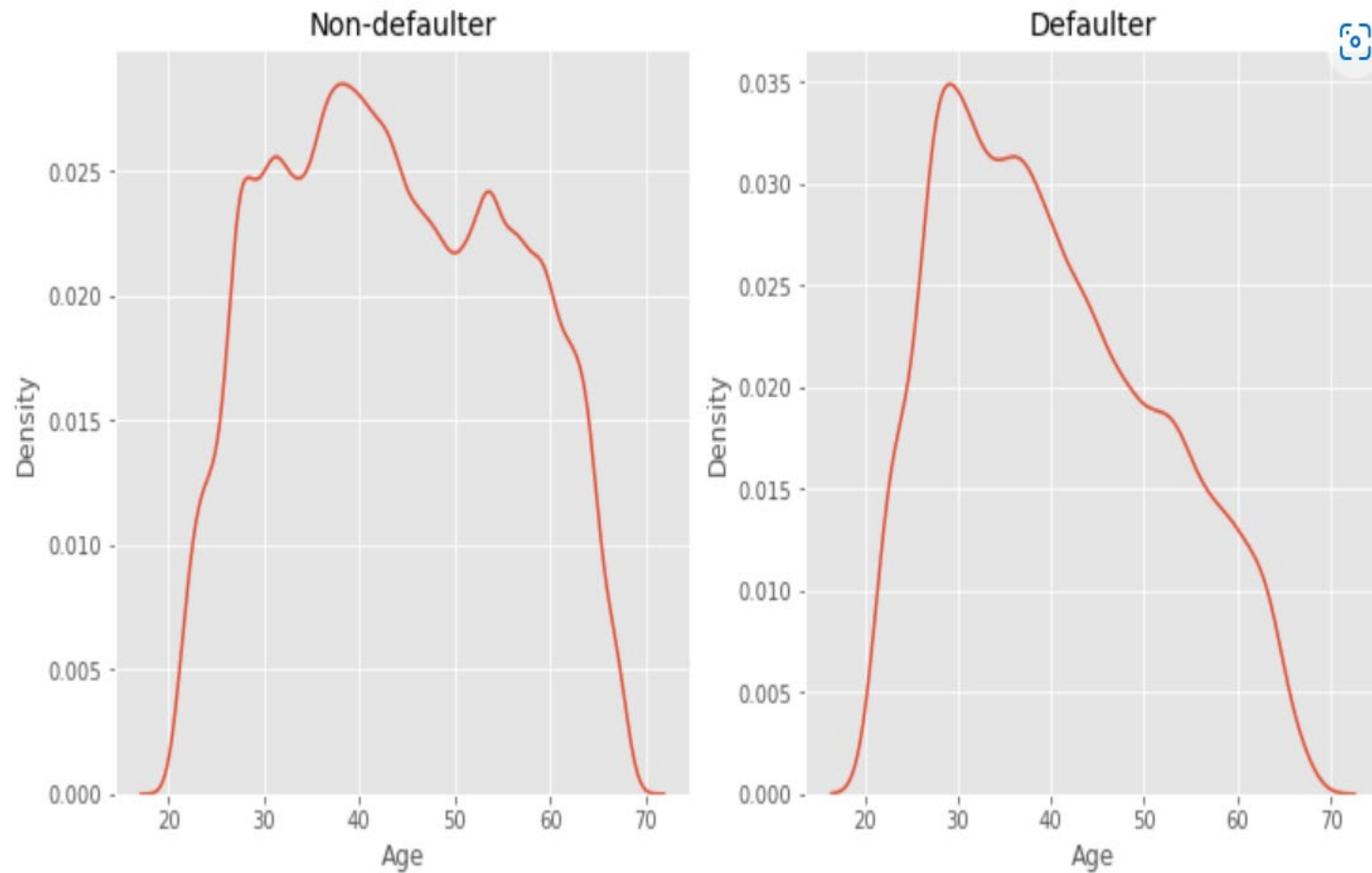
Gender



Female clients are more than male clients

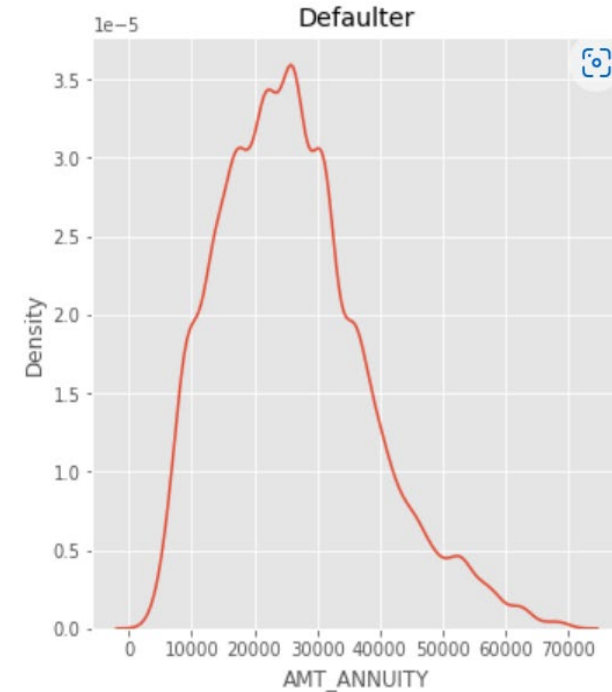
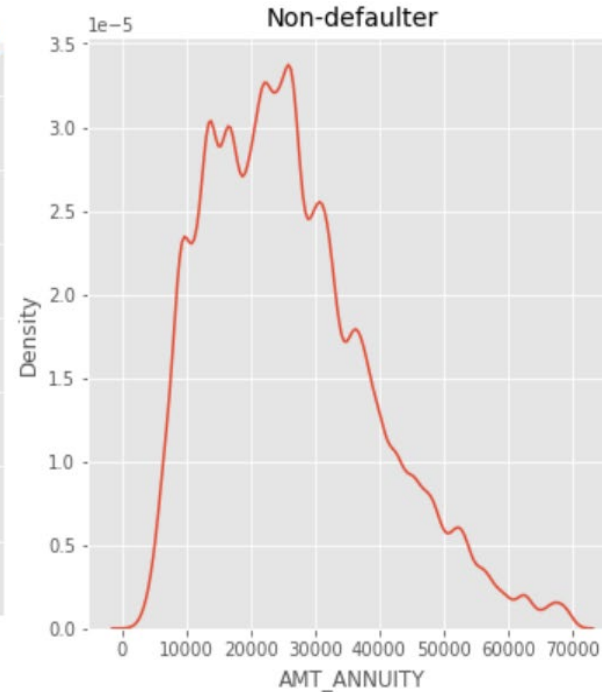
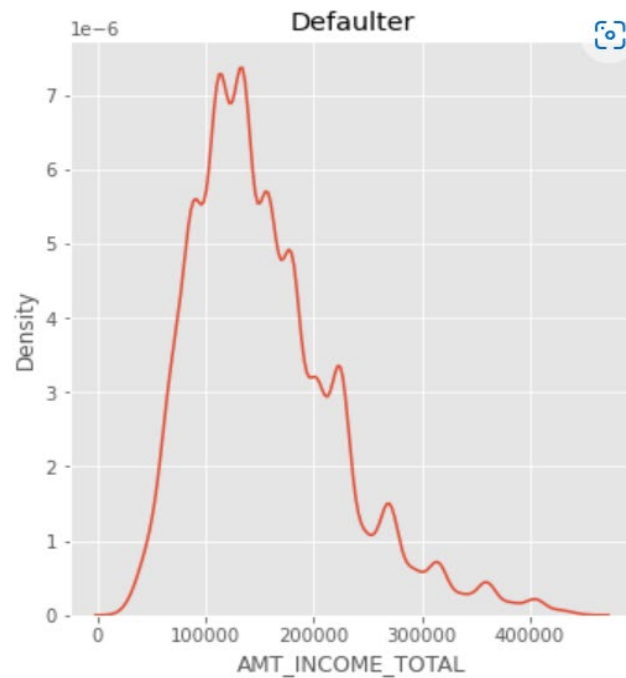
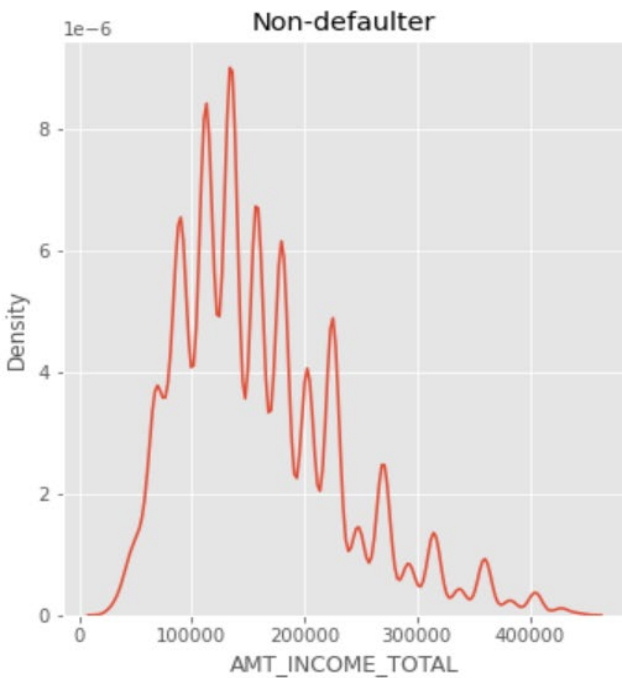
Defaulter percentage for male clients

DAYS_BIRTH



People of age 30 have higher default rate
Default cases are less for clients more than 40 years old.

Income and Annuity



For Defaulter and non-defaulter AMT_ANNUIITY distribution is similar.

Thank you!

