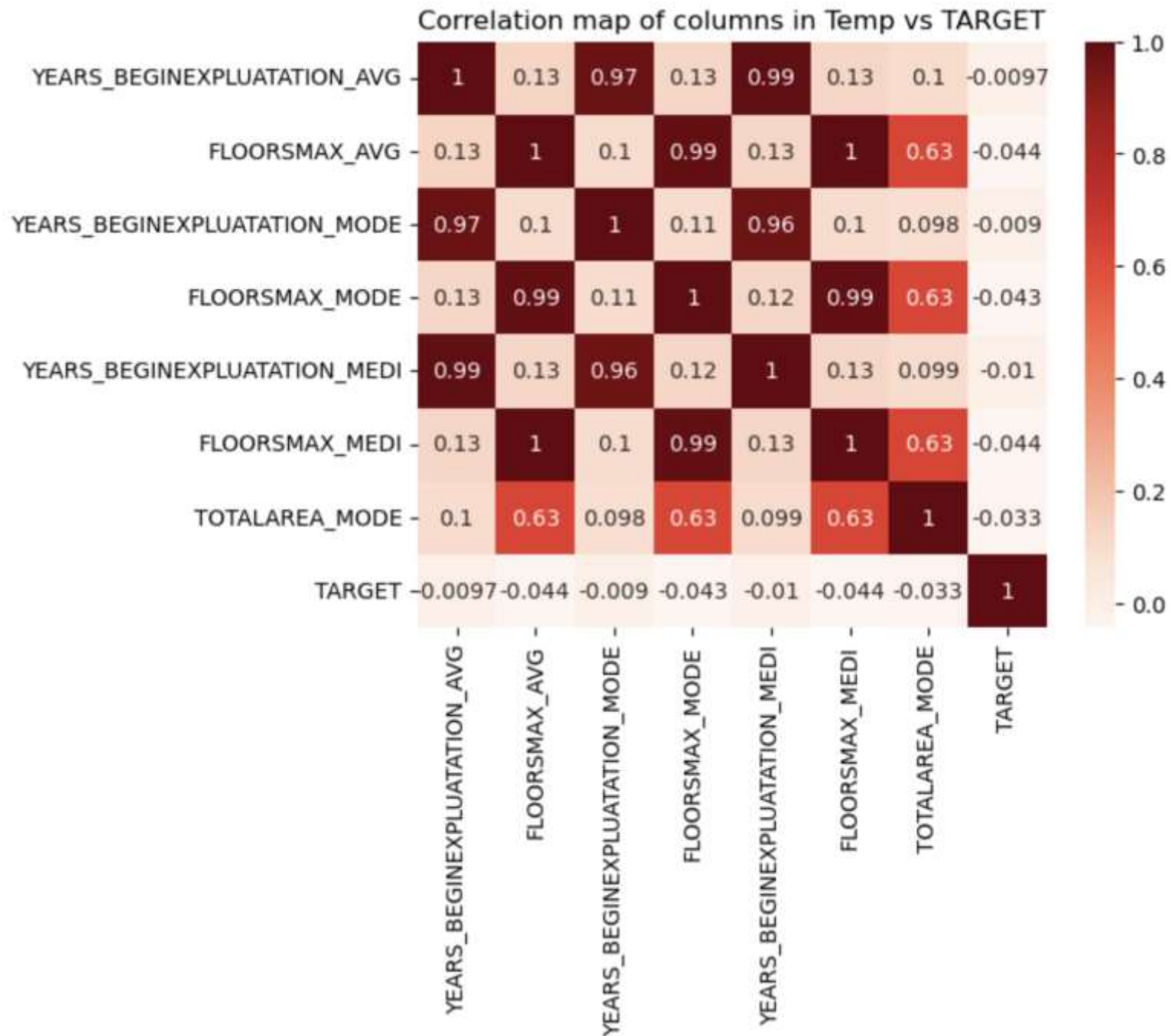


Credit EDA Assignment— Plots Presentation

- by

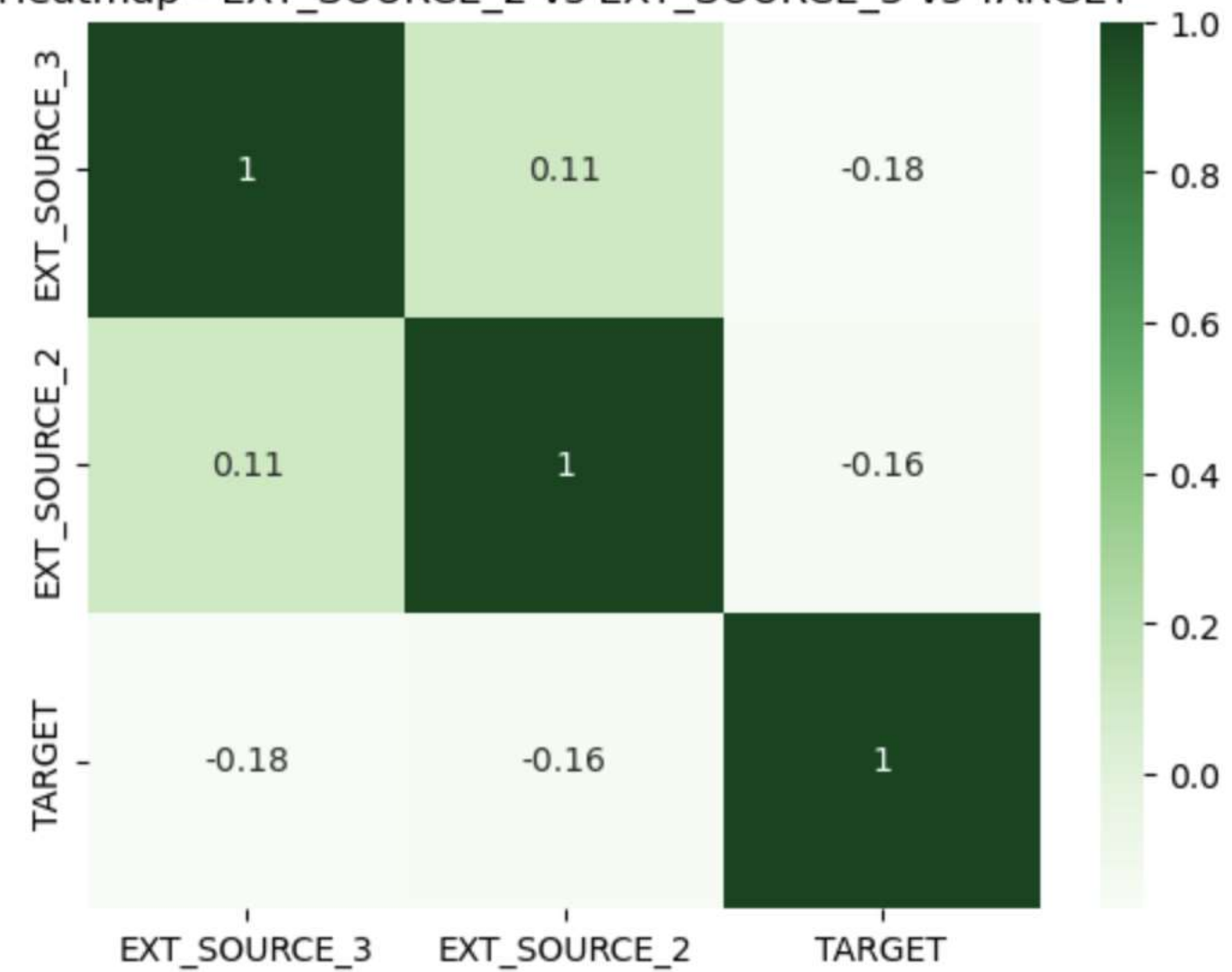
Akella Venkata Koushik

DS C49 – October 2022
batch

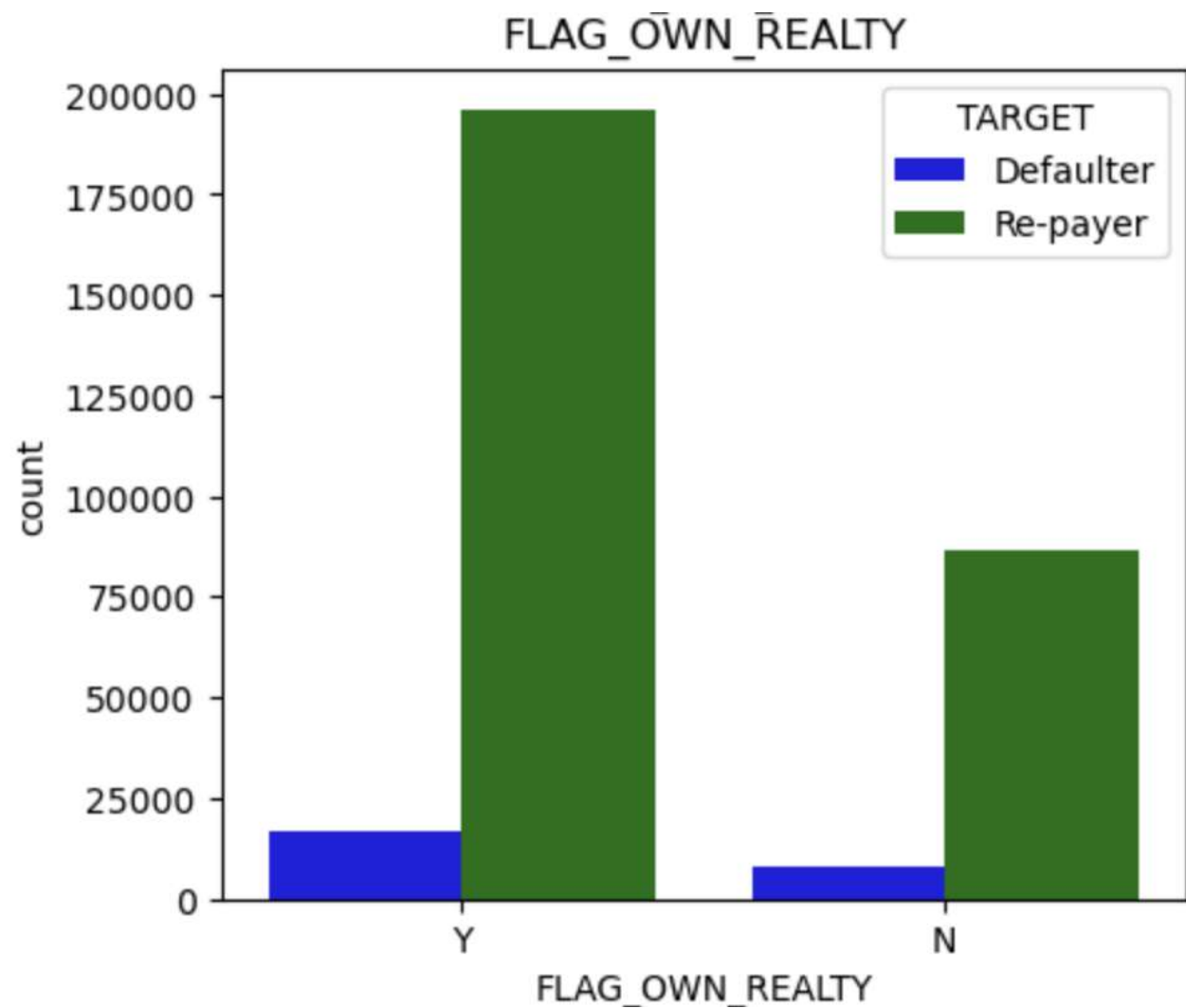
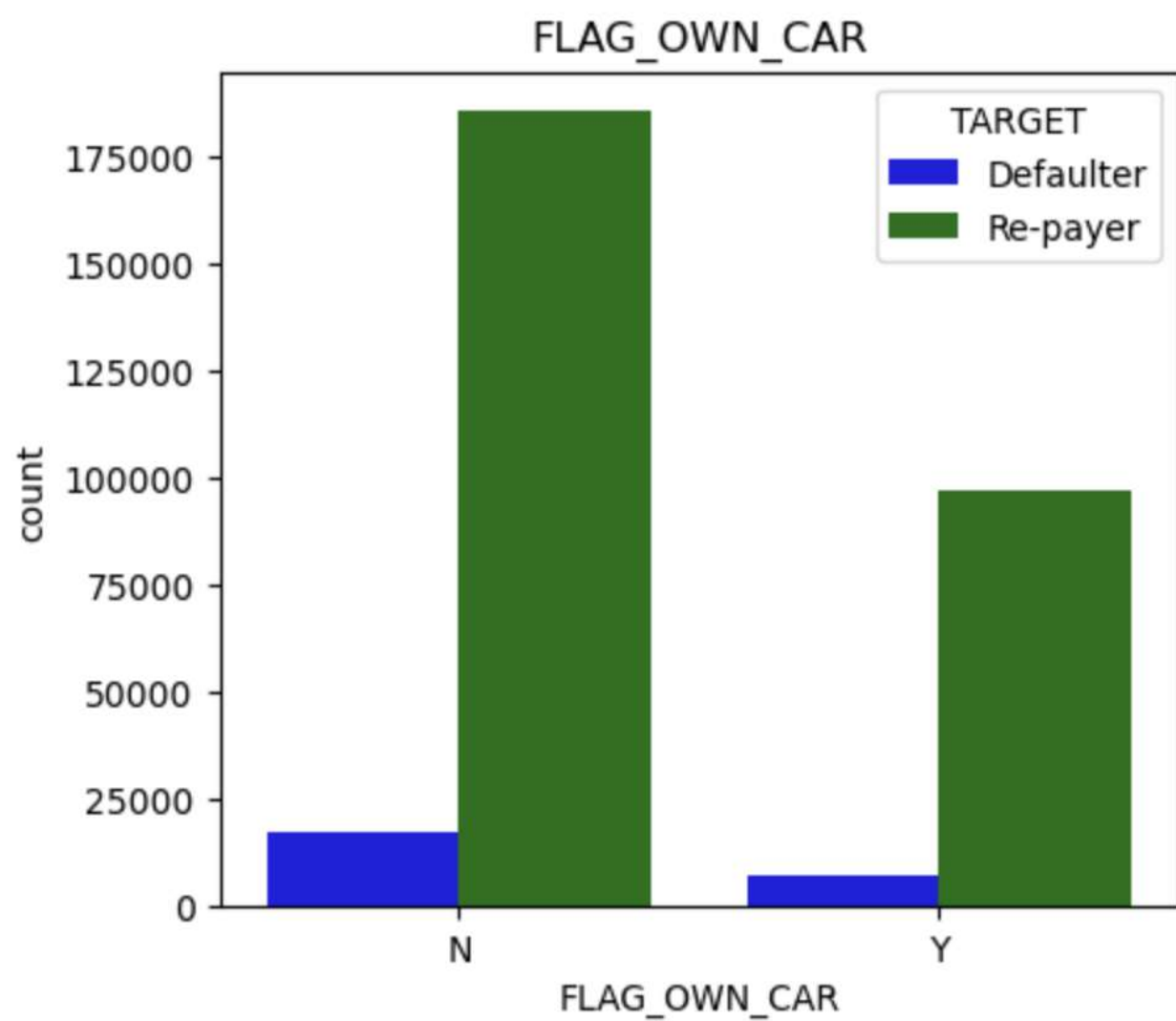


The correlation plot shows that the 8 columns, ['YEARS_BEGINEXPLUATATION_AVG', 'FLOORSMAX_AVG', 'YEARS_BEGINEXPLUATATION_MODE', 'FLOORSMAX_MODE', 'YEARS_BEGINEXPLUATATION_MEDI', 'FLOORSMAX_MEDI', 'TOTALAREA_MODE', 'EMERGENCYSTATE_MODE'] have a non-linear relationship with the Target variable, and hence, even if we impute the missing values in these columns, they will not affect the TARGET value. hence these can be dropped comfortably.

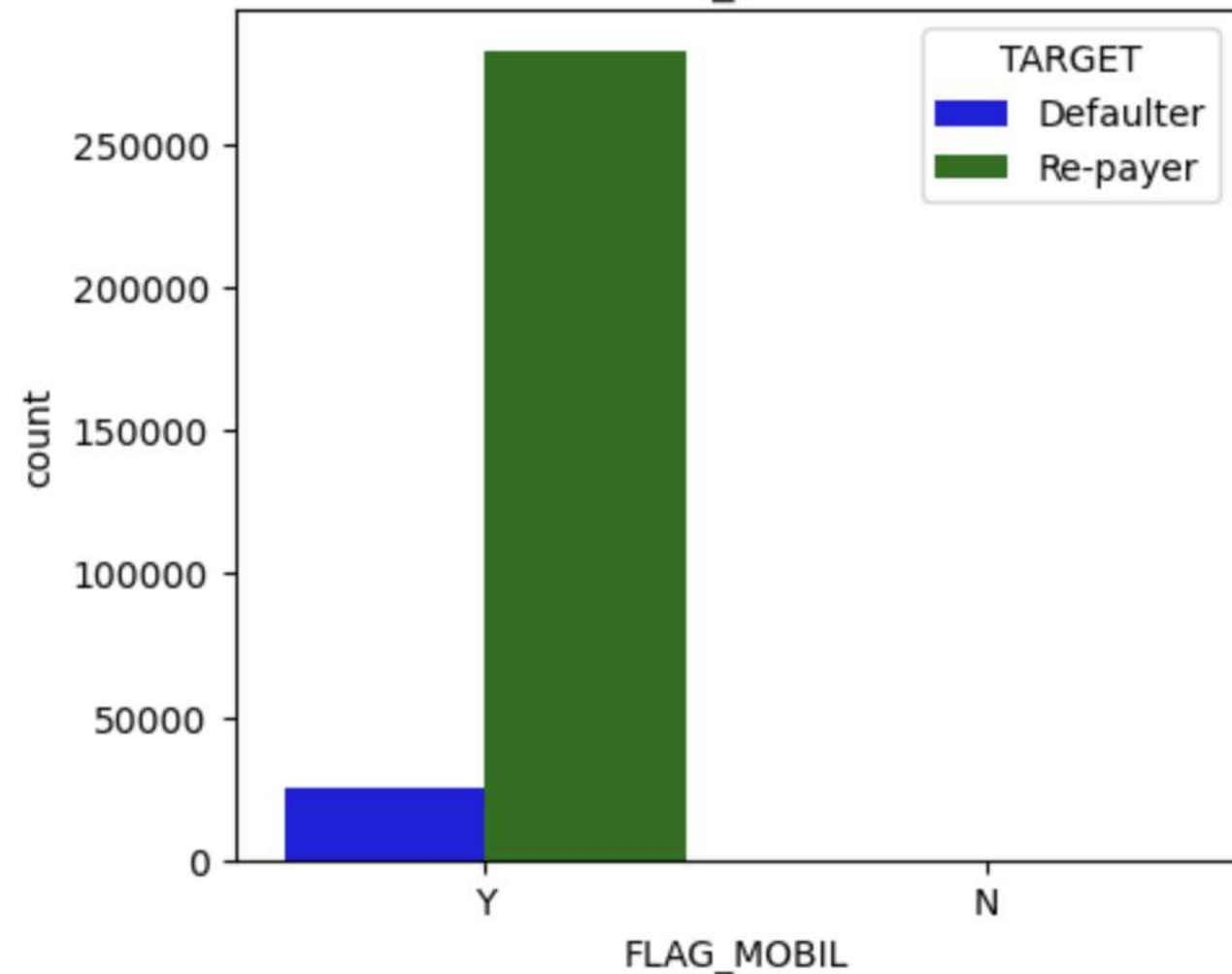
Heatmap - EXT_SOURCE_2 vs EXT_SOURCE_3 vs TARGET



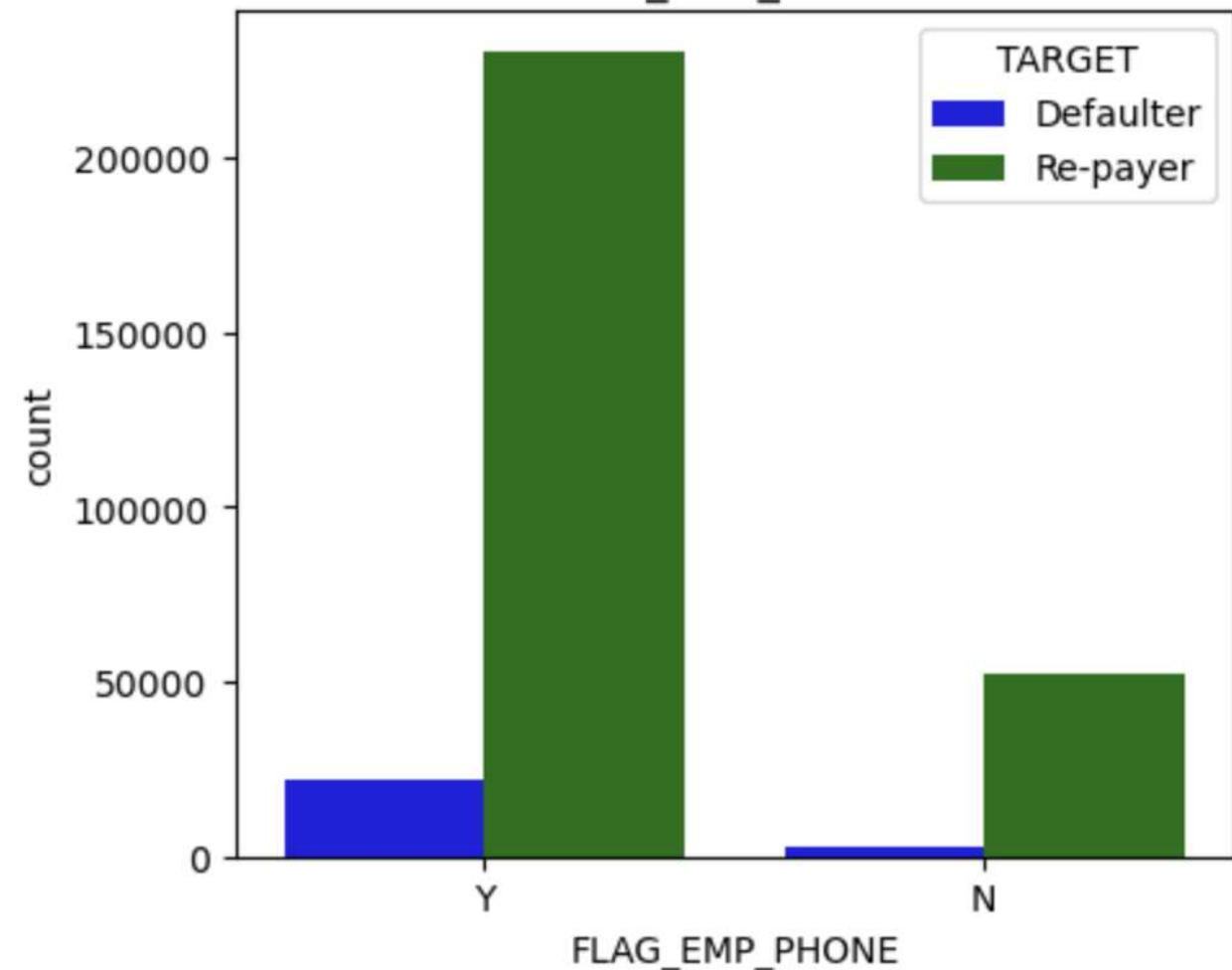
these two columns, 'EXT_SOURCE_3', 'EXT_SOURCE_2' also clearly have a non-linear relation with TARGET variable and hence I feel these have little to no effect on the TARGET variable. and hence, I am dropping these columns



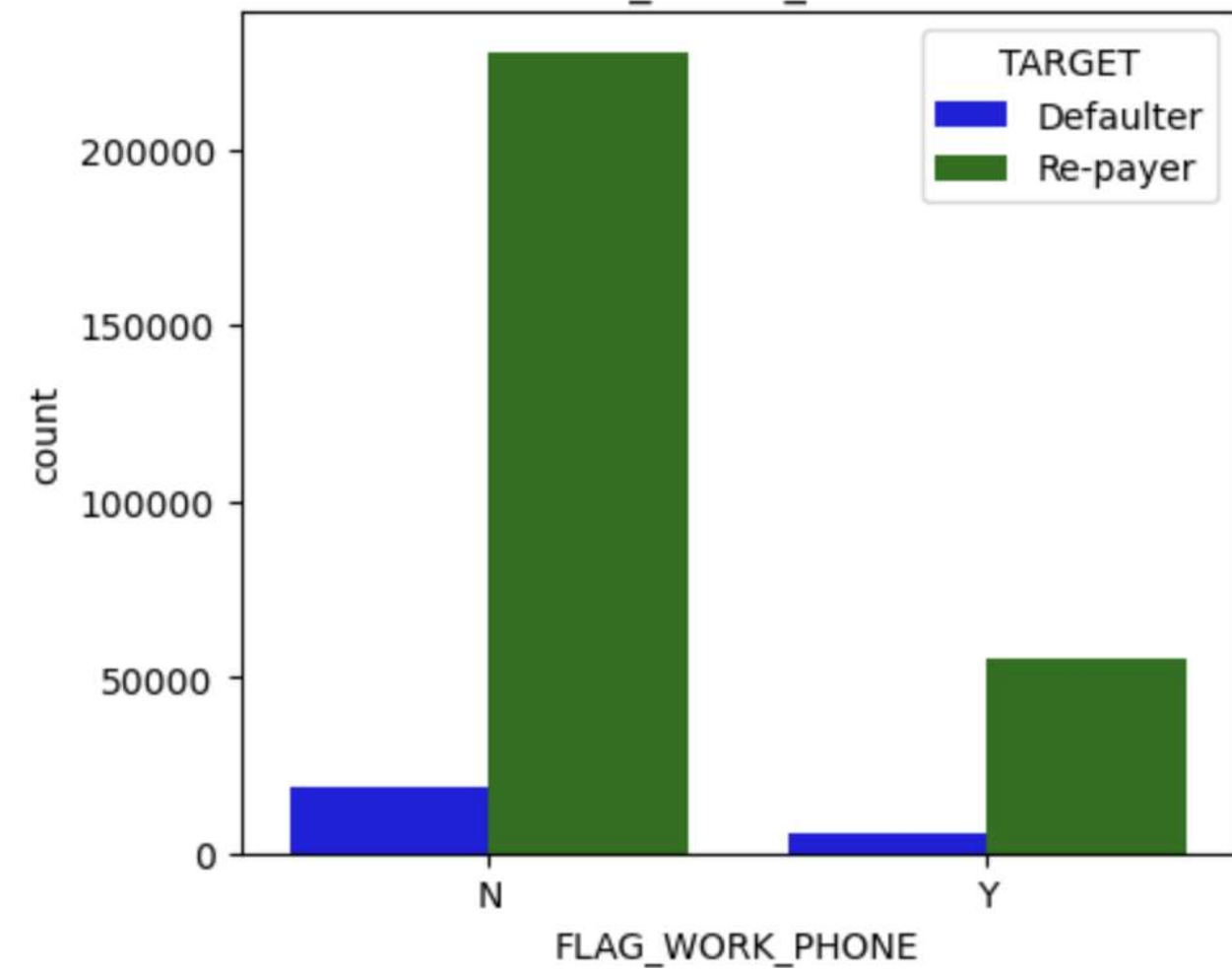
FLAG_MOBIL



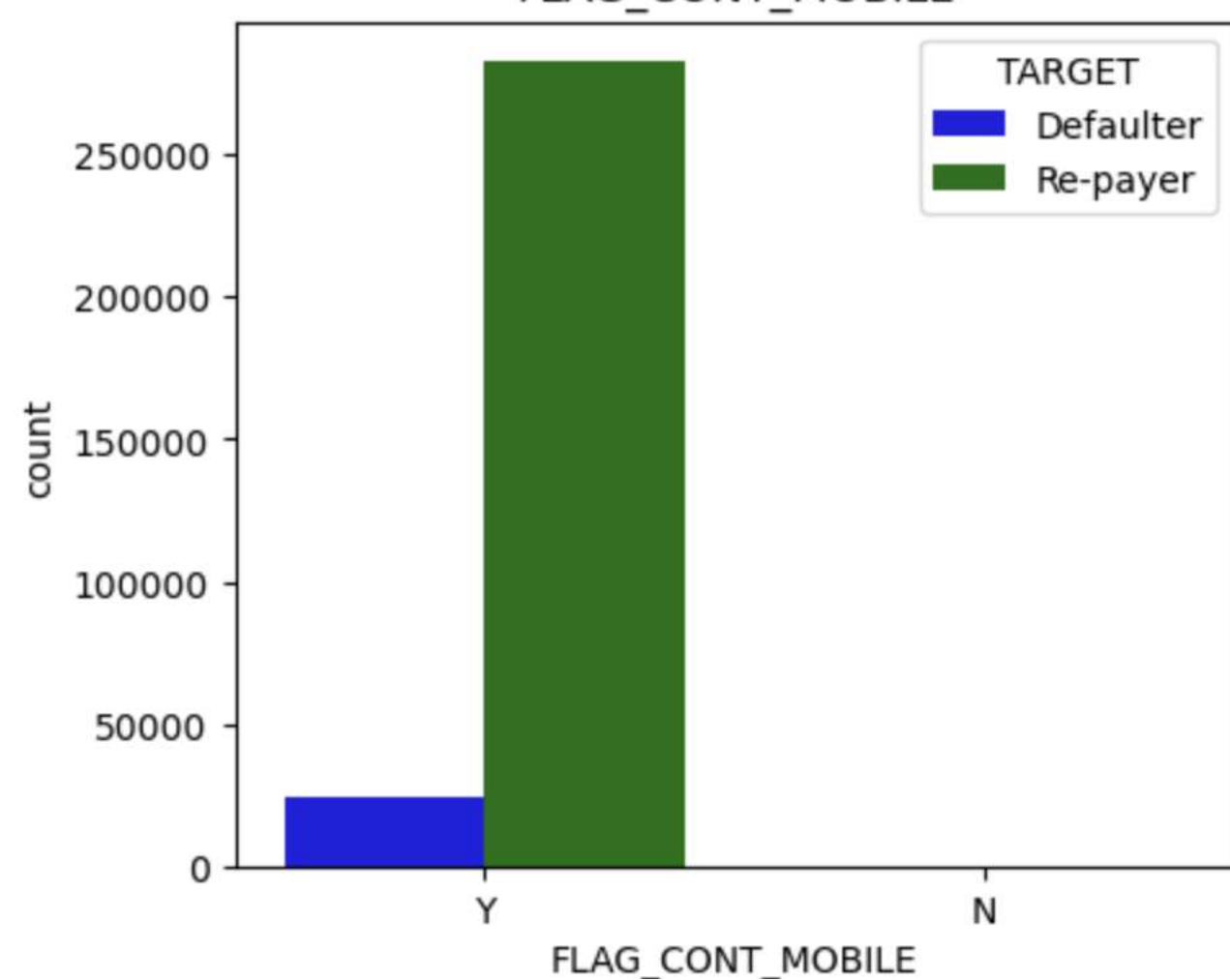
FLAG_EMP_PHONE

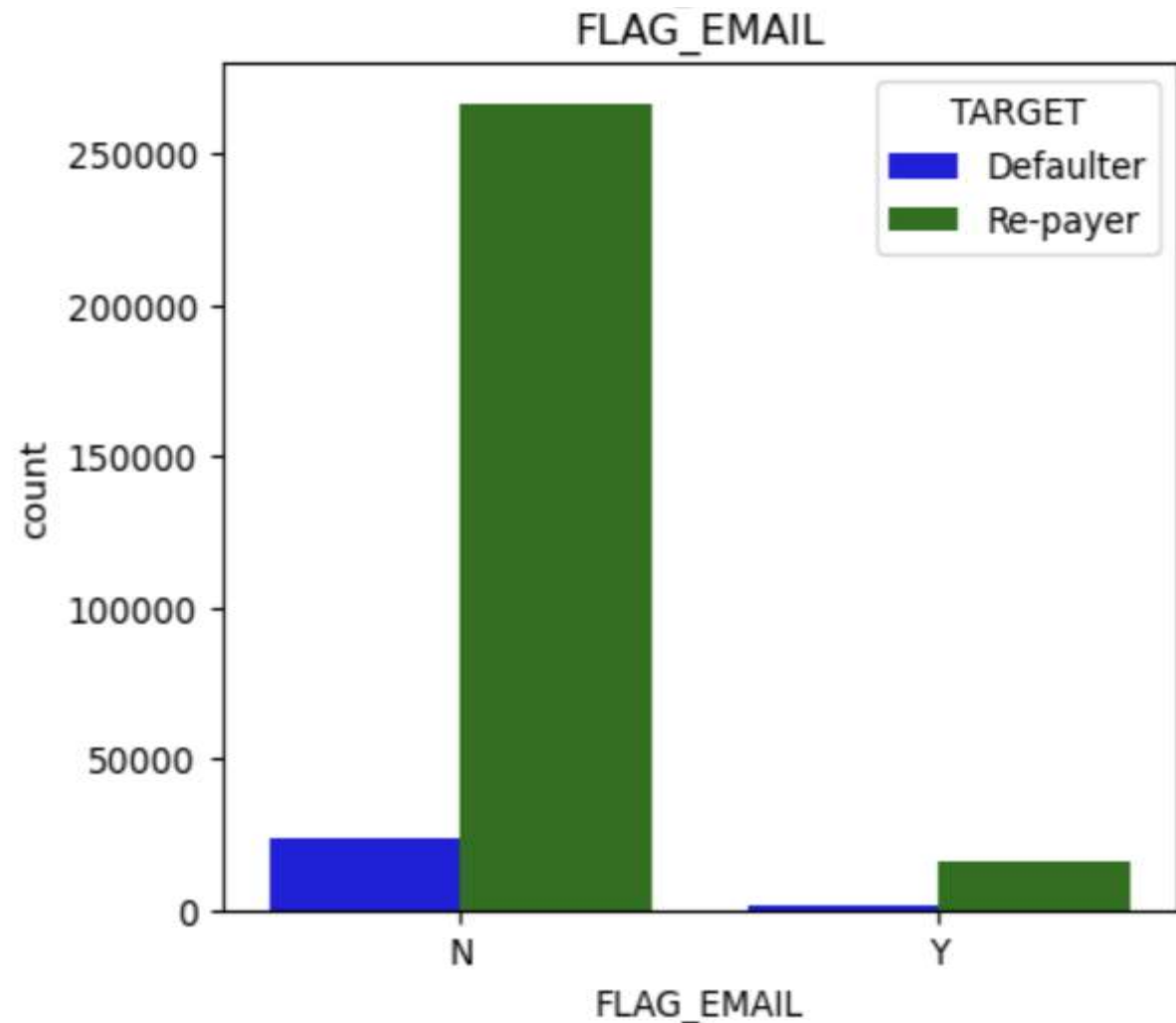
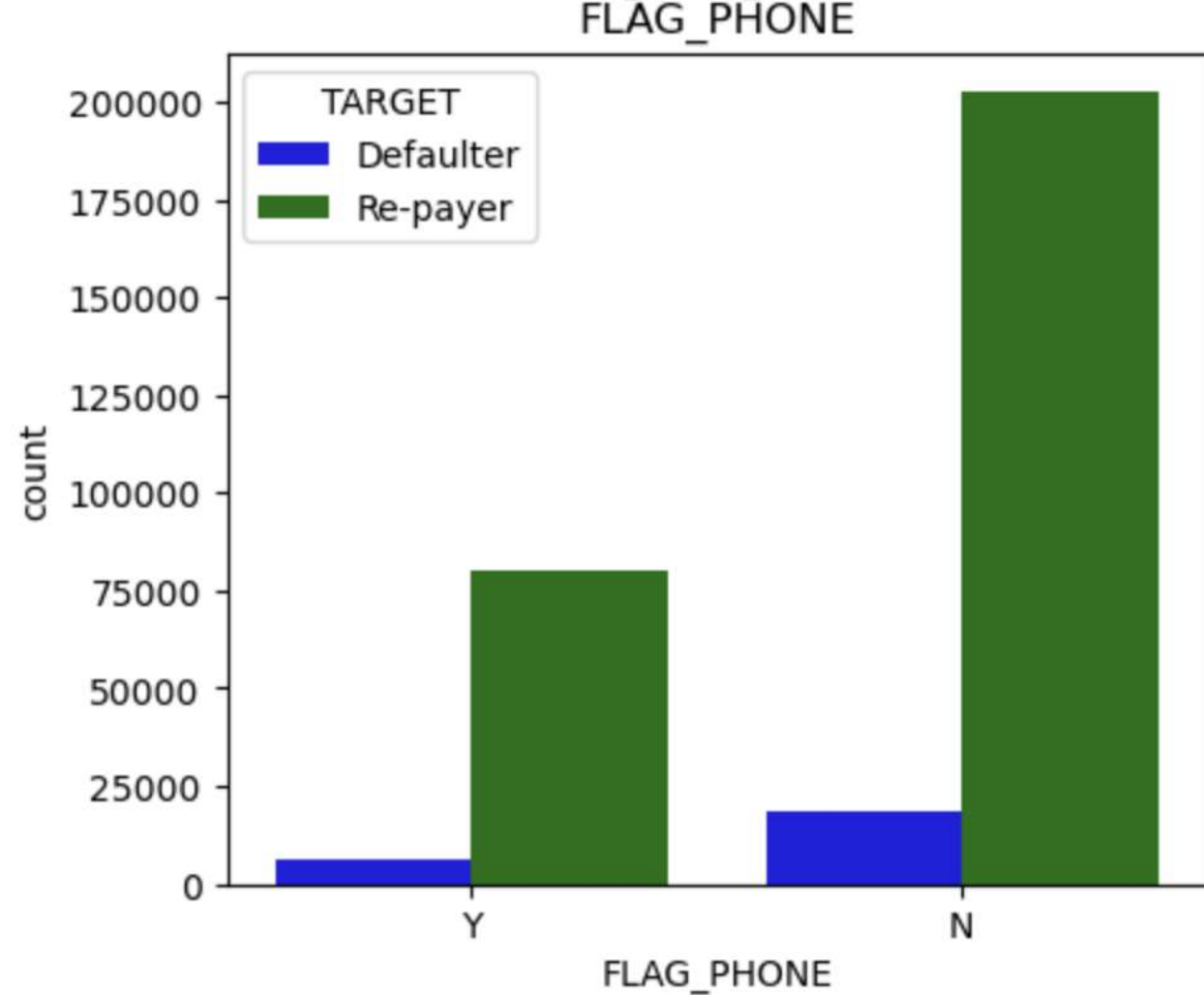


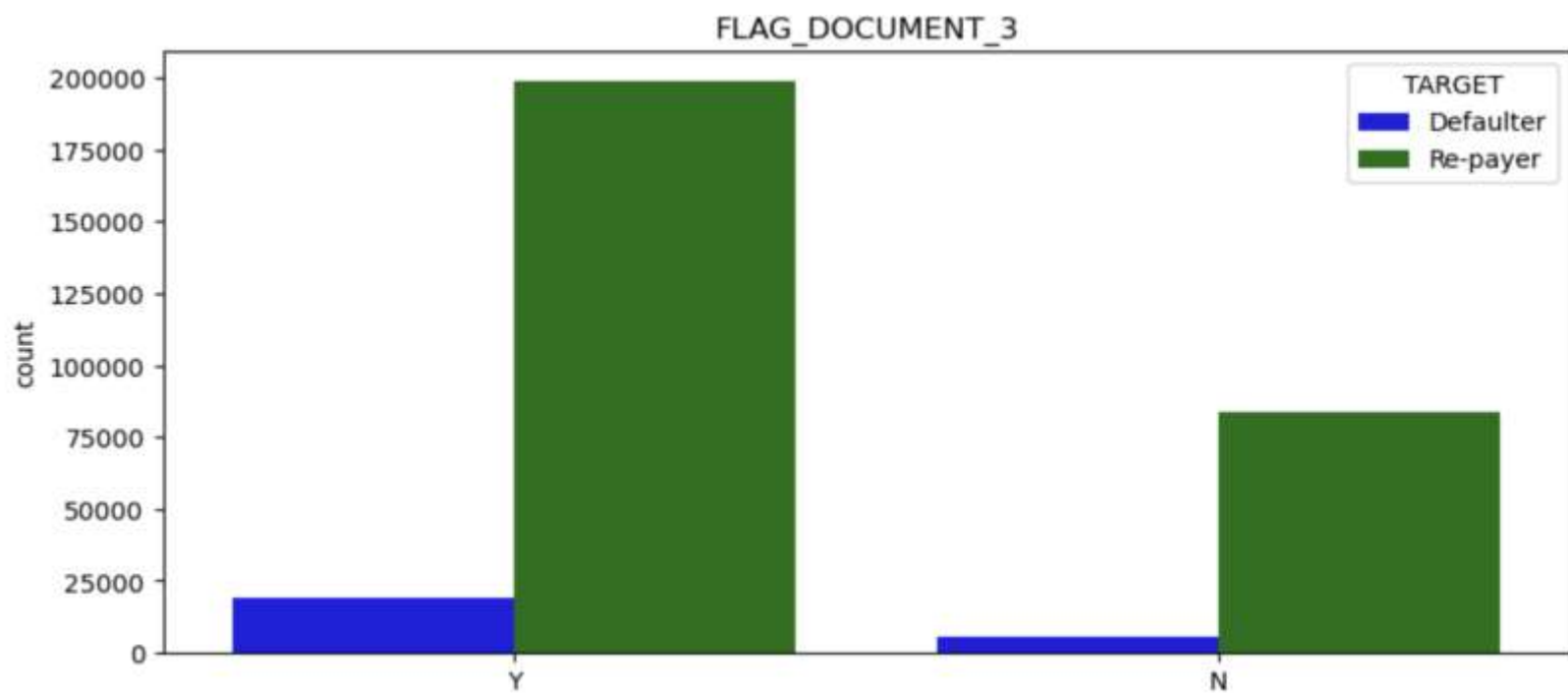
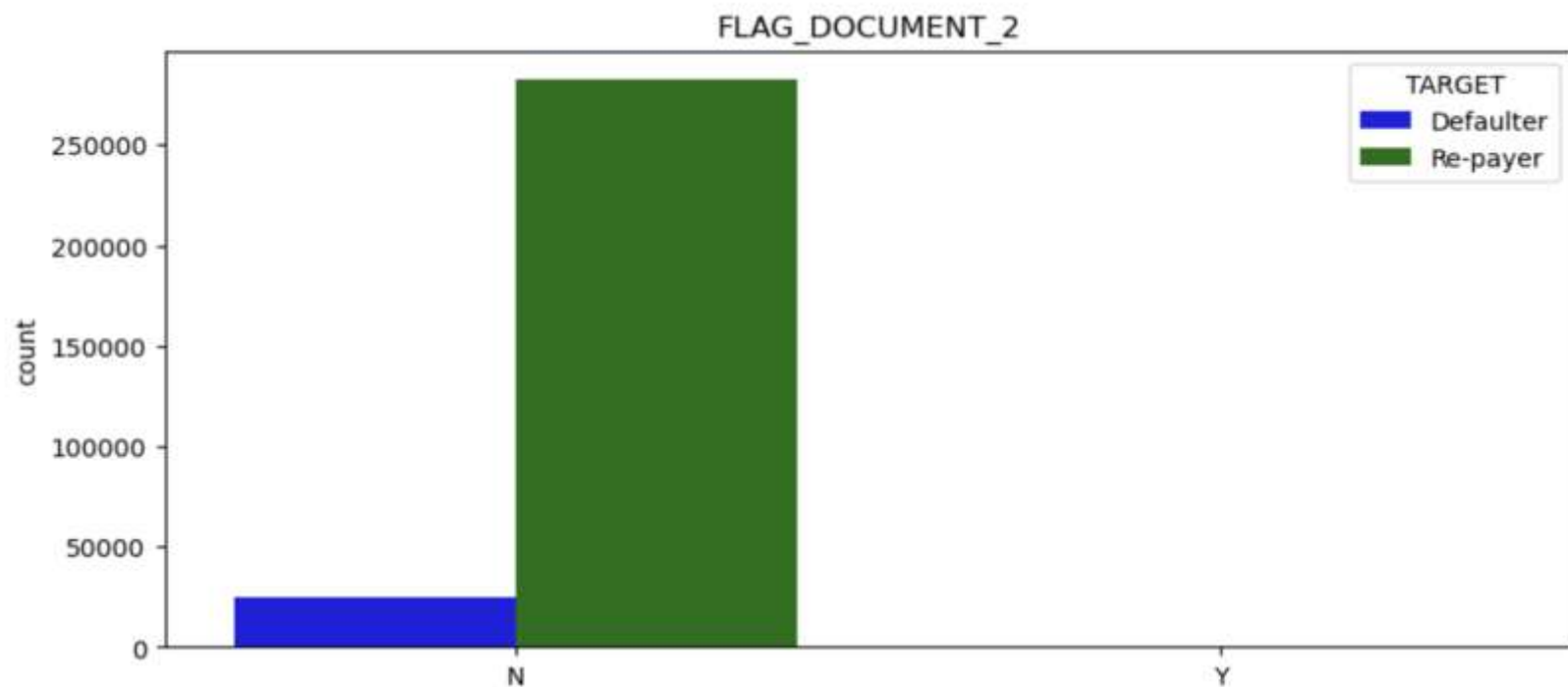
FLAG_WORK_PHONE

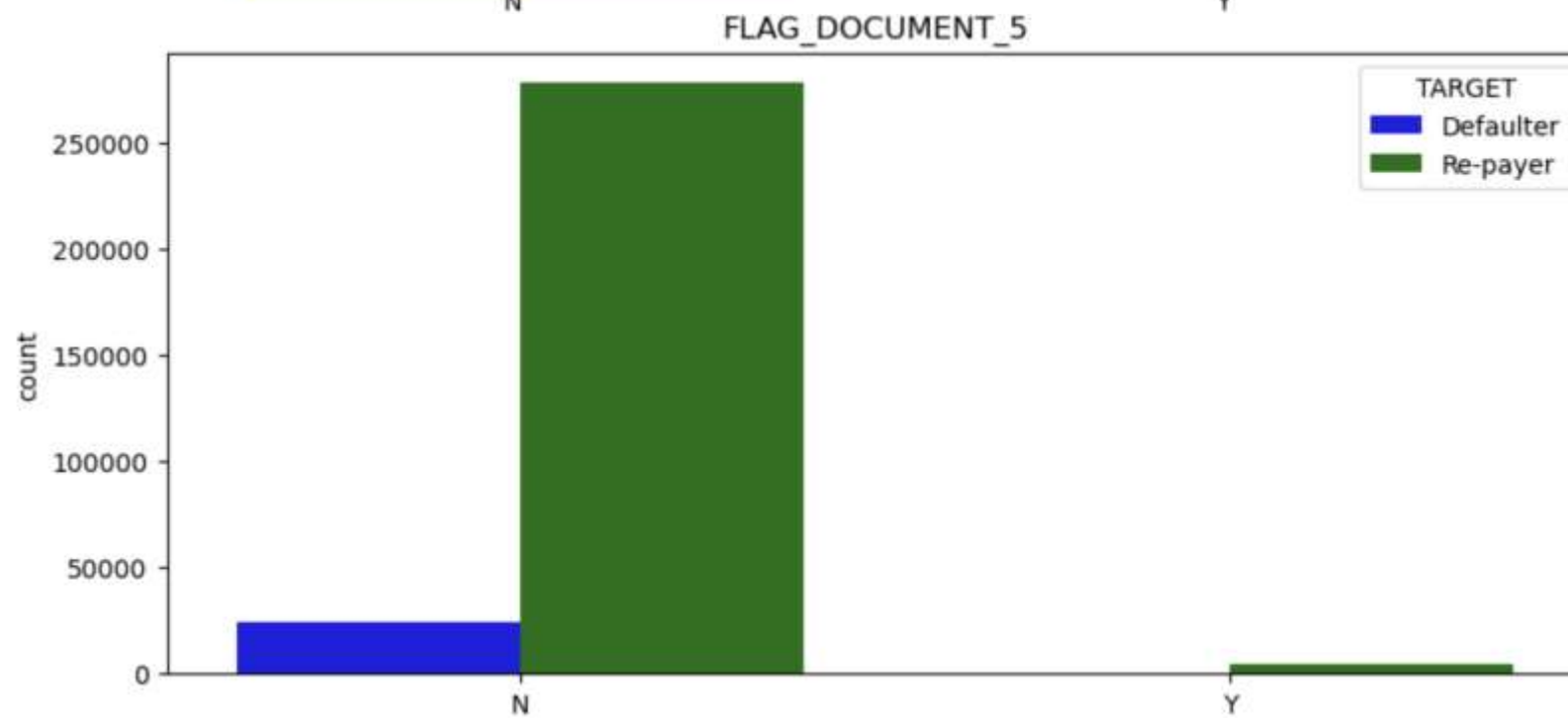
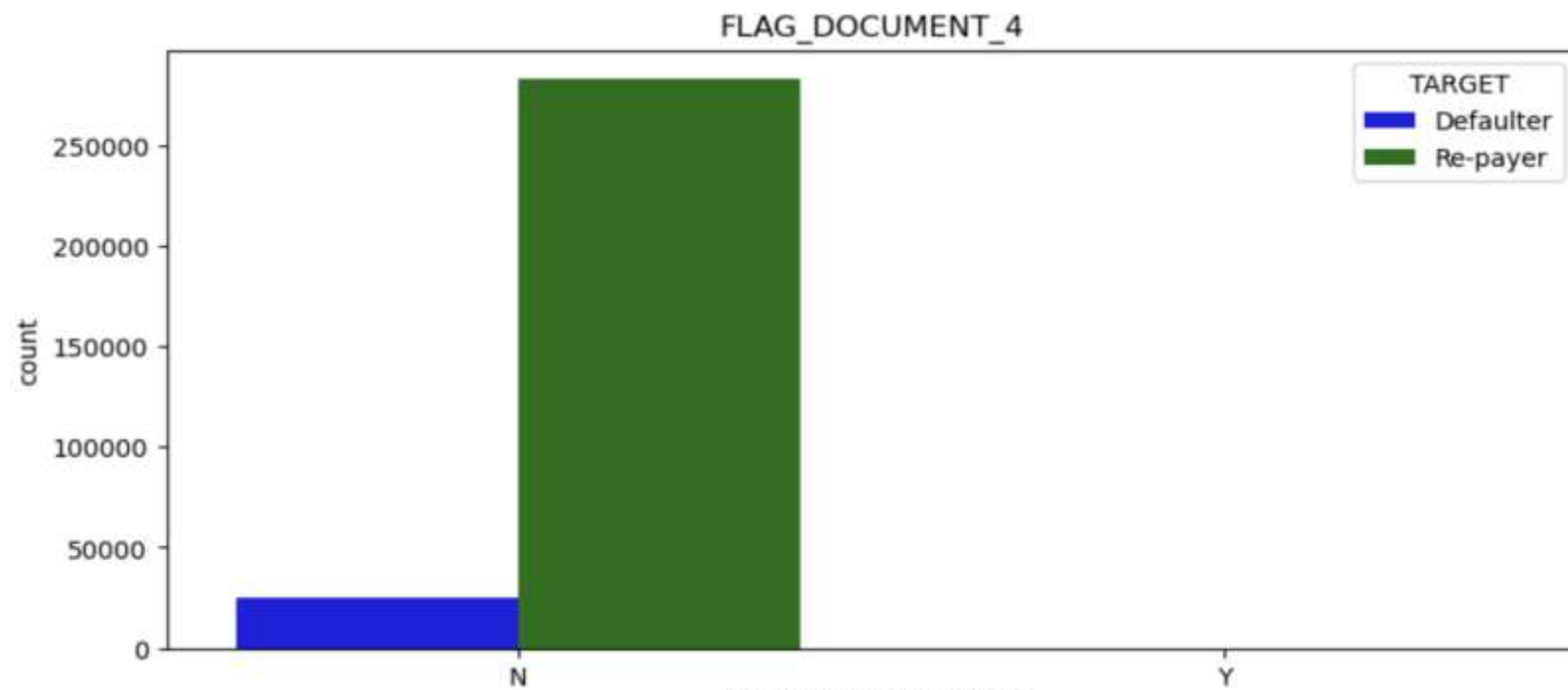


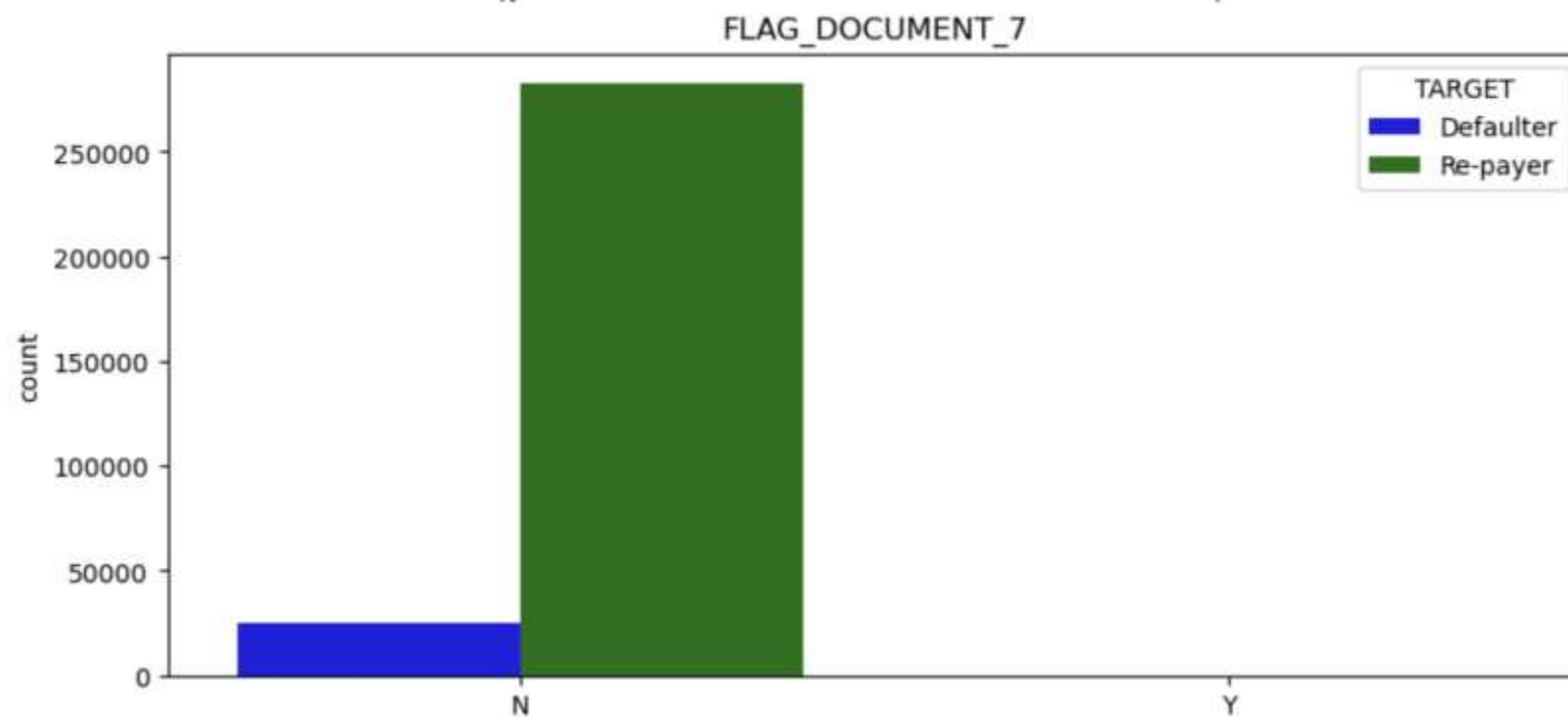
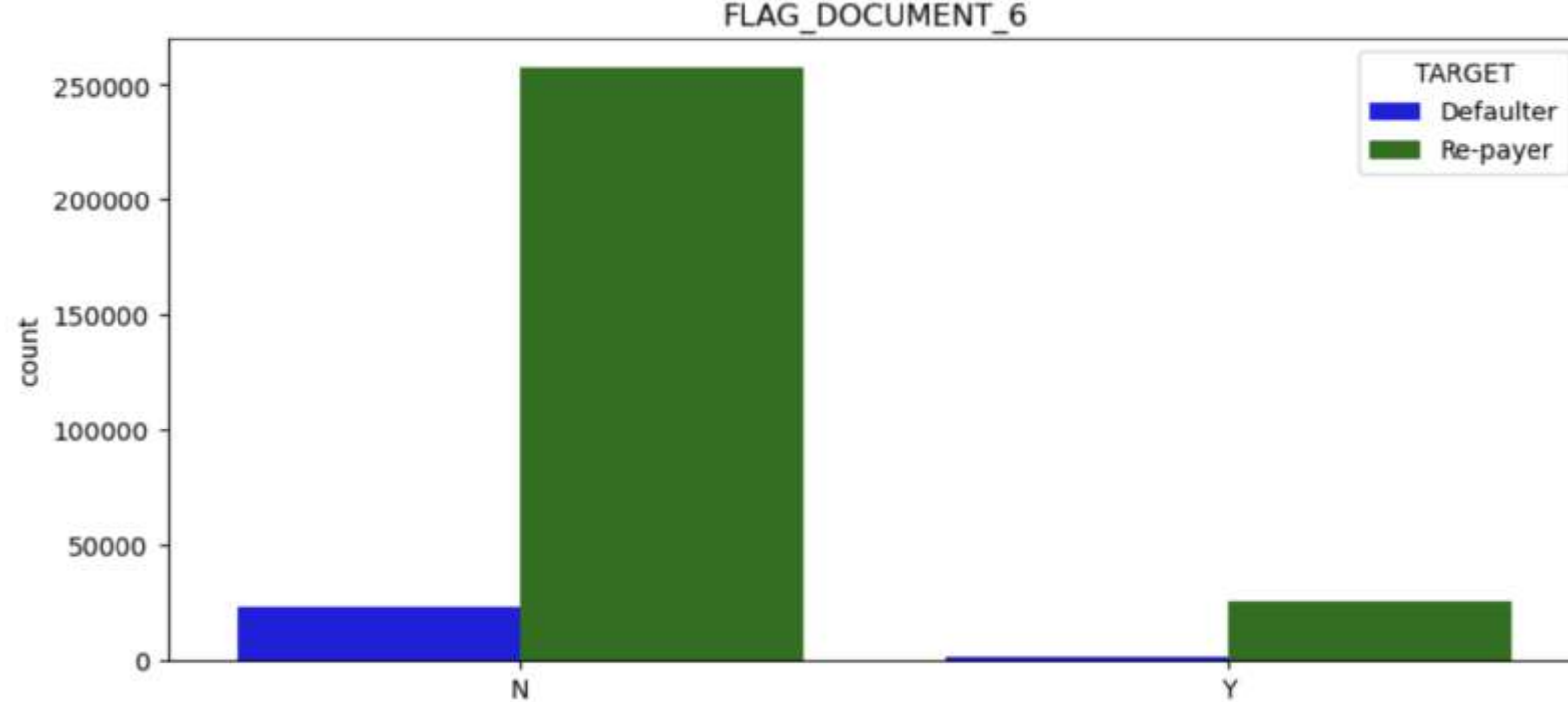
FLAG_CONT_MOBILE

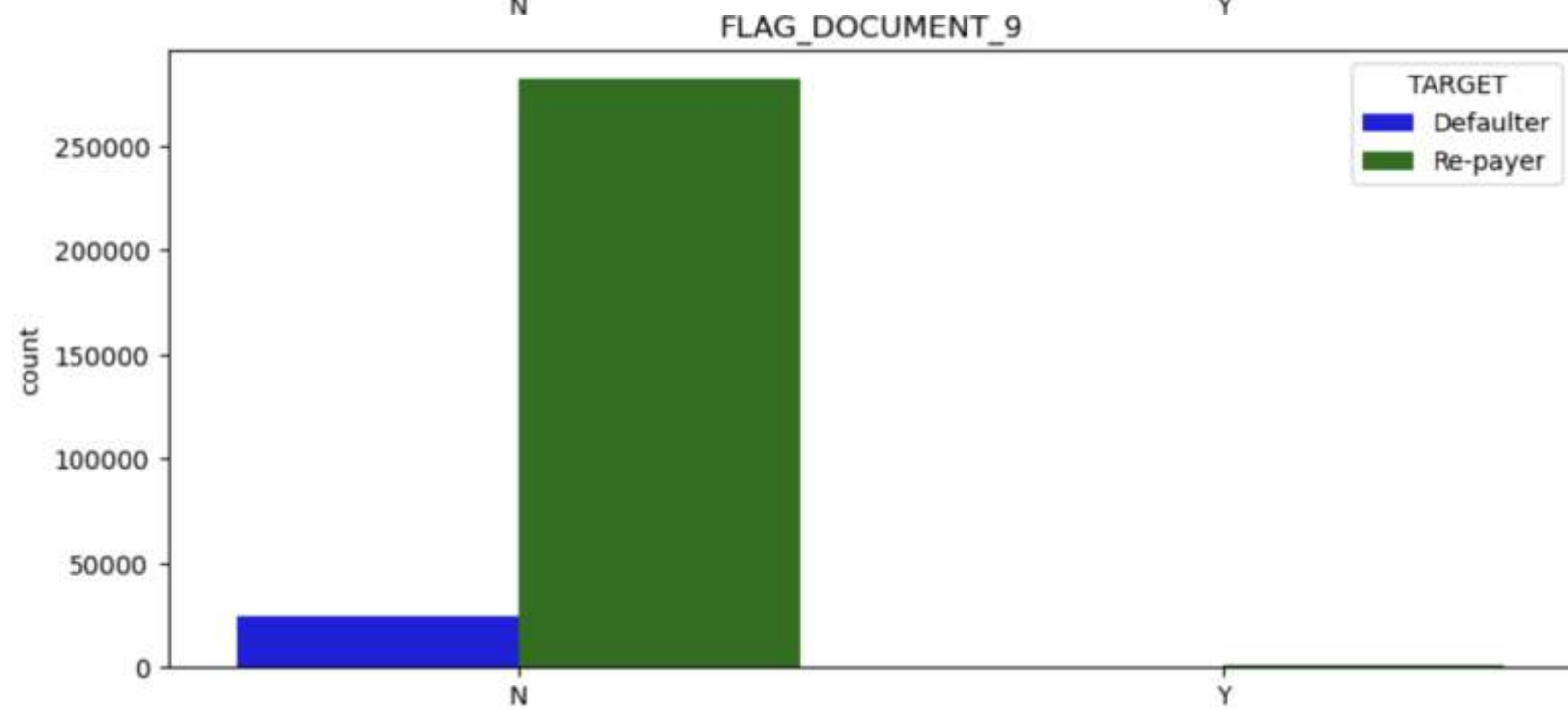
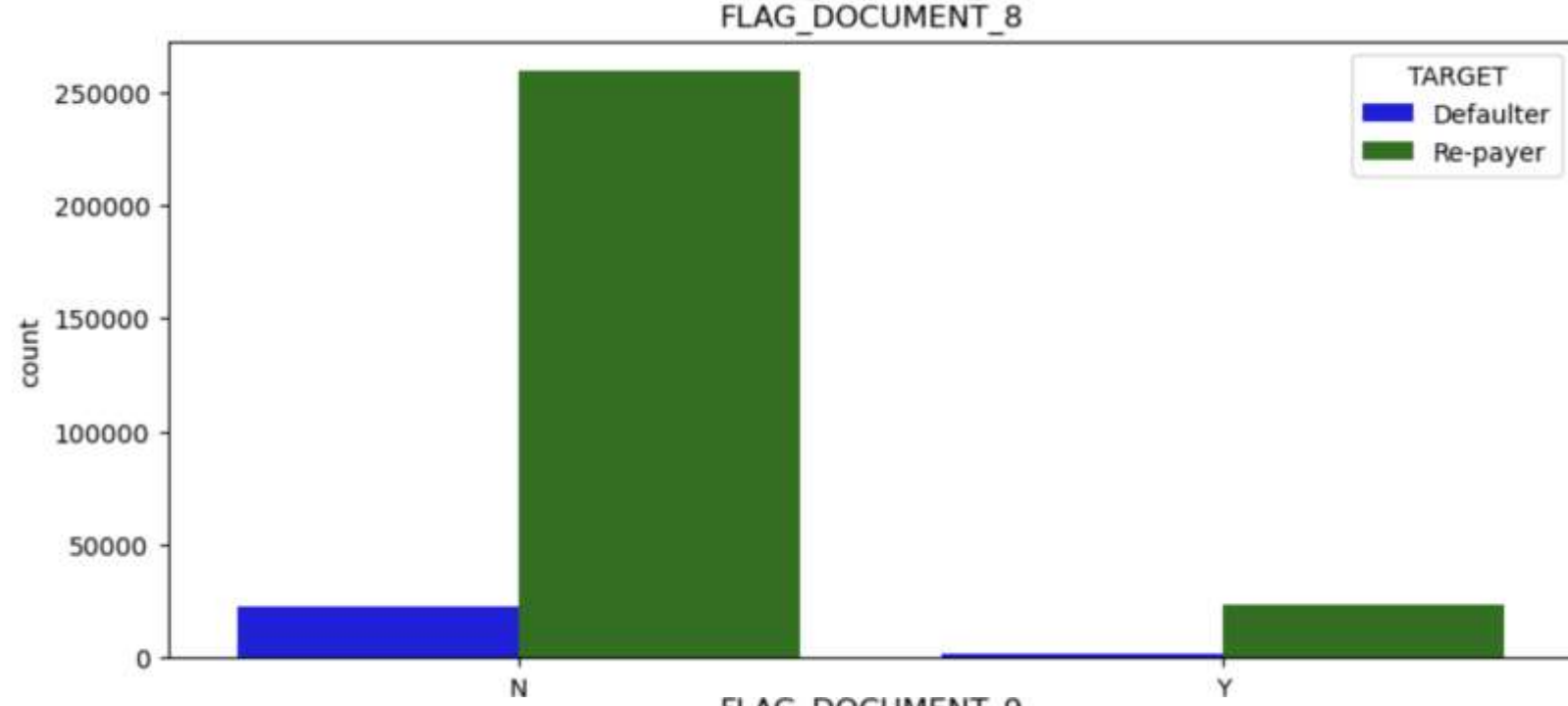


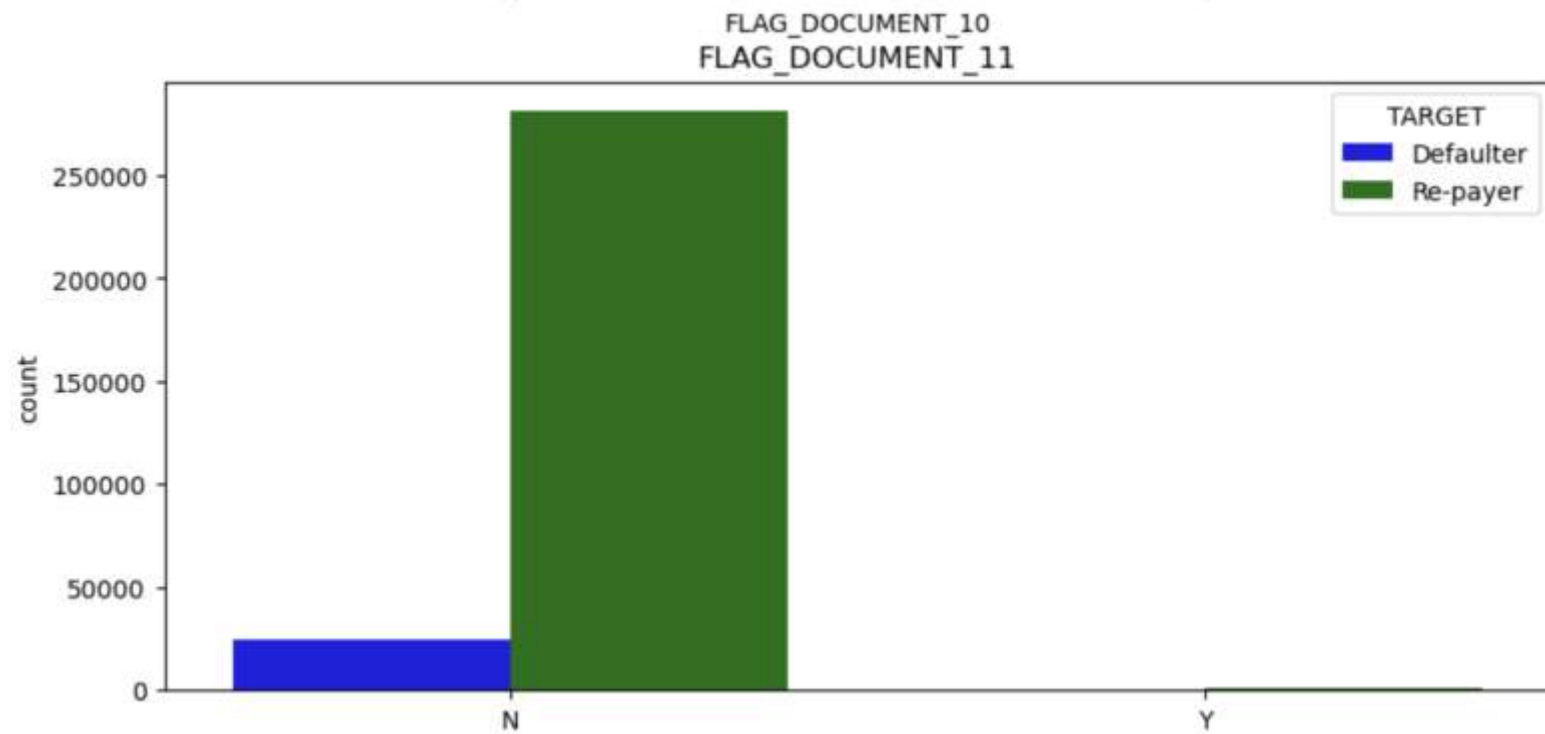
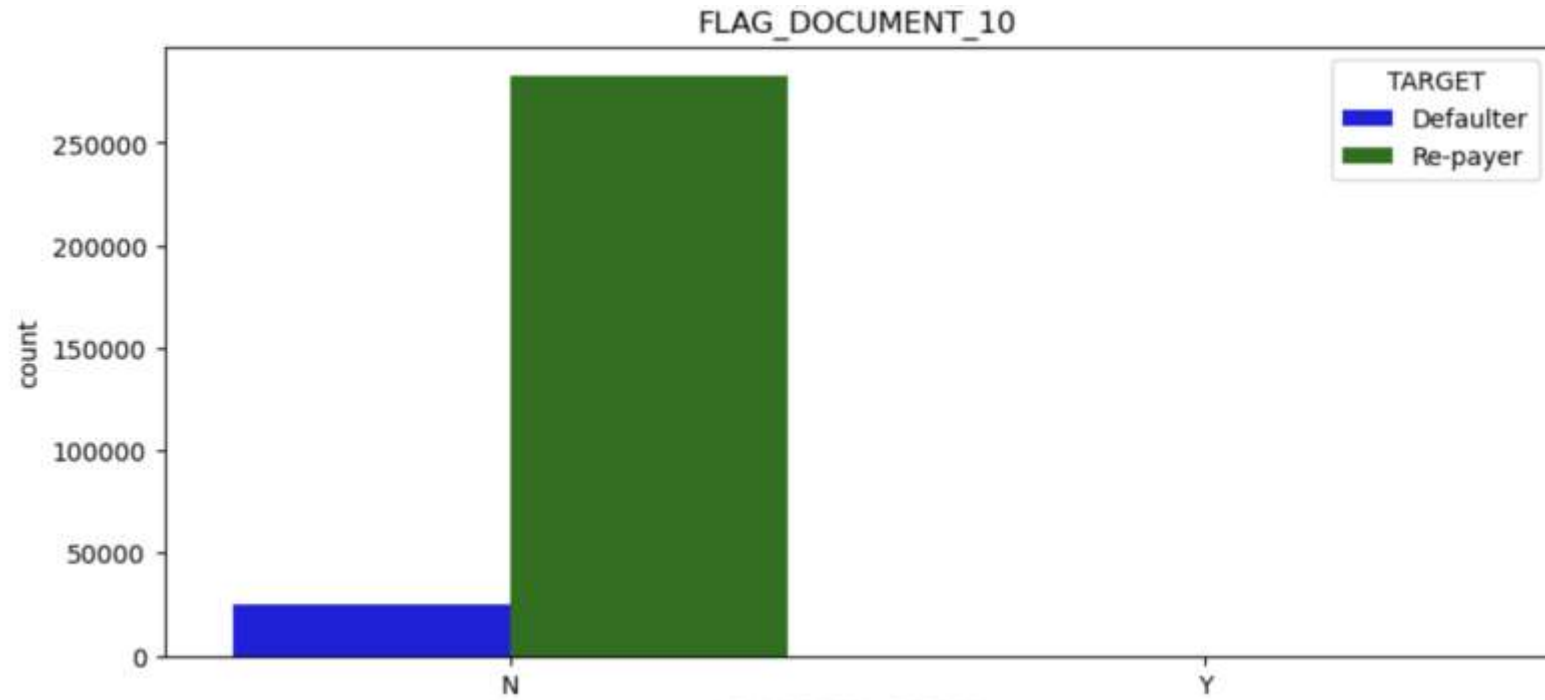


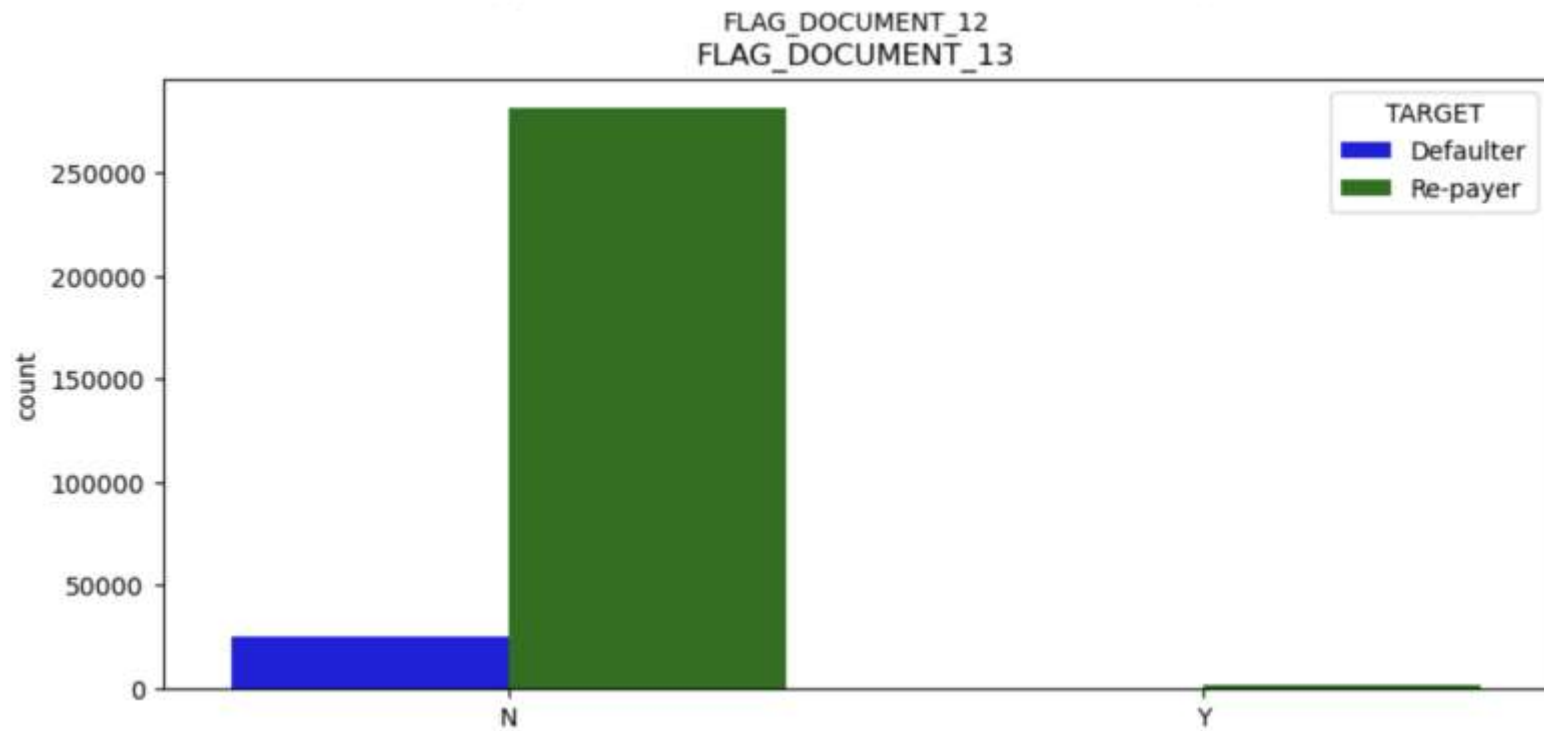
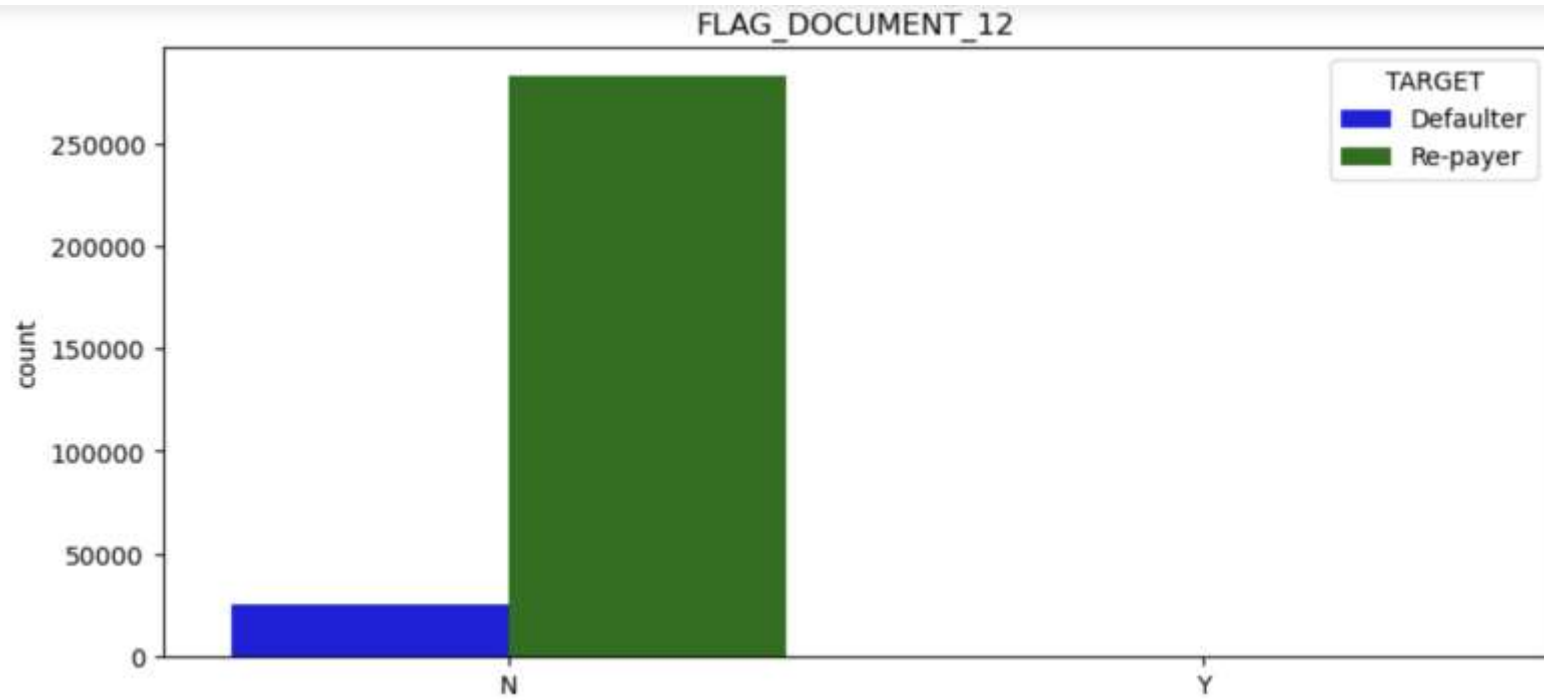


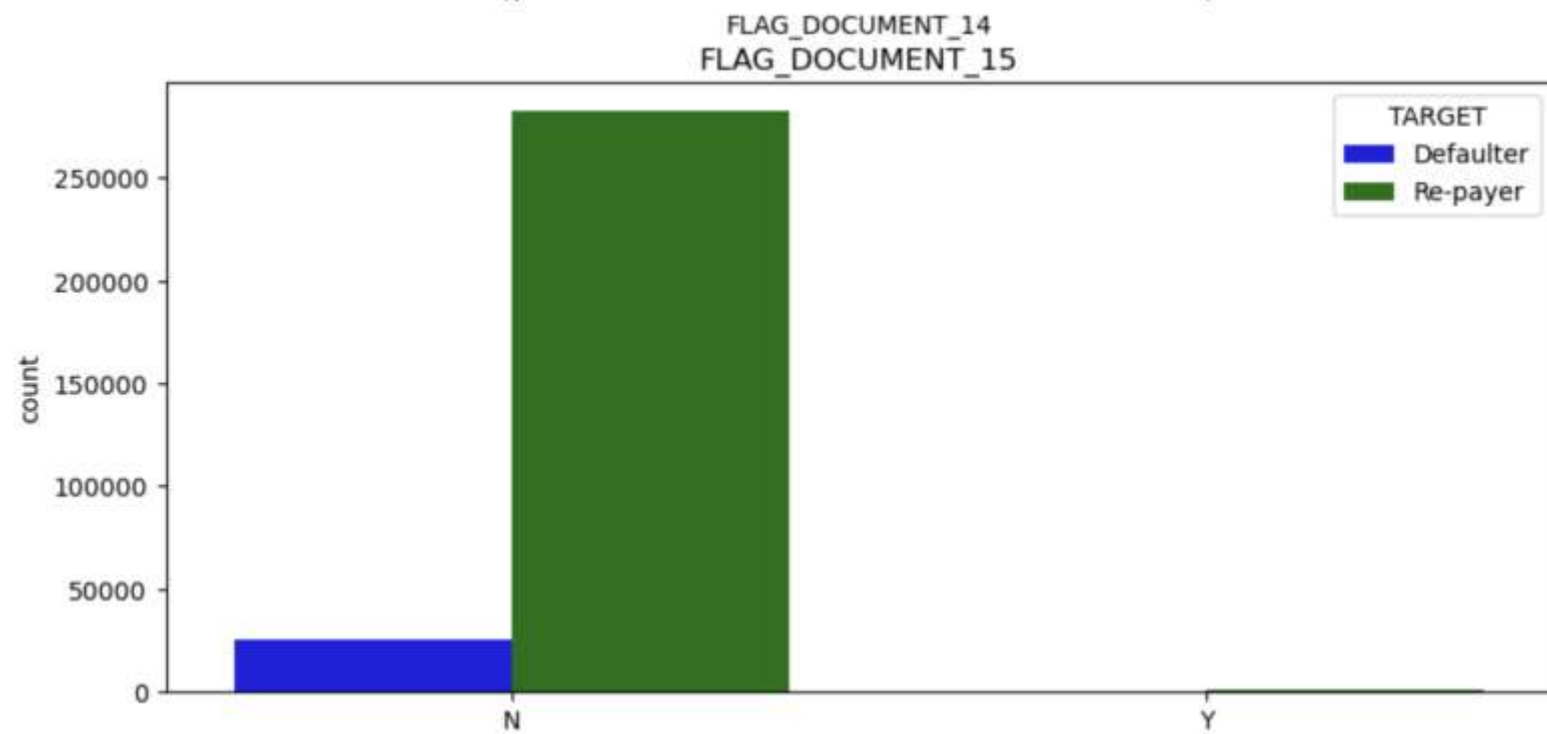
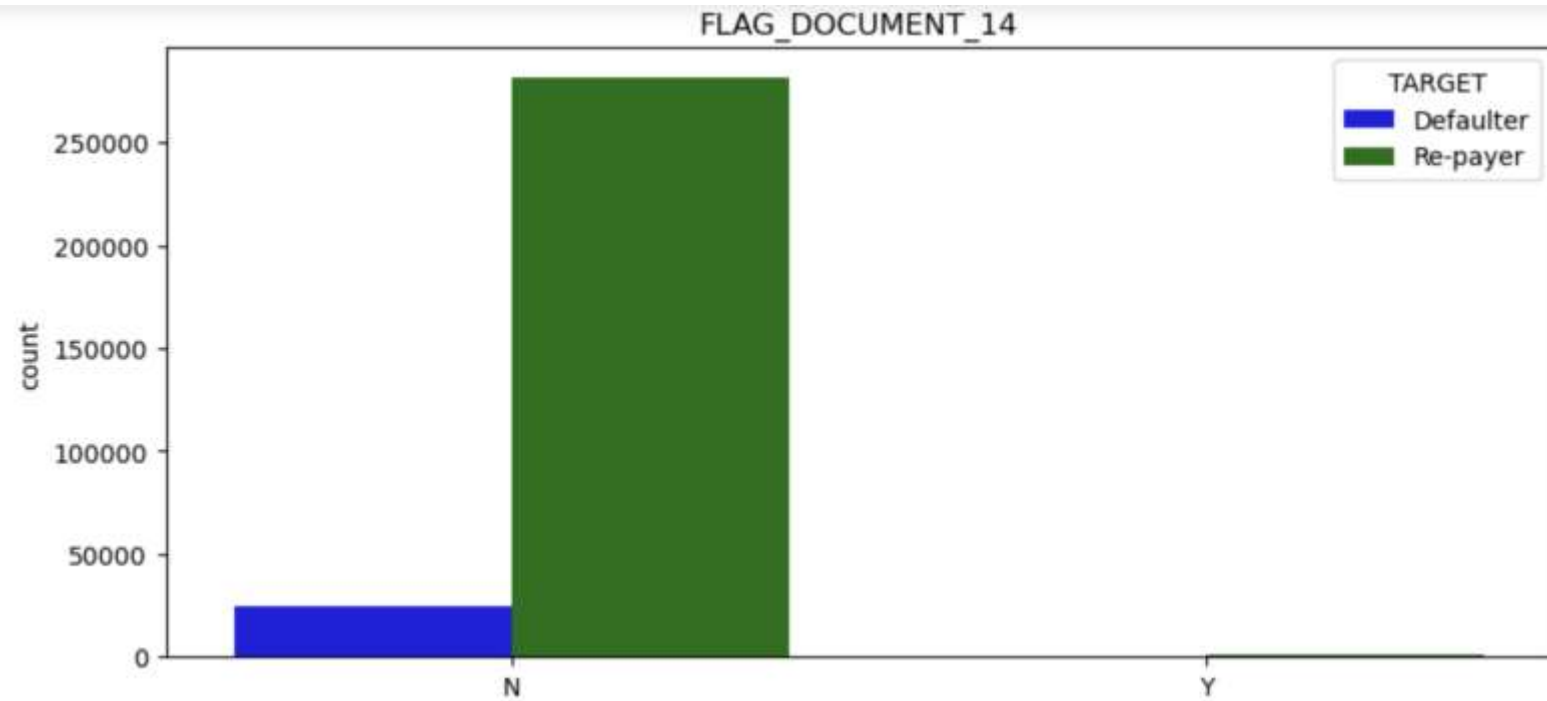


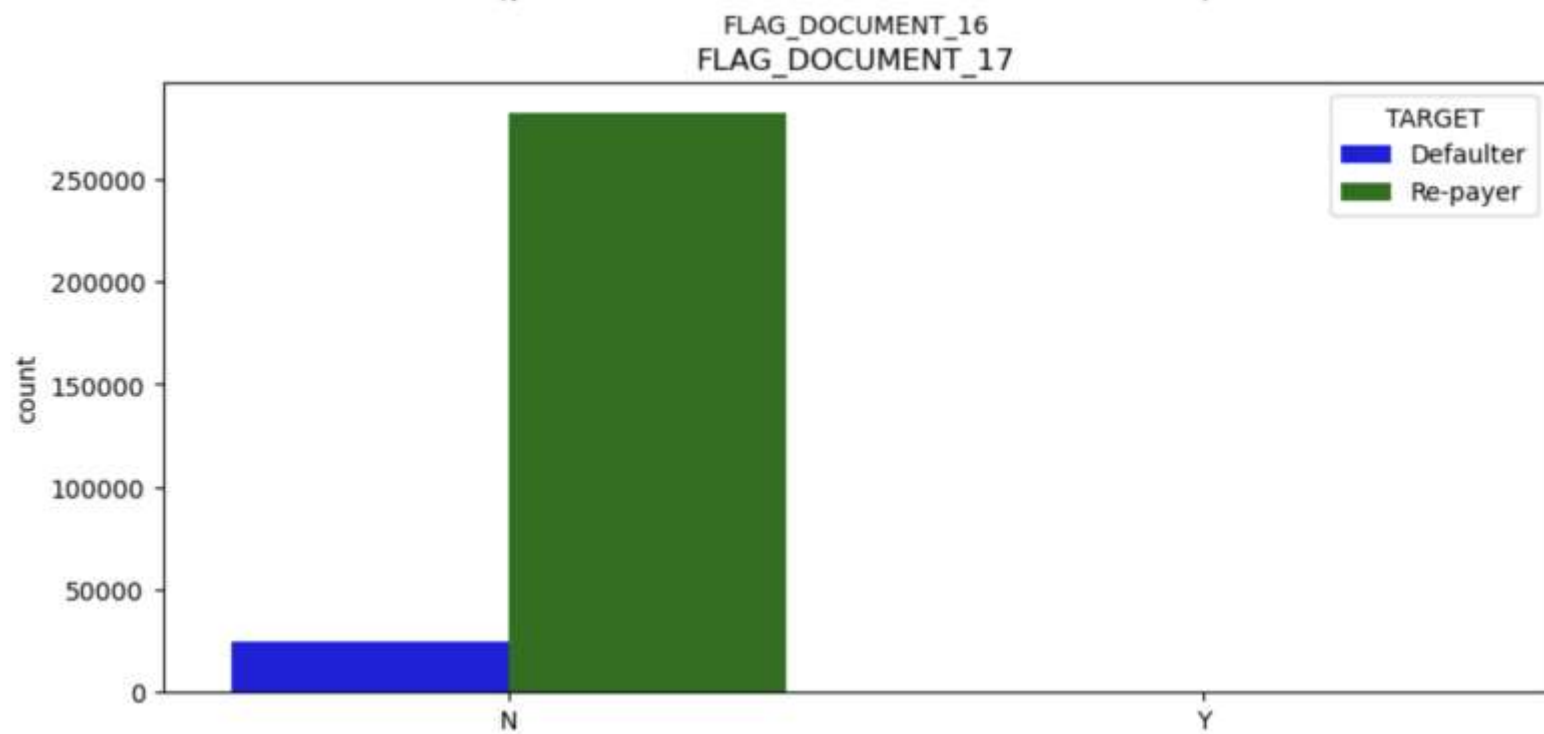
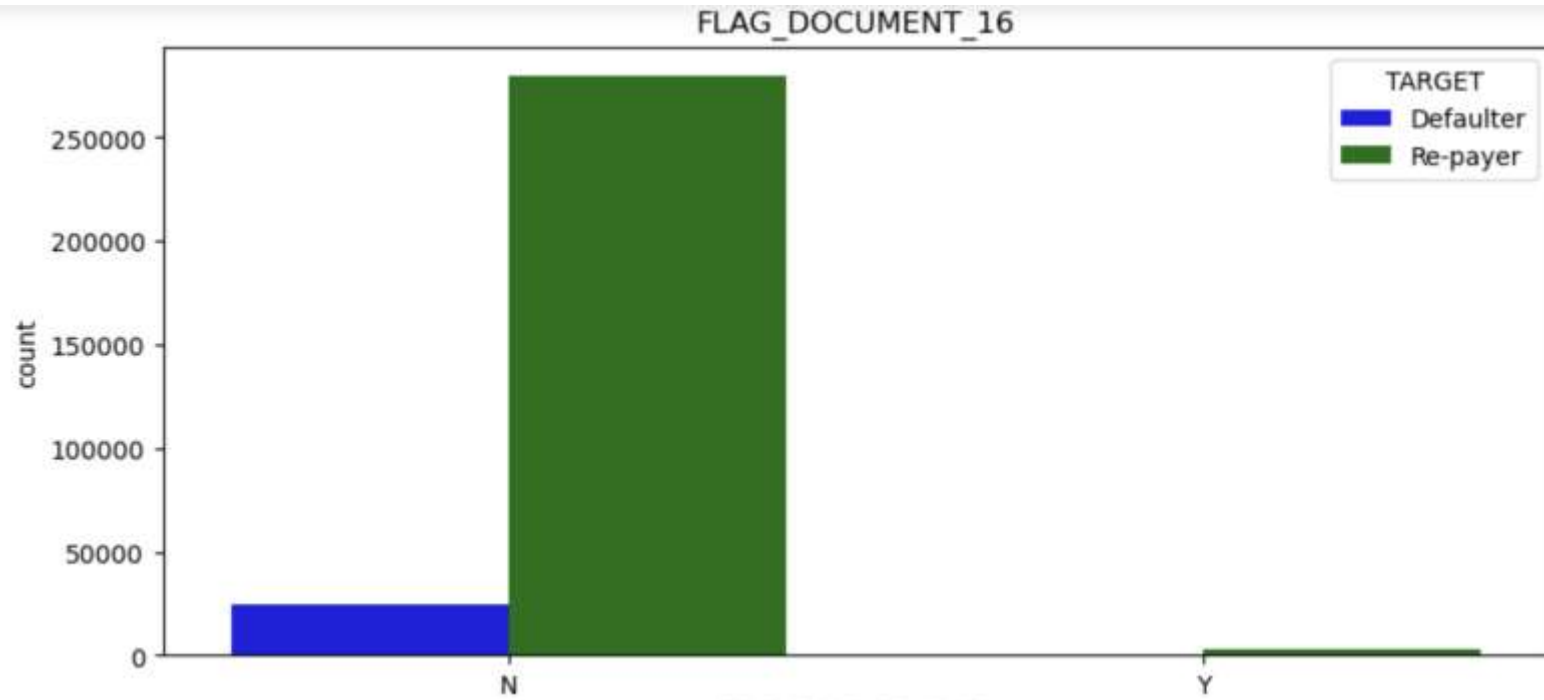


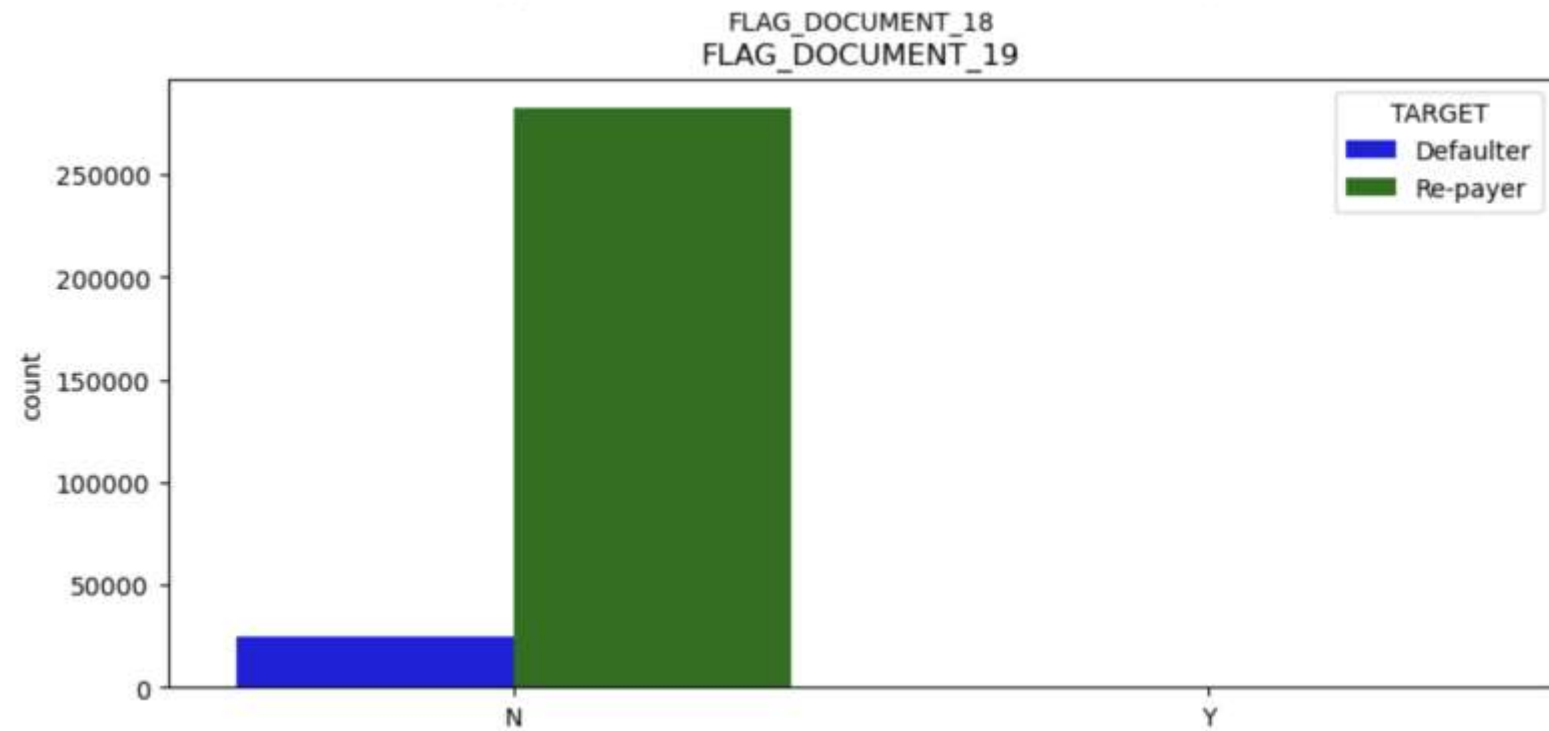
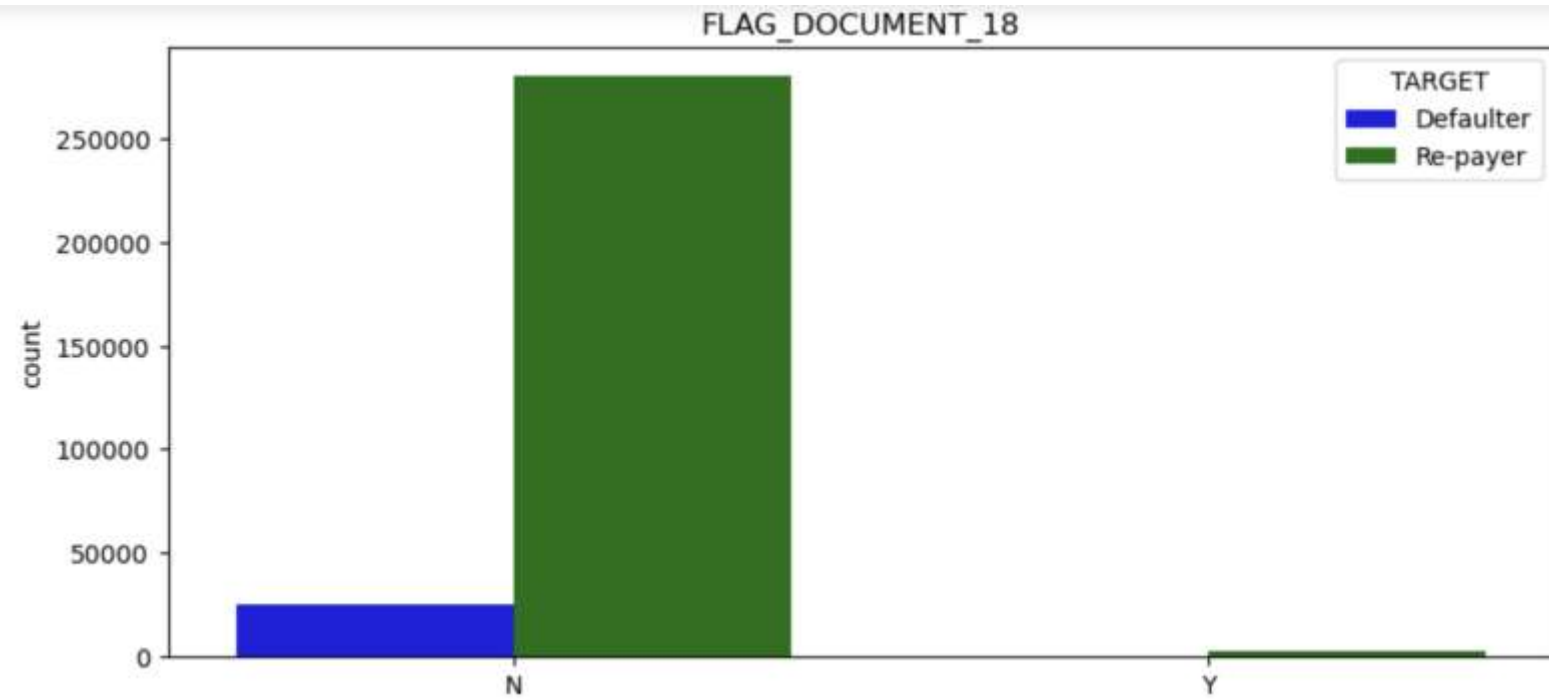


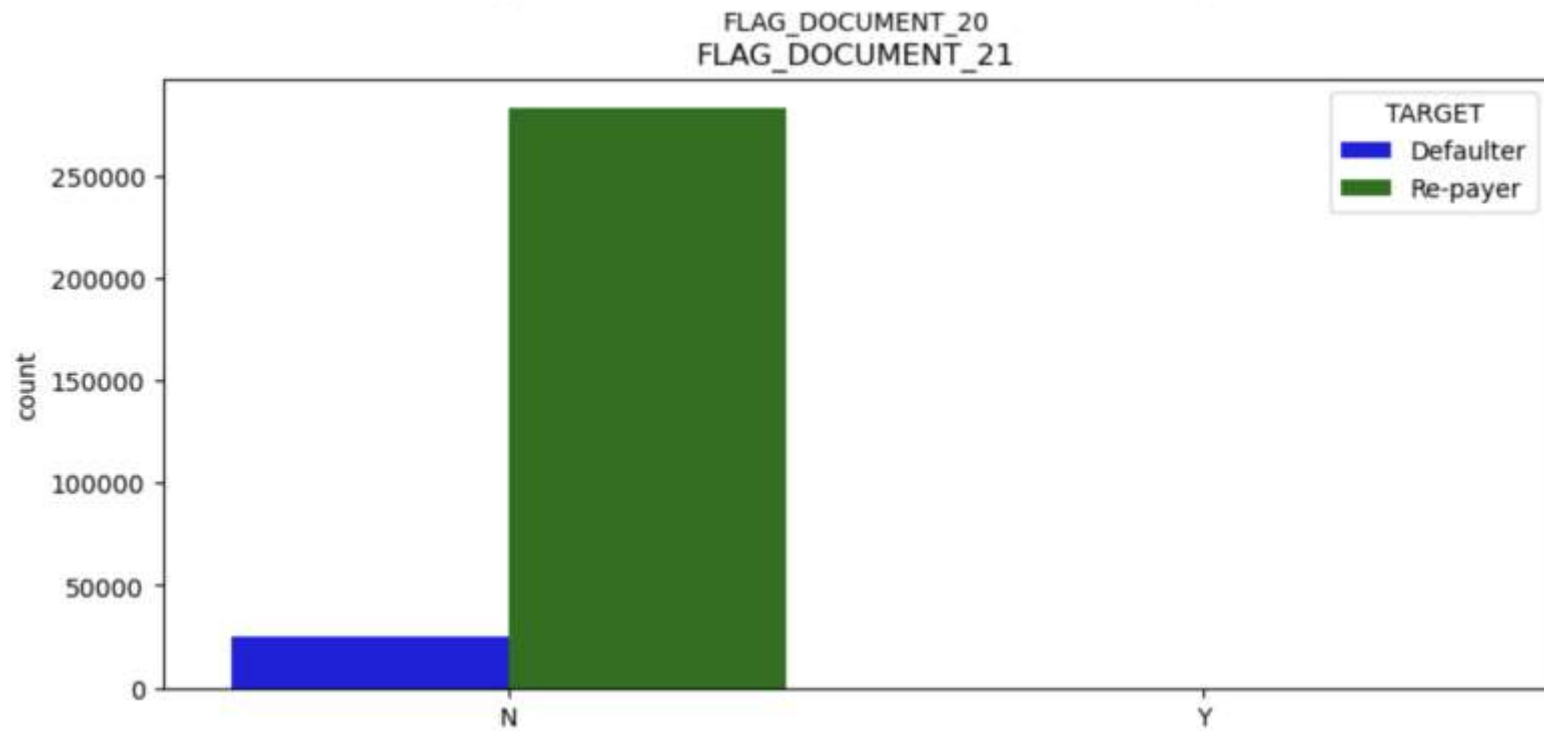
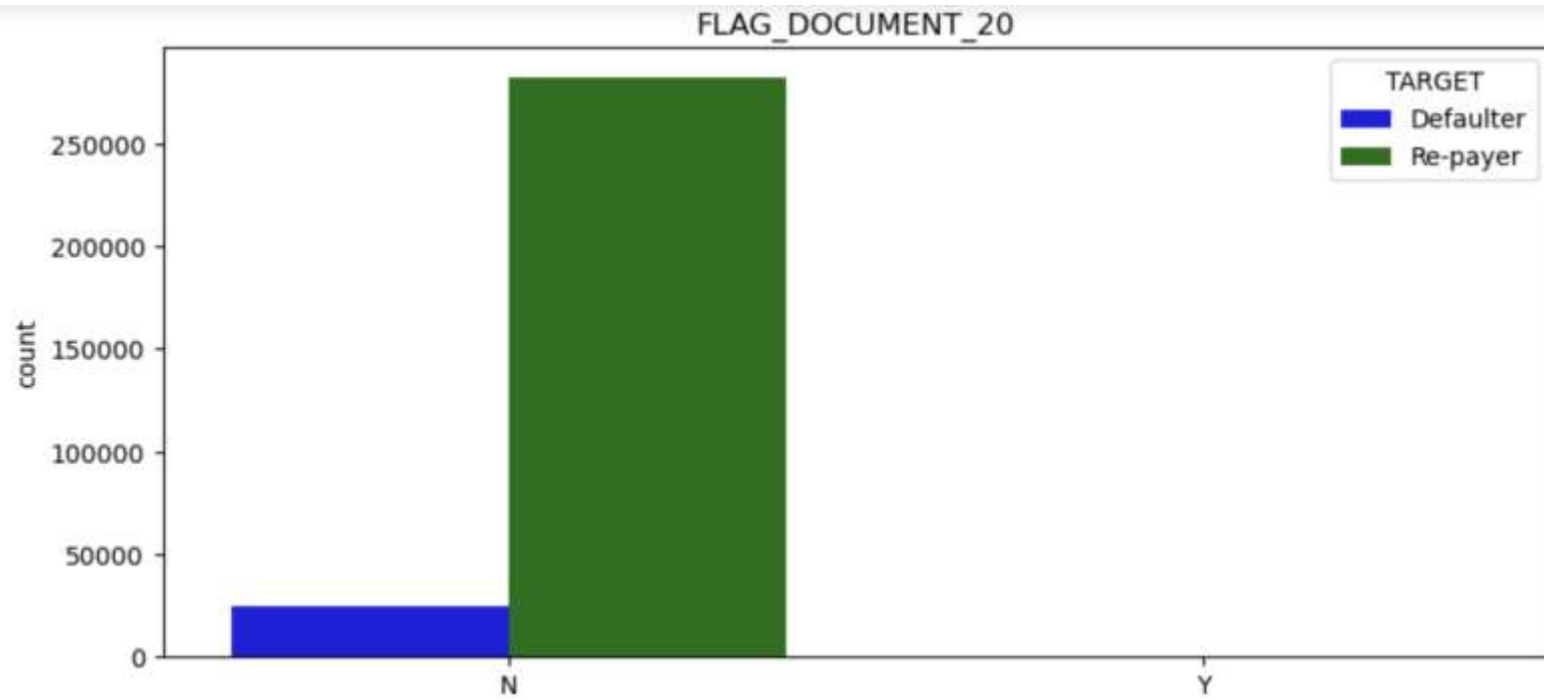










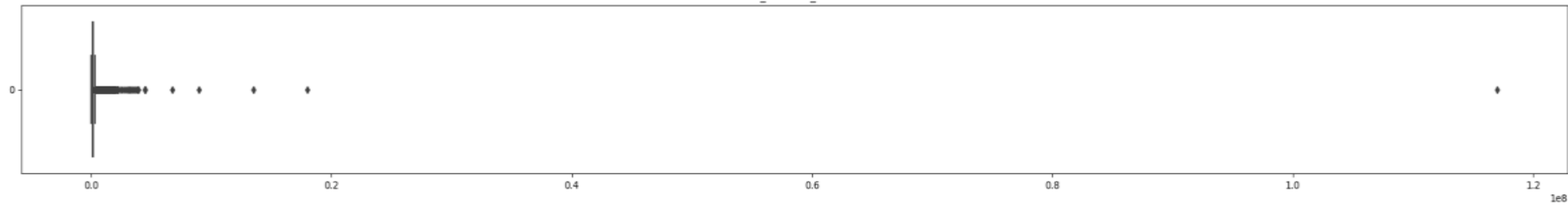


Observations from above plots in slides 4-17:

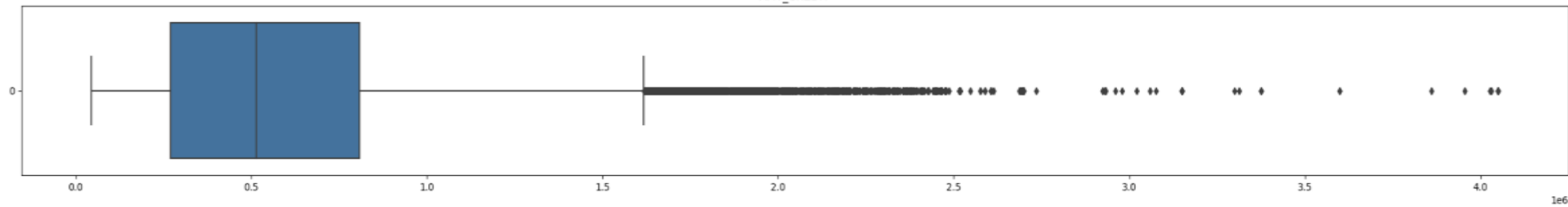
FLAG_OWN_REALTY, FLAG_MOBIL, FLAG_EMP_PHONE, FLAG_CONT_MOBIL, FLAG_DOCUMENT_3 for 'Y' has more re-payers than defaulters. i.e., applicants

1. applicants having own house/flat have been re-payers than defaulters.
2. applicants who provided mobile number have more re-payers than defaulters.
3. applicants whose phones are reachable easily are re-payers than defaulters.
4. applicants who provided Document-3 are likely to be re-payers than defaulters

AMT_INCOME_TOTAL



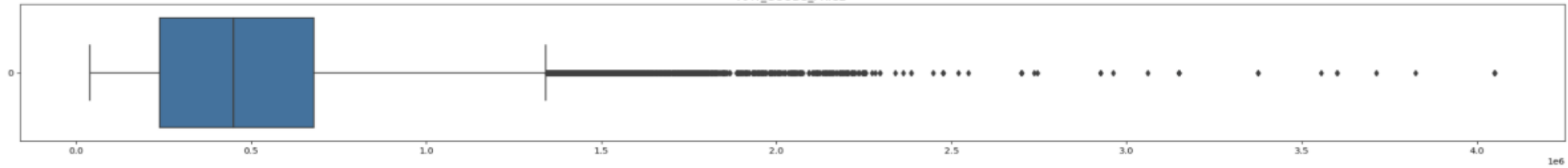
AMT_CREDIT



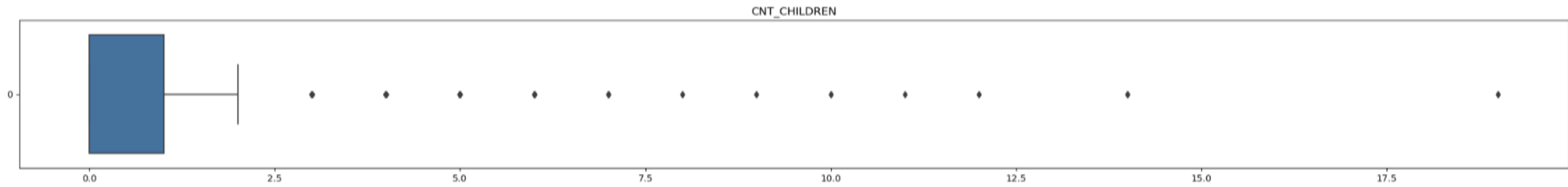
AMT_ANNUITY



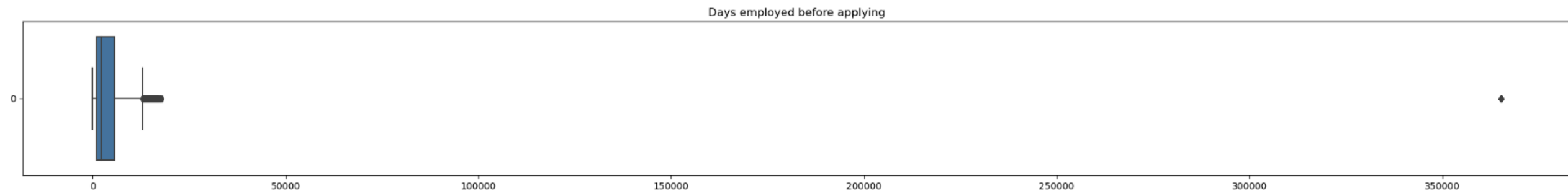
AMT_GOODS_PRICE



'AMT_INCOME_TOTAL', 'AMT_CREDIT', 'AMT_ANNUITY',
'AMT_GOODS_PRICE' have a lot of outliers in these columns, we can either drop the rows that corresponding to these outliers and proceed with the analysis, but in our case, since I think, these are some of the significant columns, I wish to keep them intact for our analysis.

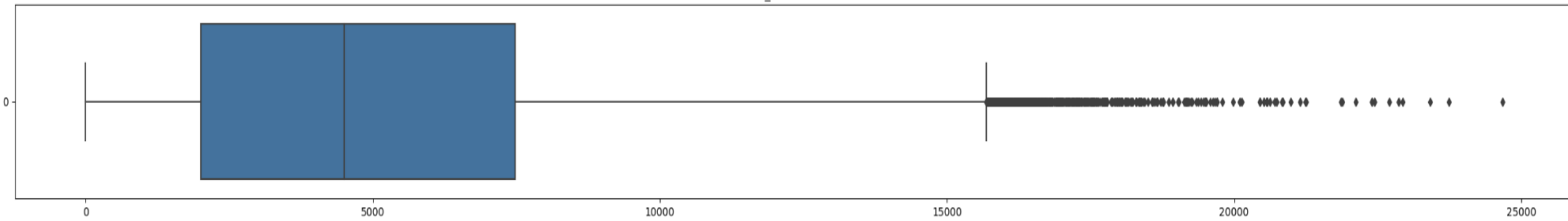


- As evident from the description of CNT_CHILDREN column and also its corresponding boxplot, 1st quartile is completely missing from the boxplot as most of the data is present in the second quartile ie, above 75%.
- Also, as can be seen from the boxplot, there are a very few outliers corresponding to number of children usually greater than 4 which is a rarity. hence this data can be termed as reliable

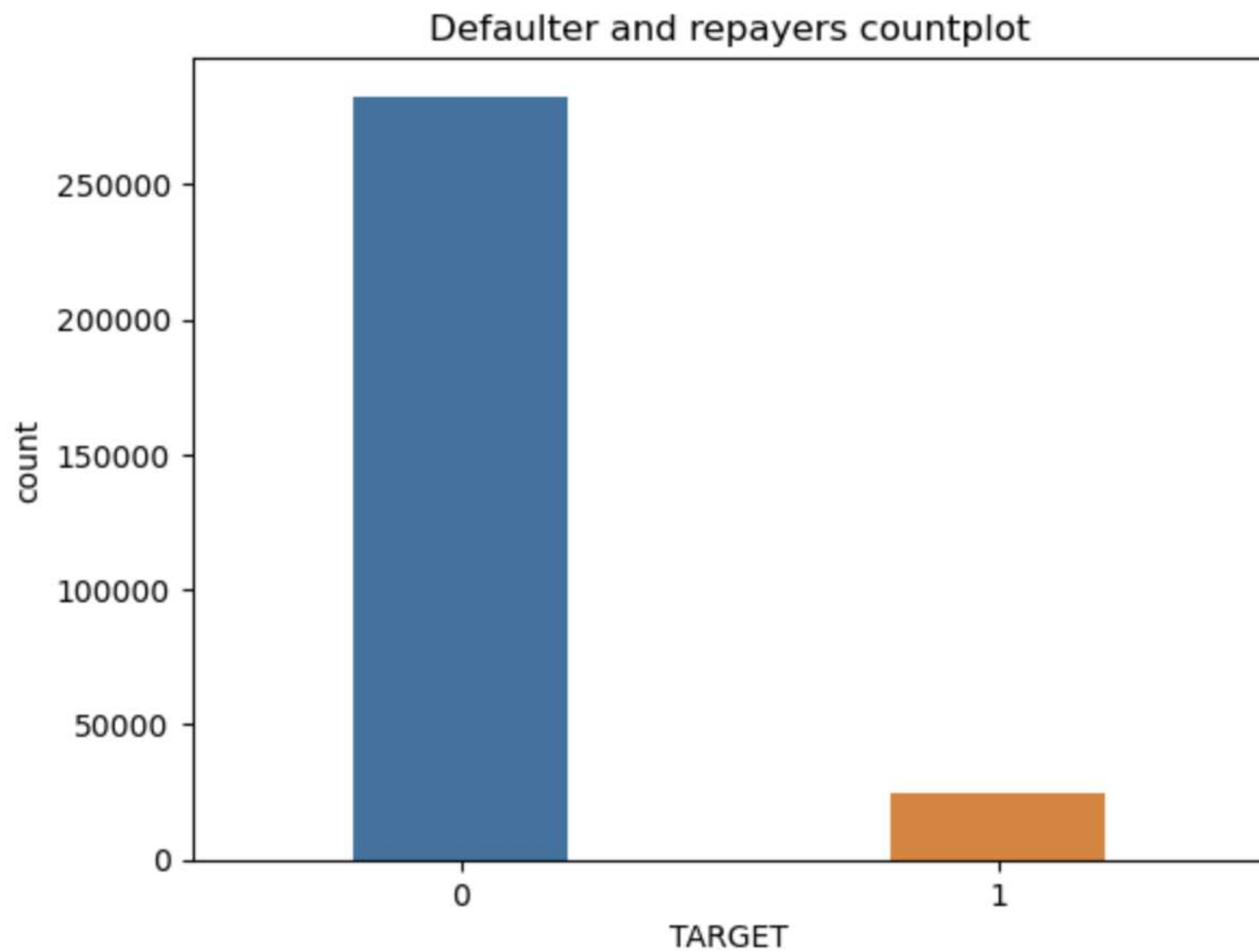


- one of the outlier point is above 350000 days which is roughly 958 years of employment which is not possible and hence we may conclude that it is a wrong data

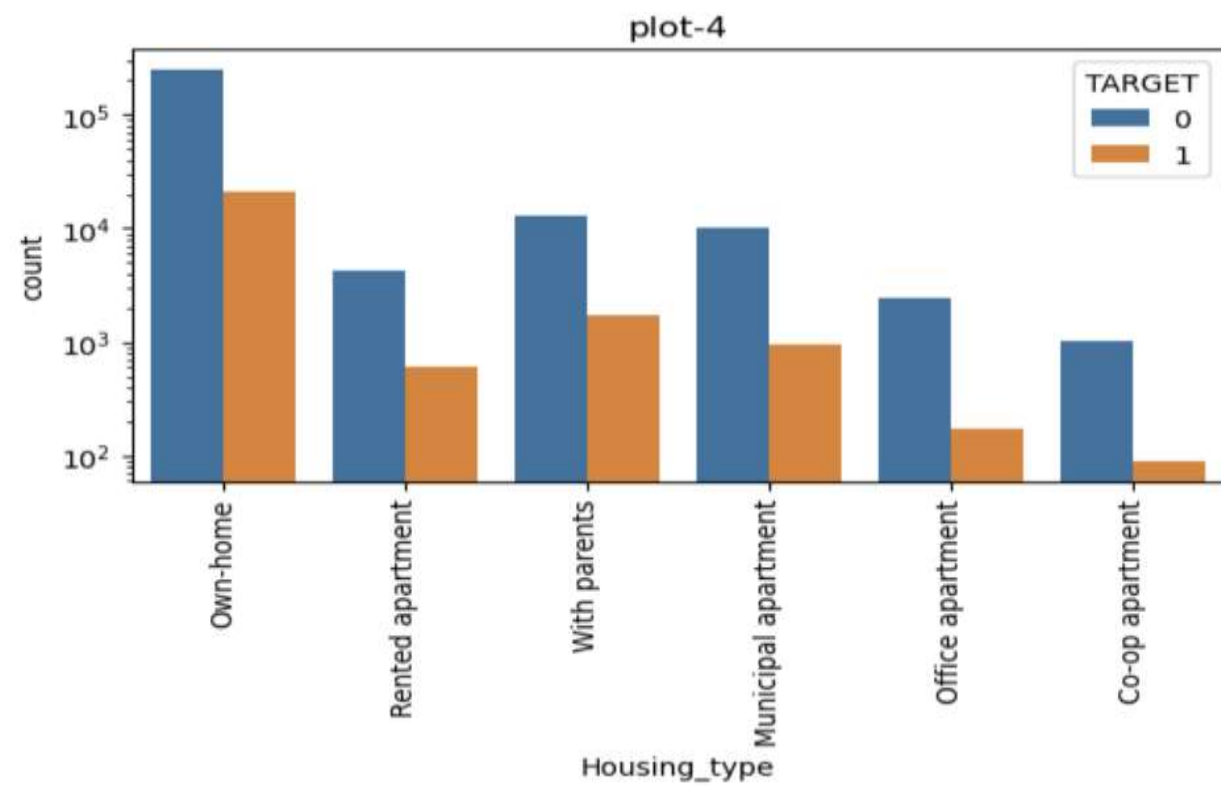
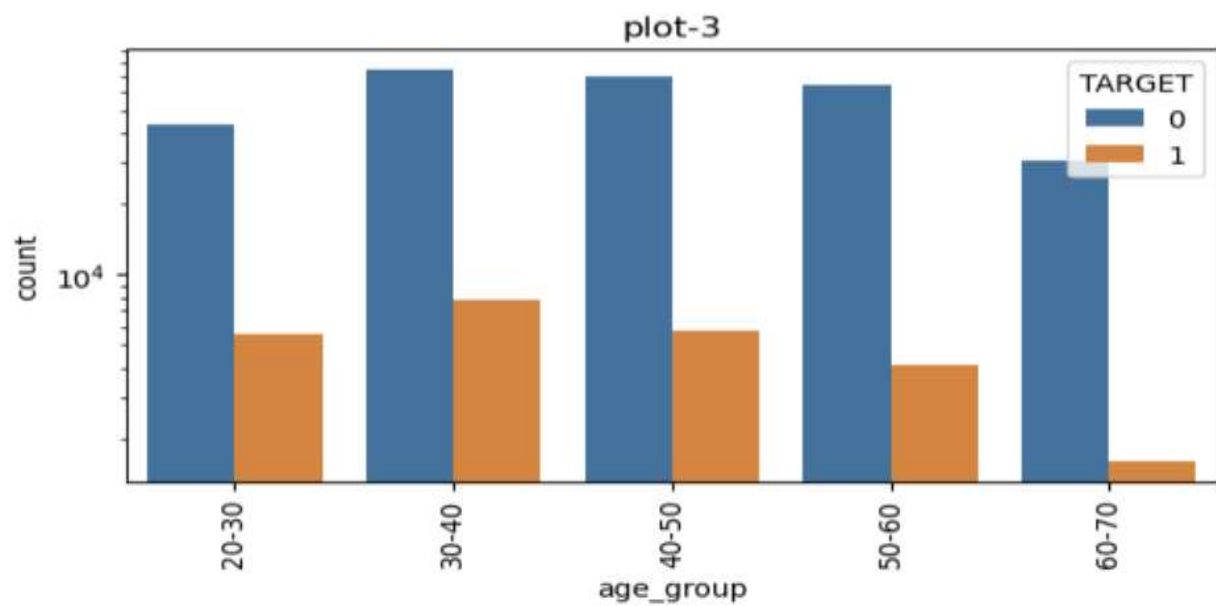
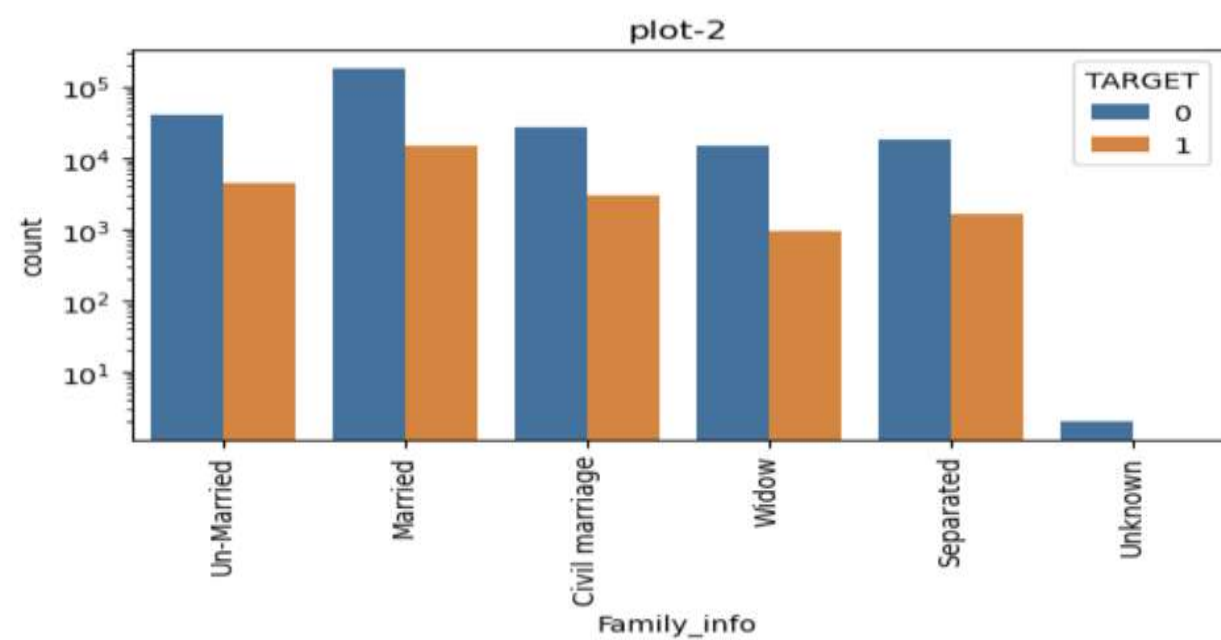
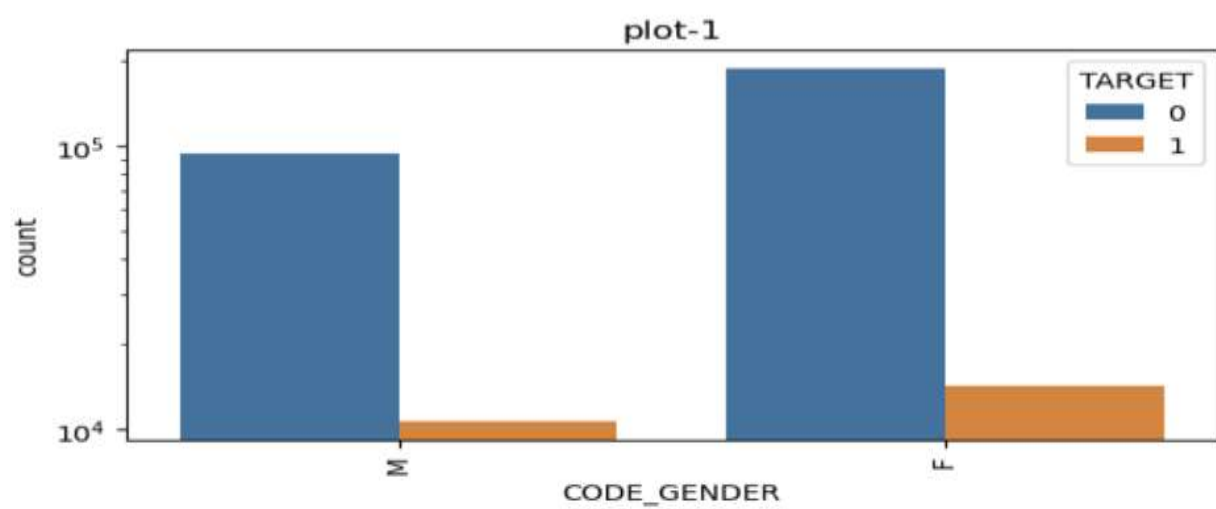
DAYS_REGISTRATION

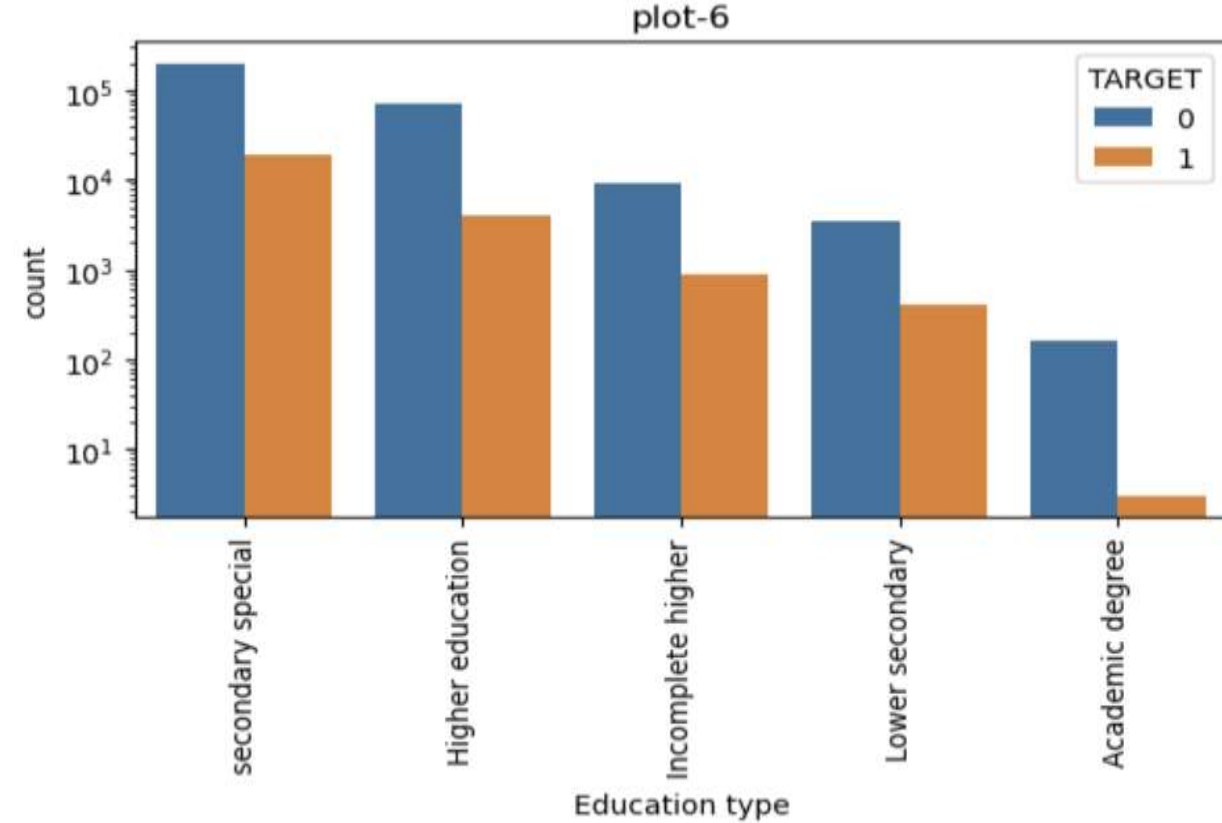
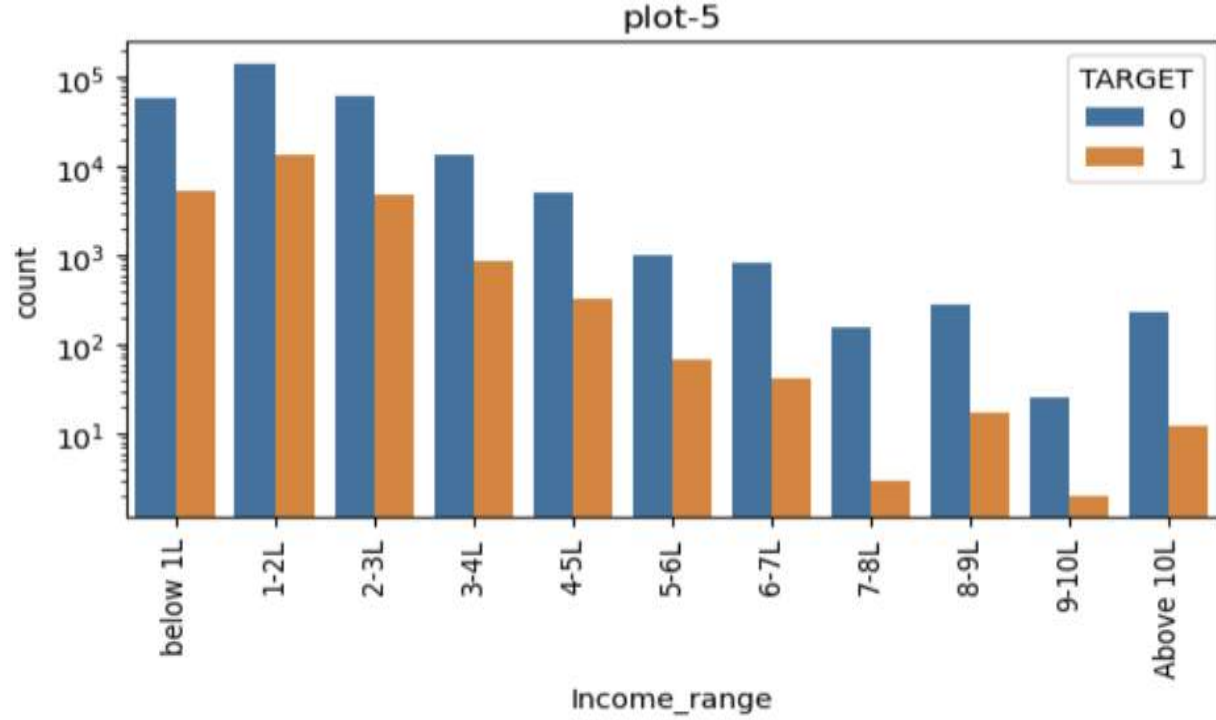


as the max age of applicant is 69 years, the years before which that applicant has changed his registration should be less than 68 years or 25000 days which is true as evident from the boxplot above, hence the data is reliable



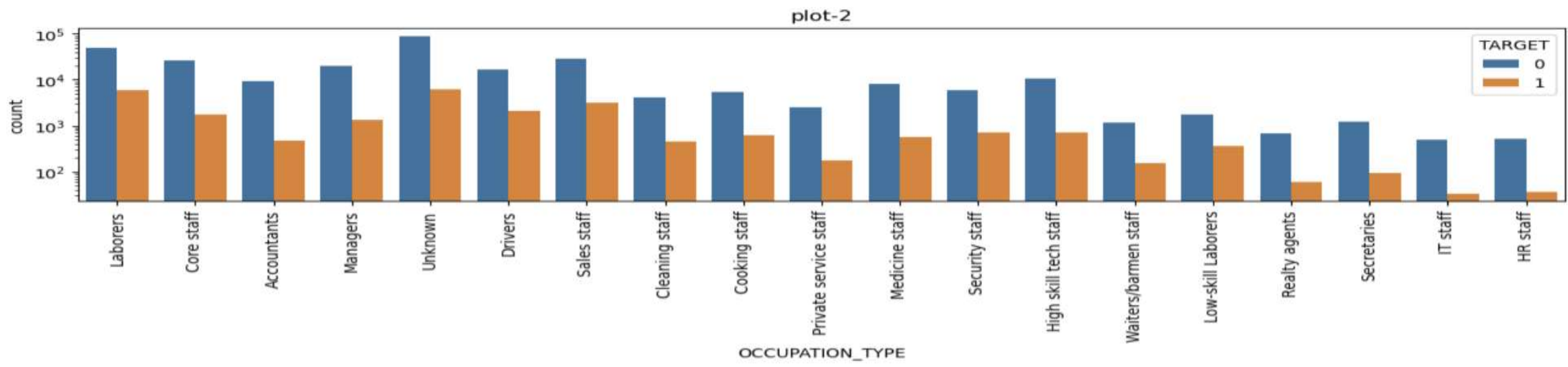
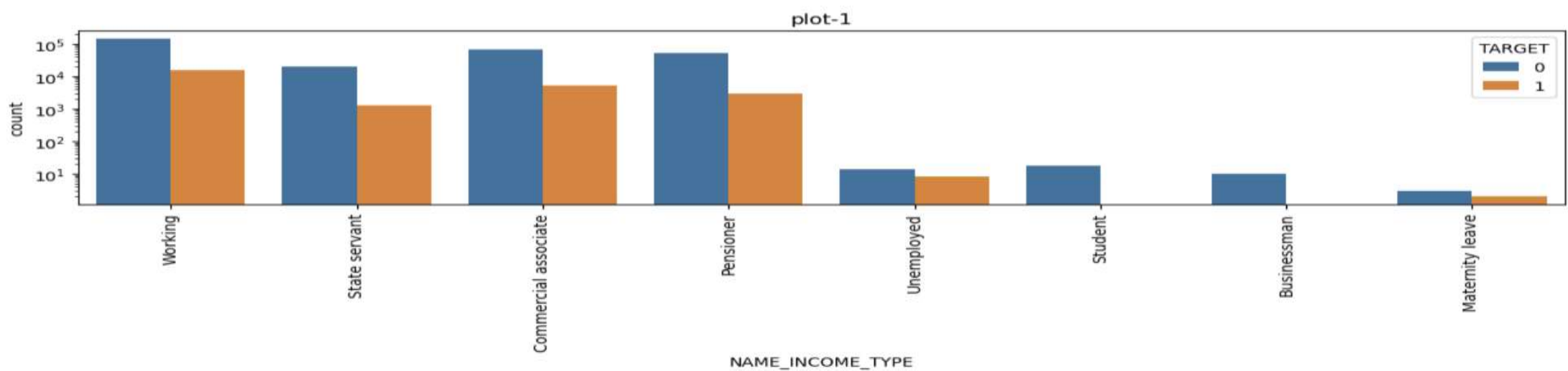
This count shows that there is a huge variation in number of defaulters and repayers

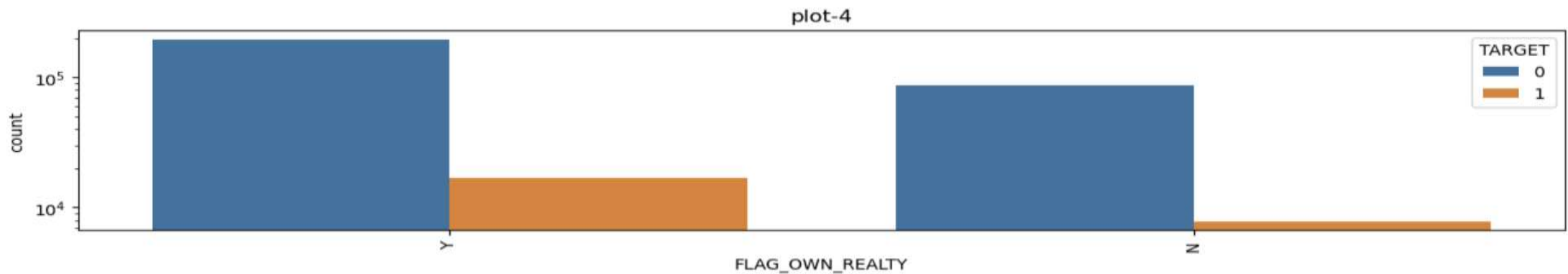
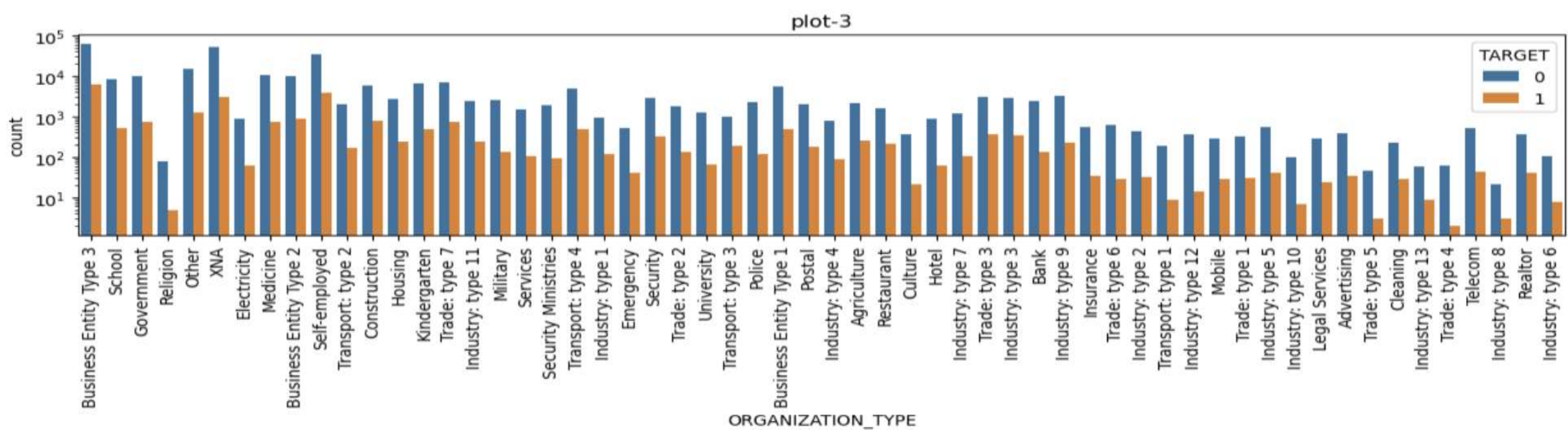




NOTE: inference from the plots above

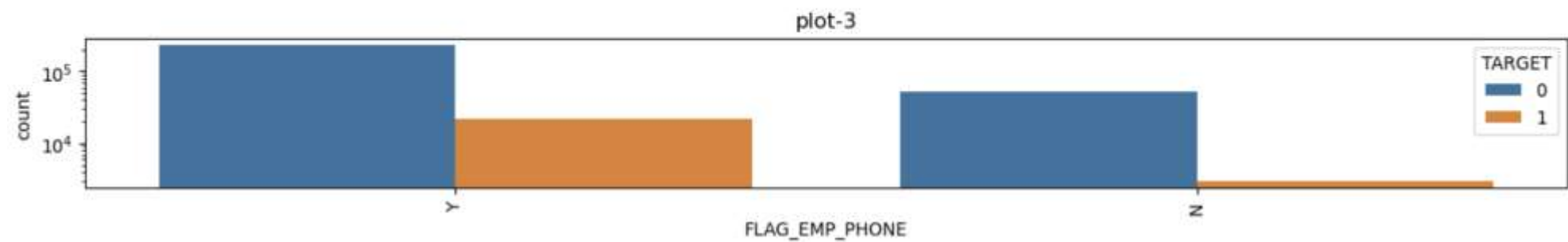
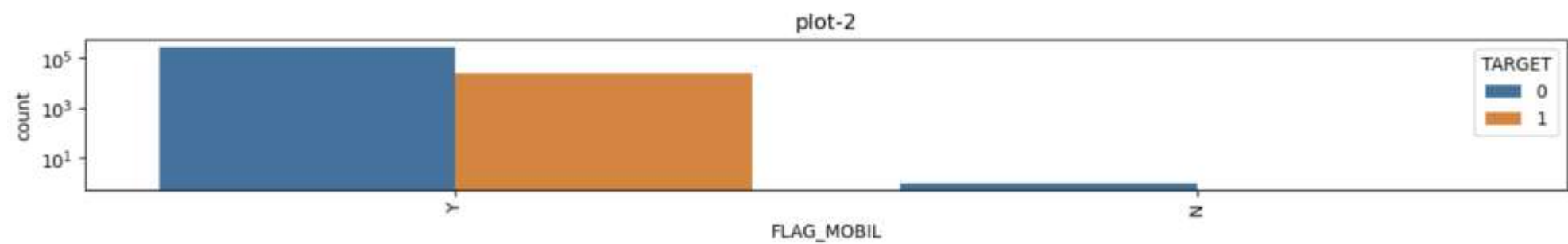
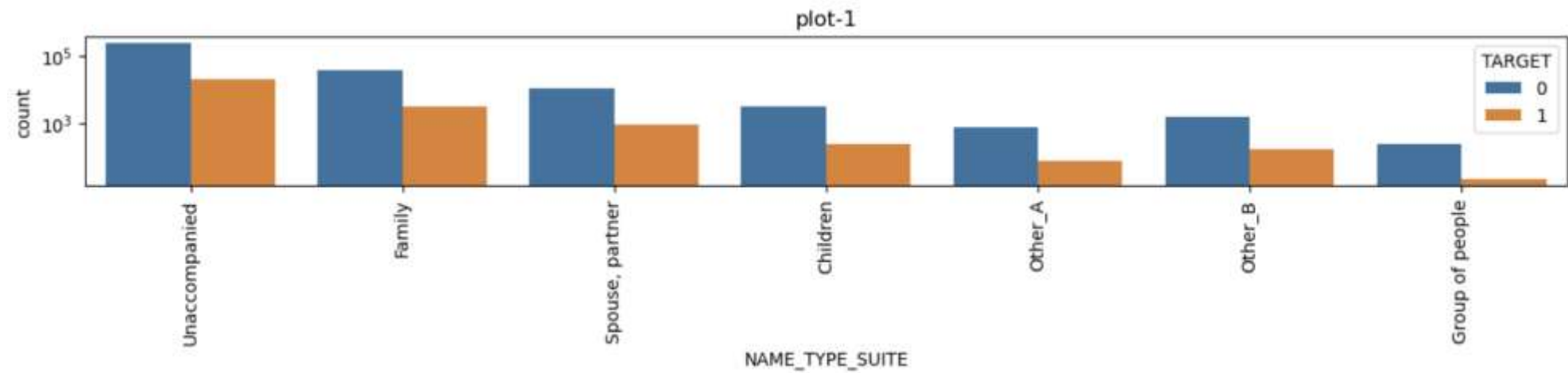
- plot-1: female clients are likely to default than male clients (XNA - is present because as instructed, I haven't handled the missing values and left them as it is)
- plot-2: married clients are likely to default than other clients
- plot-3: clients with age-group between 30-40 are more defaulters, then comes 40-50, and then 50-60
- plot-4: clients with own house/apartment are likely to default than those who stay in rented house or with their parents
- plot-5: clients with income in the range 1-2L are likely to default than clients with other income groups
- plot-6: clients with secondary special education are likely to default.





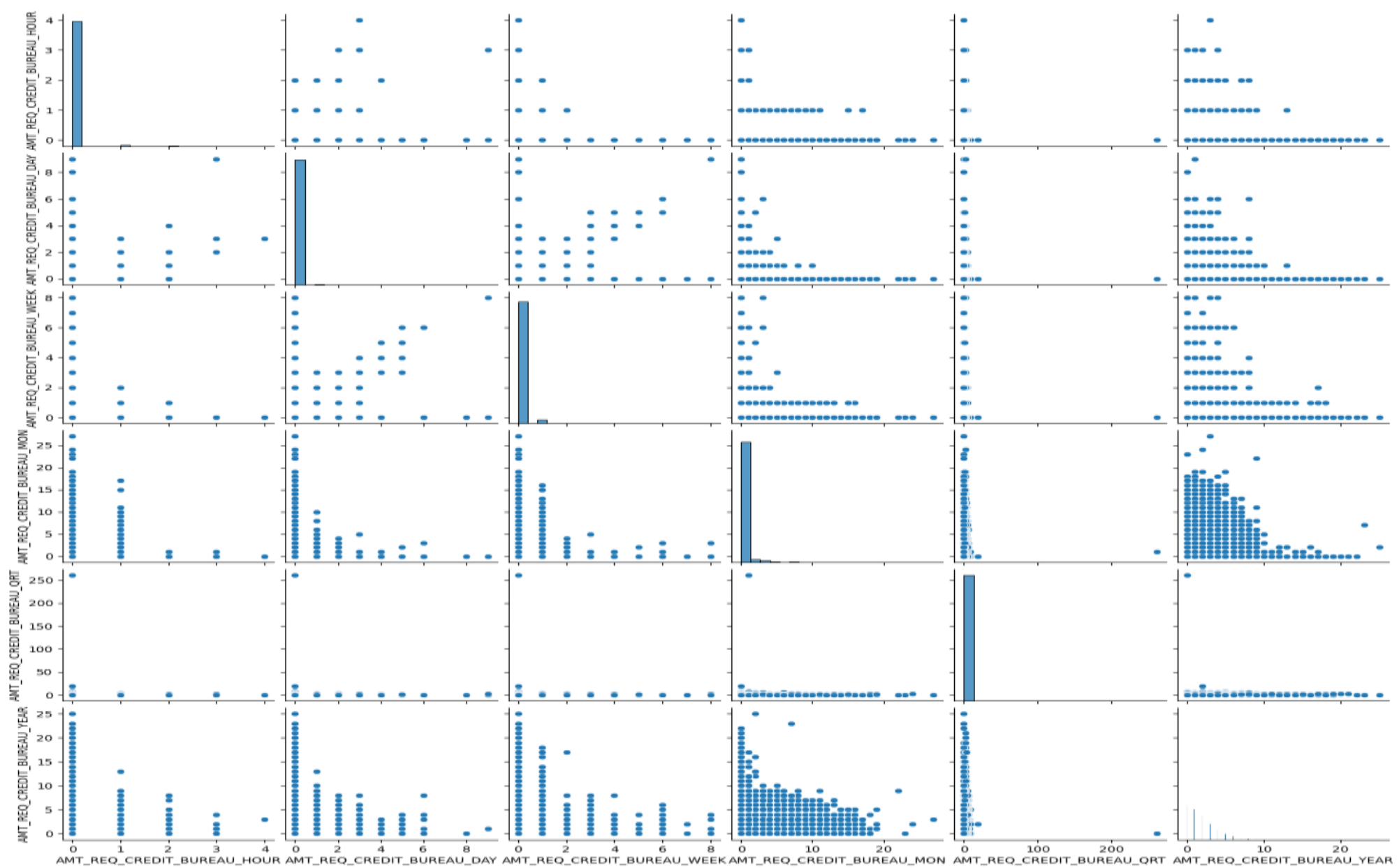
NOTE: inferences from the above plots

- plot-1: working clients are likely to default.
- plot-2: clients who haven't specified their occupation type are likely defaulters. of those who mentioned their occupation, labourers are likely defaulters.
- plot-3: clients who have business entity of type-3 are likely defaulters apart from those who haven't specified their organization type.
- plot-4: clients who have their own Real-estate property are mostly turning defaulters than those who don't own any real-estate

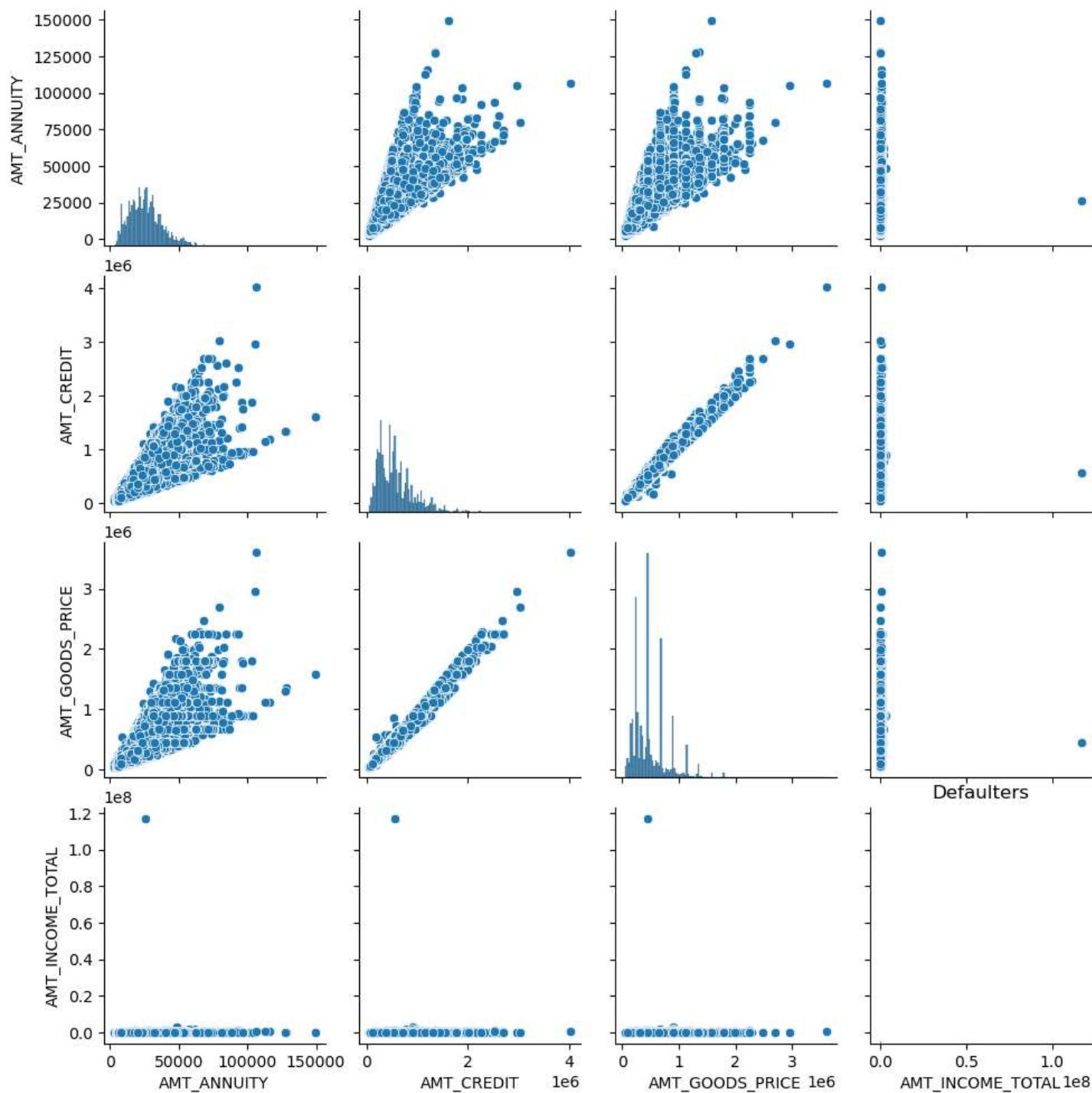


NOTE: inferences from the above plots based on contact info

- plot-1: clients who were not accompanied by anyone are likely to become defaulters.
- plot-2: clients who provided their phone number are defaulting more than those who don't.
- plot-3: clients who provided their work phone number are defaulting more than those who don't.



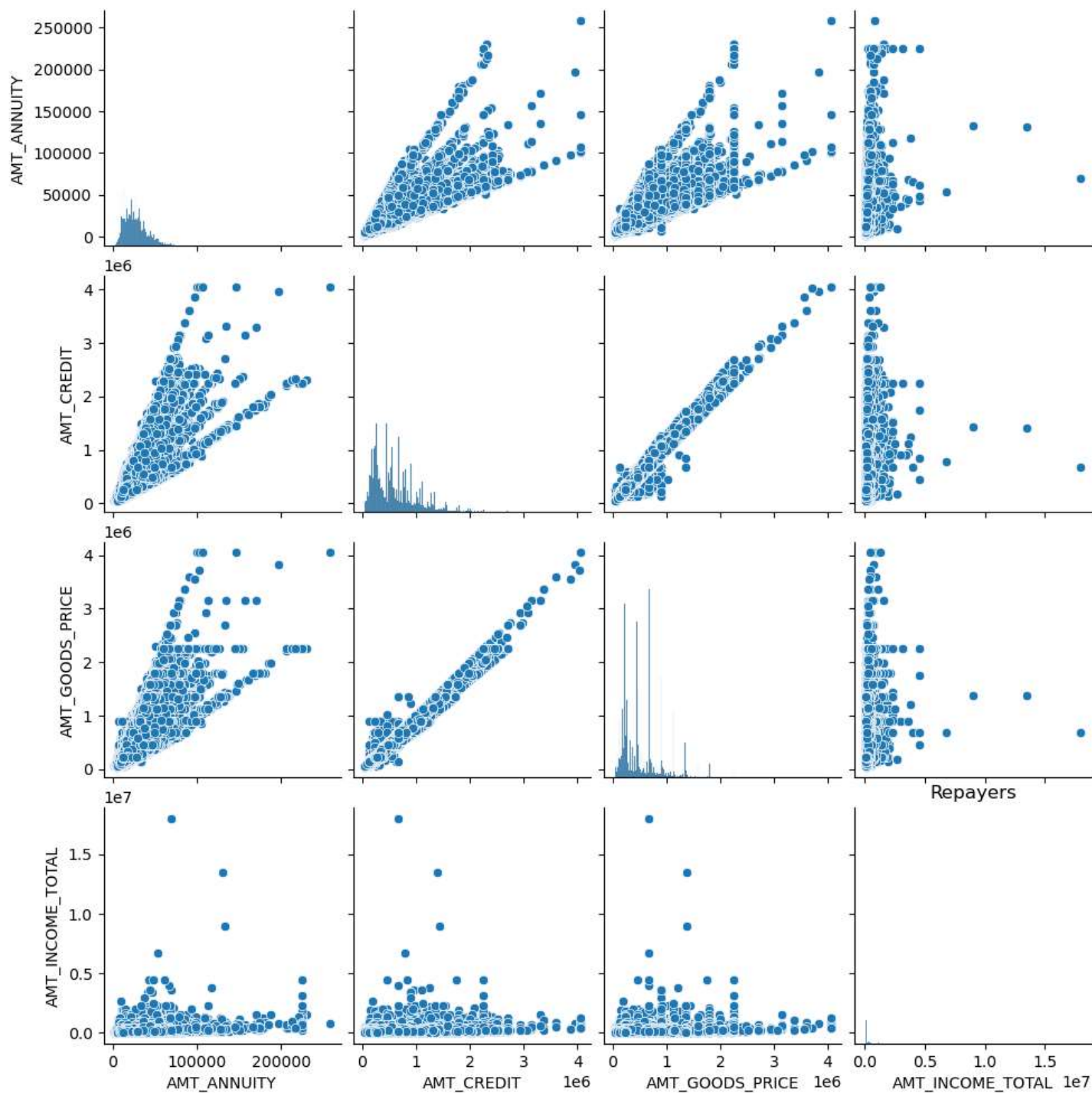
'AMT_REQ_CREDIT_BUREAU_HOUR', 'AMT_REQ_CREDIT_BUREAU_DAY', 'AMT_REQ_CREDIT_BUREAU_WEEK', 'AMT_REQ_CREDIT_BUREAU_MON', 'AMT_REQ_CREDIT_BUREAU_QRT', 'AMT_REQ_CREDIT_BUREAU_YEAR', have no proper relation between them.



for Defaulters:

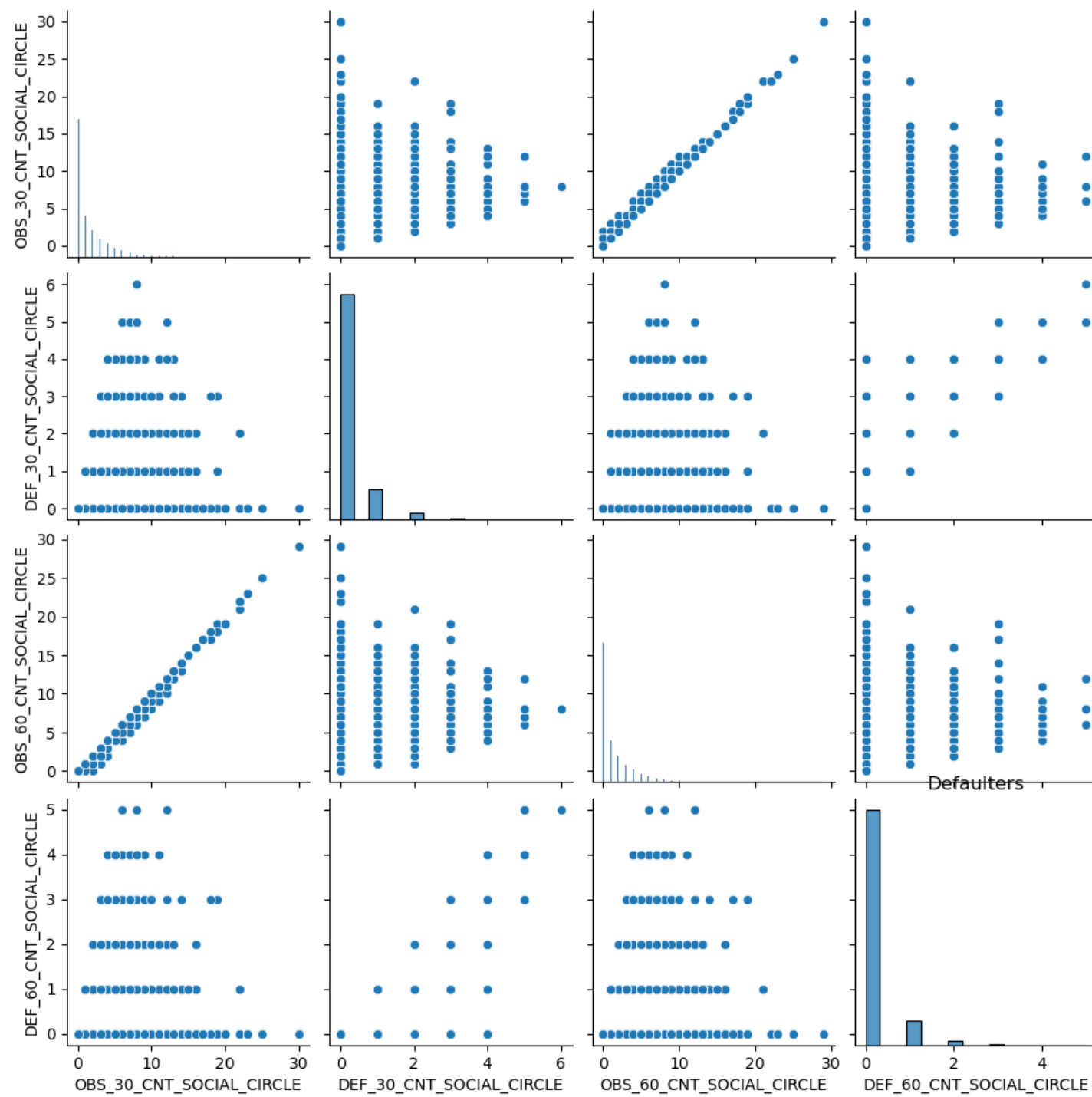
* goods_price and Credit amount have some linear relation among them

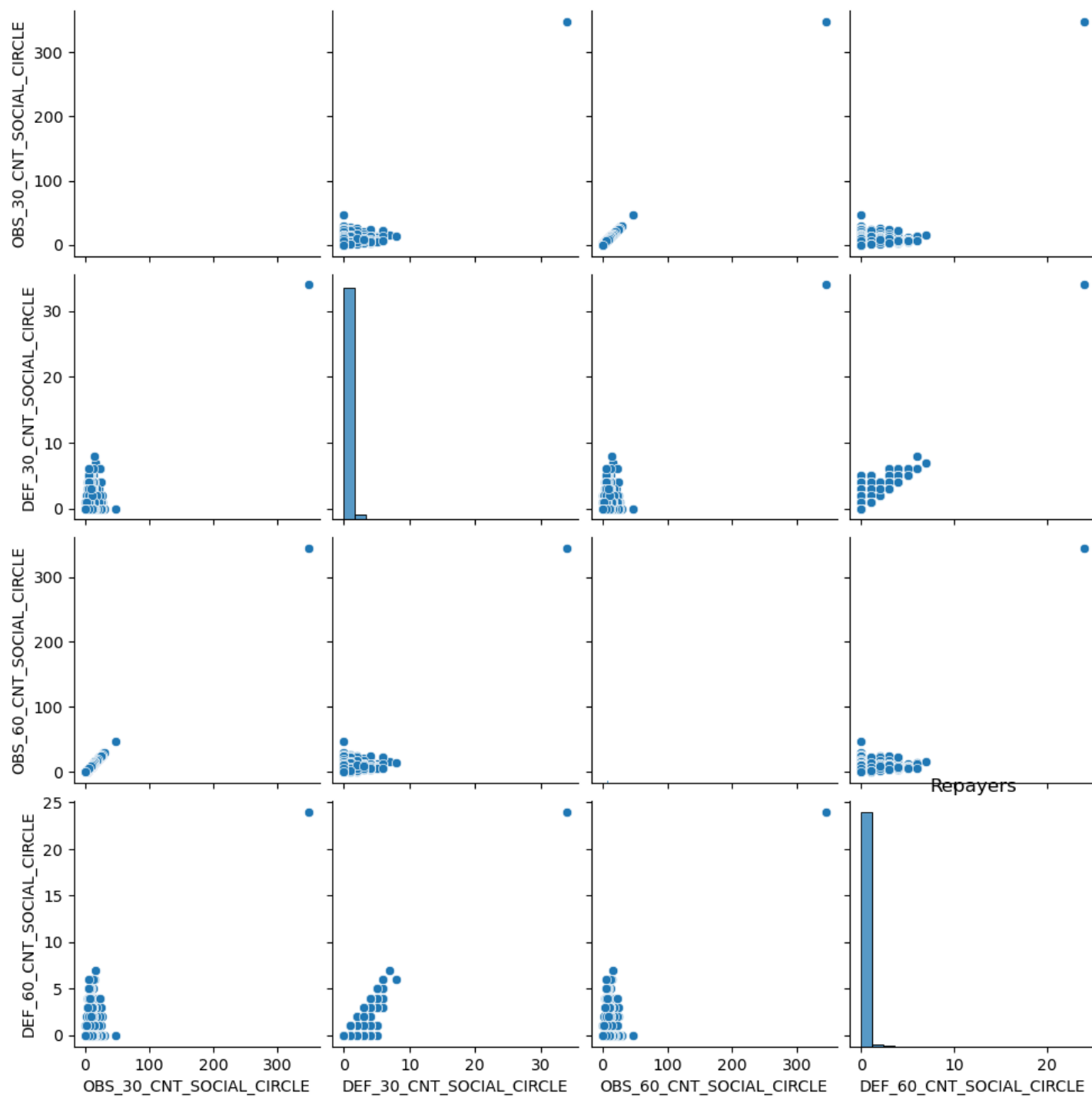
* goods-price and credit amount show some similar sort of relationship with Annuity amount



for repayers:

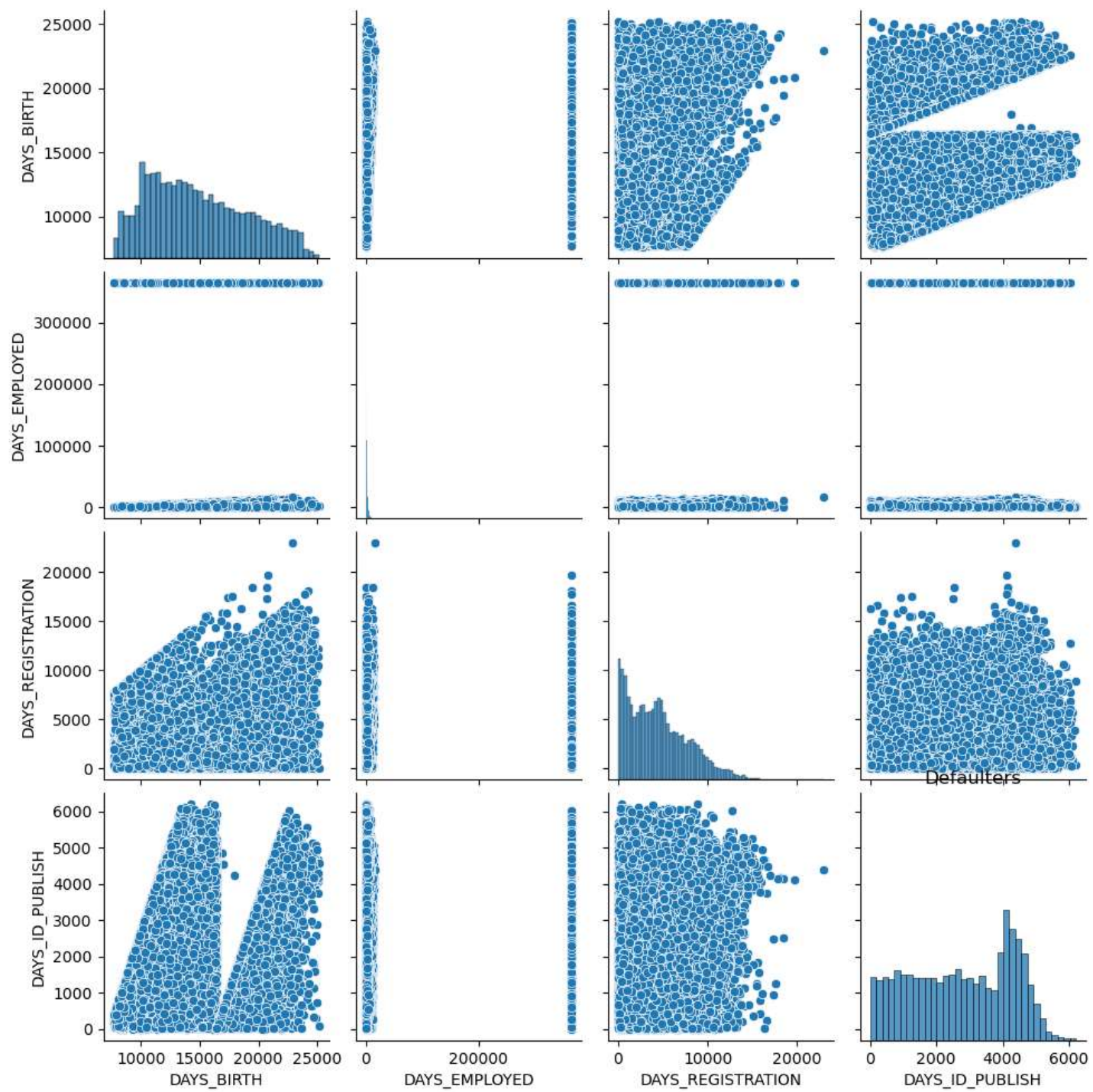
- similar kind of patterns are observed in case of both Repayers and Defaulters. but the scatter points in Repayers are more as the data is imbalanced.

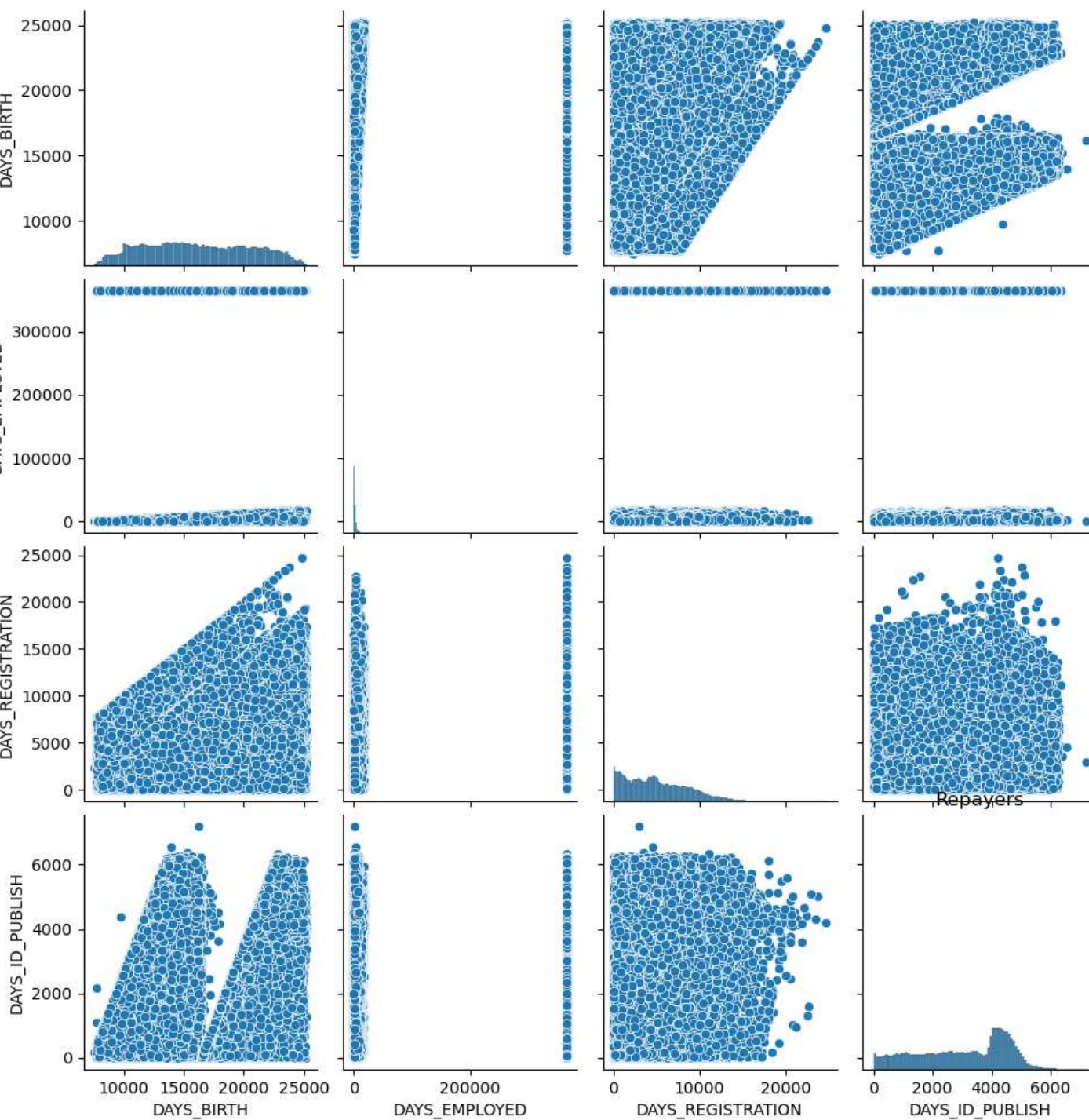




NOTE:

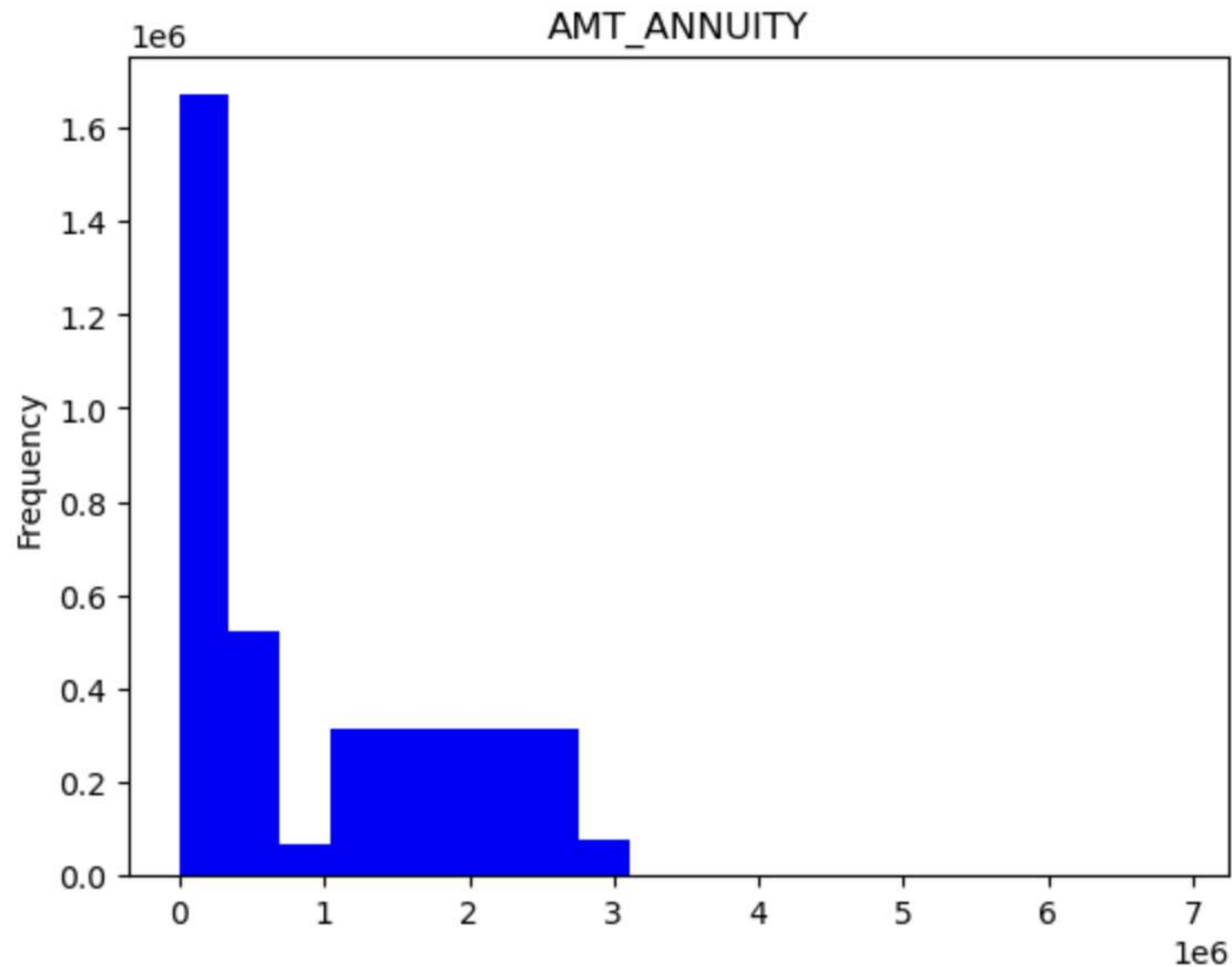
- in both the cases of Defaulters and Repayers, observable social surroundings of the client 30 DPD and 60 DPD show a linear relation.





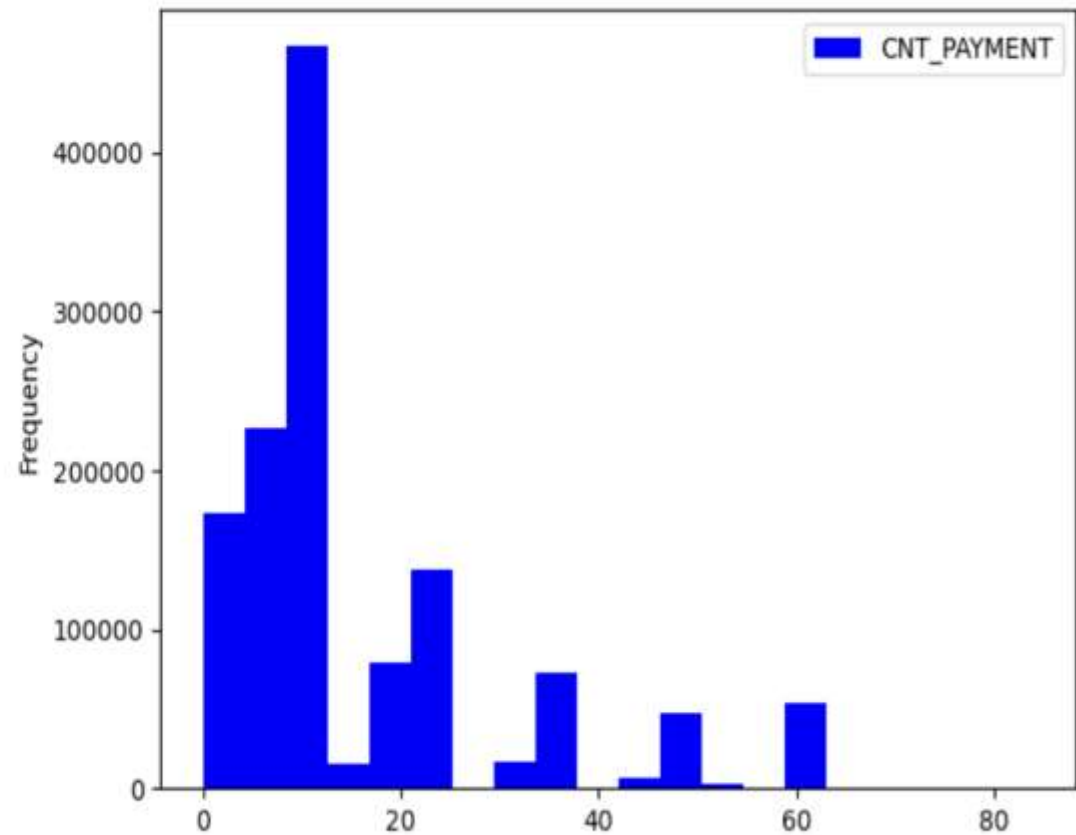
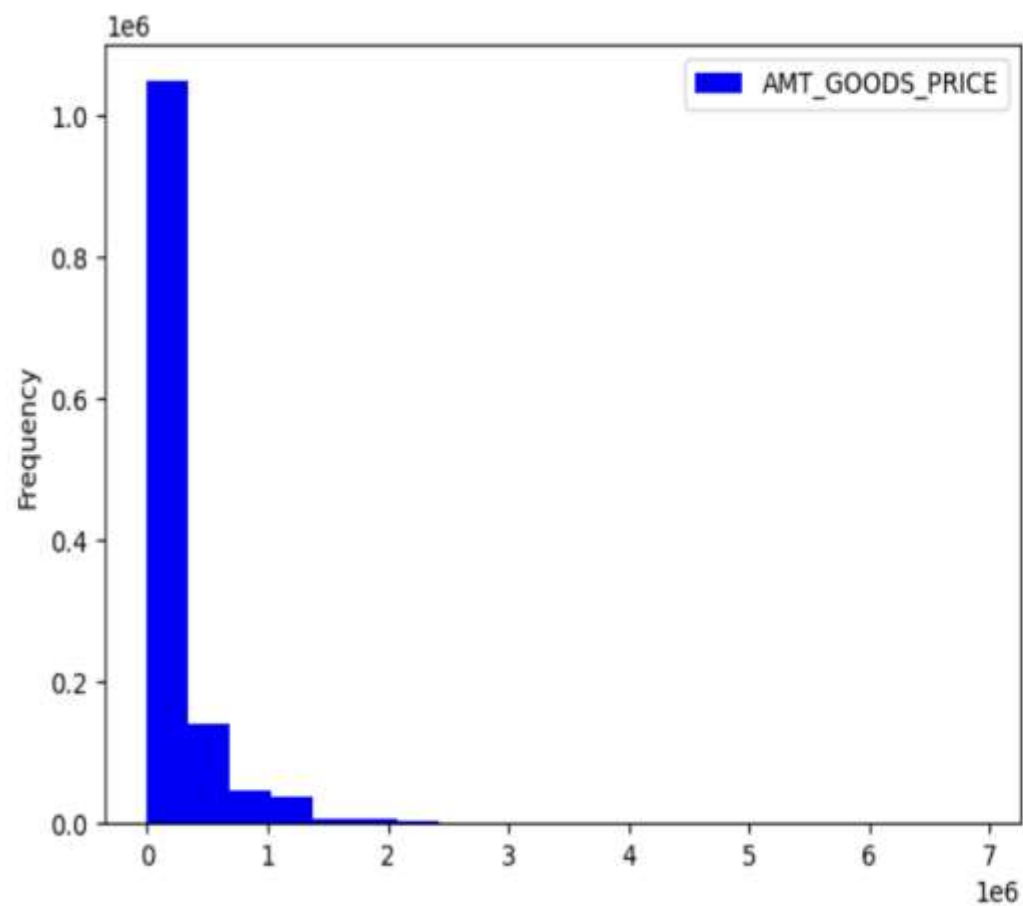
NOTE:

- from the above two pairplots for Defaulters and Repayers, it has been found that there is not specific relation between 'DAYS_BIRTH','DAYS_EMPLOYED', 'DAYS_REGISTRATION', 'DAYS_ID_PUBLISH'



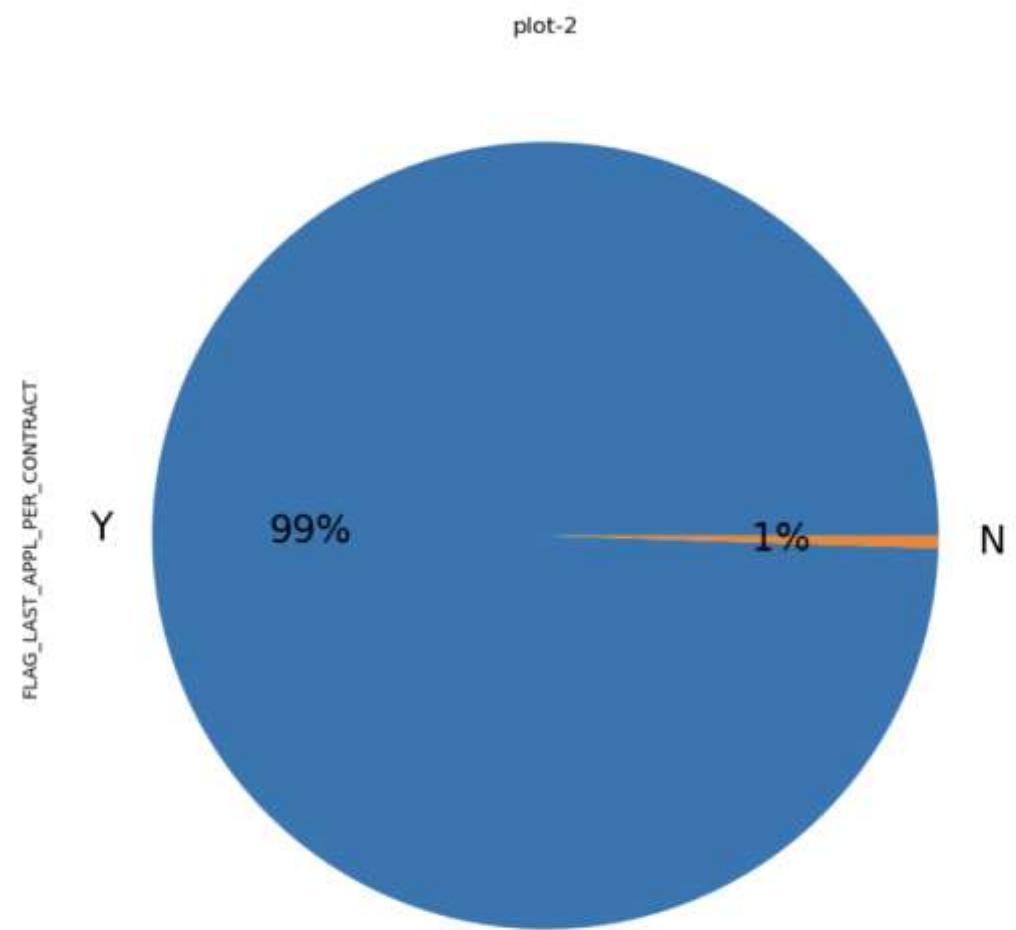
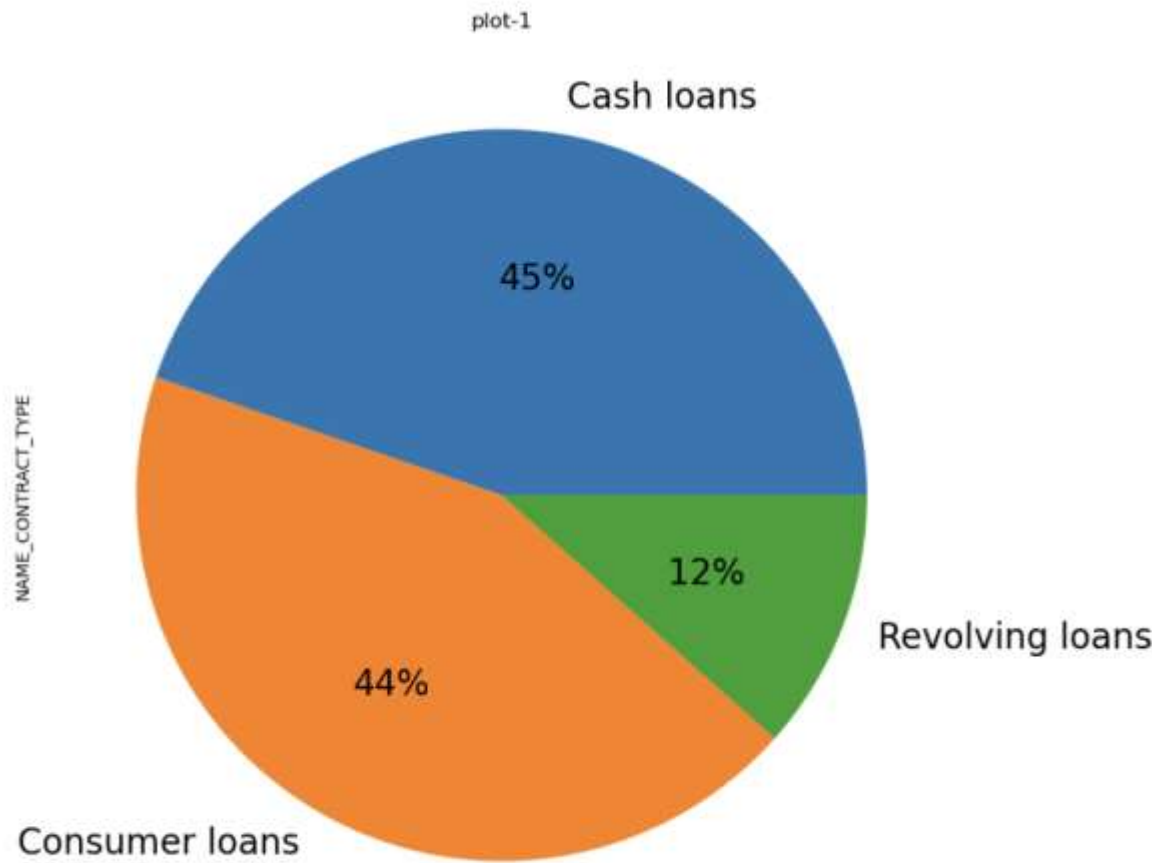
NOTE:

the hist plot of AMT_ANNUIITY doesn't show any clear trend of Gaussian curve so there is no clear way that I am able to decide how to deal with null values in AMT_ANNUIITY. hence I am also dropping this column.

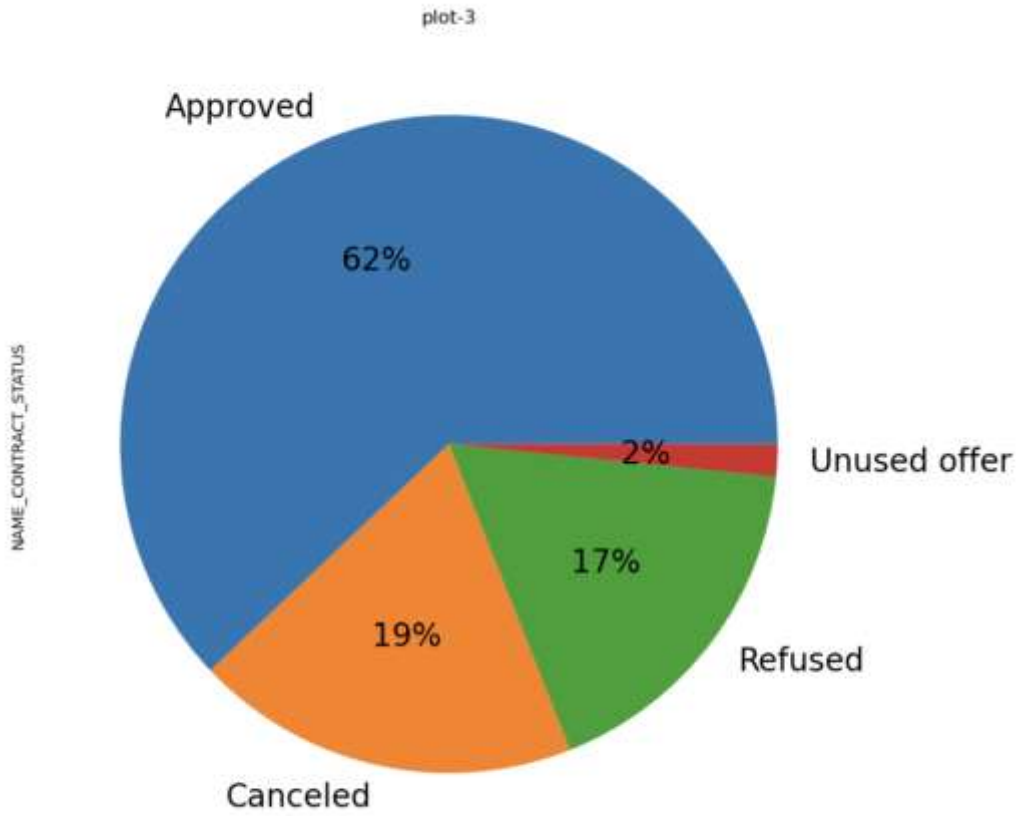


NOTE:

since a similar kind of trend is observed in case of `AMT_GOODS_PRICE` and `CNT_PAYMENT`, as in case of `AMT_ANNUITY`, we are going to drop these two columns as well

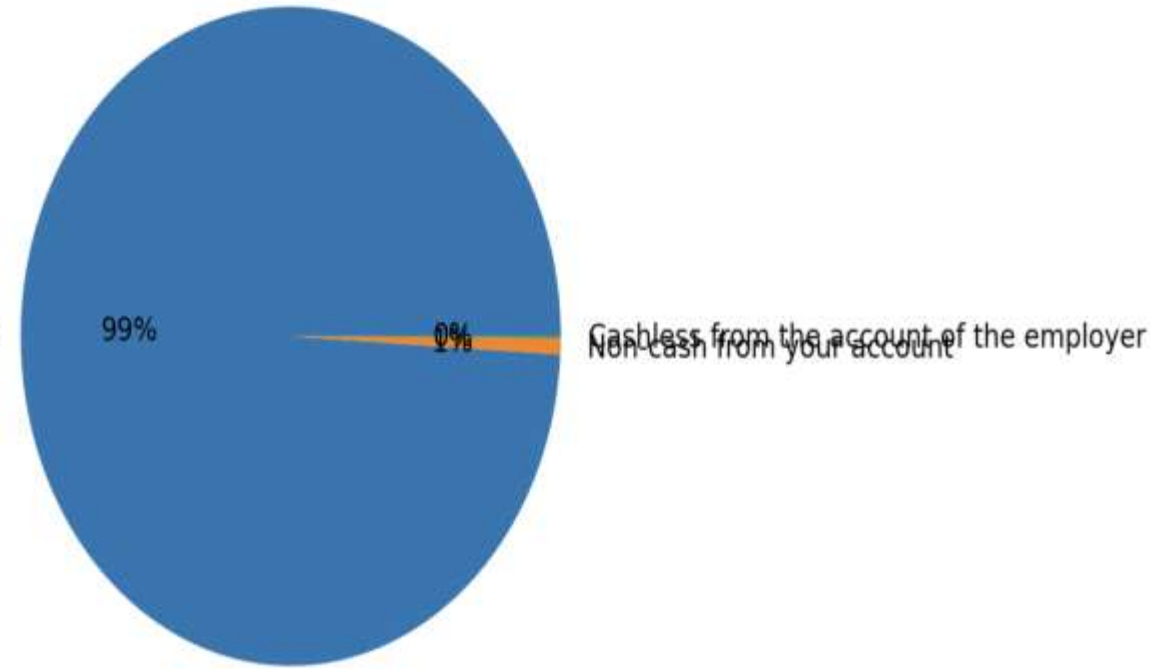


- plot-1: % of cash loans > consumer loans> revolving loans
- plot-2: FLAG_LAST_APPL_PER_CONTRACT is too unbalanced data and can be dropped



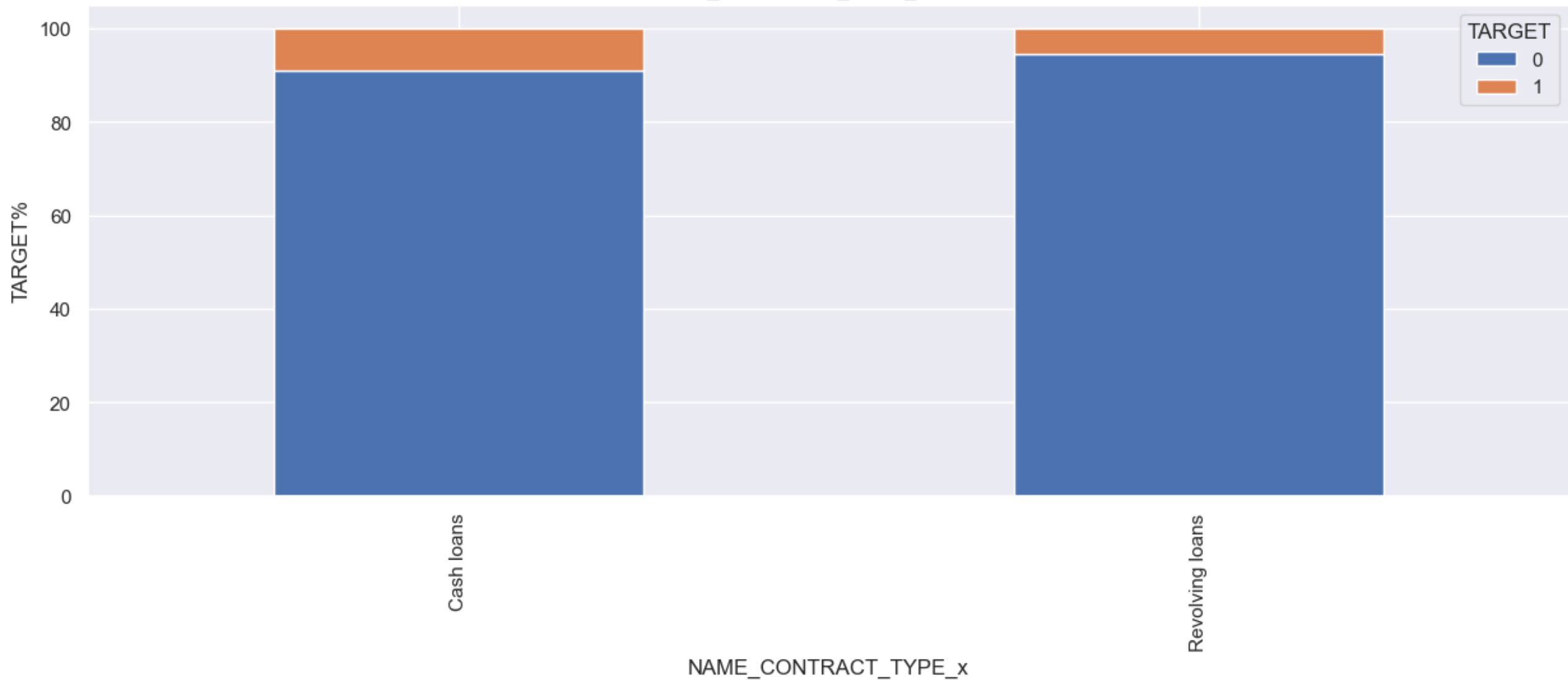
NAME_PAYMENT_TYPE

Cash through the bank

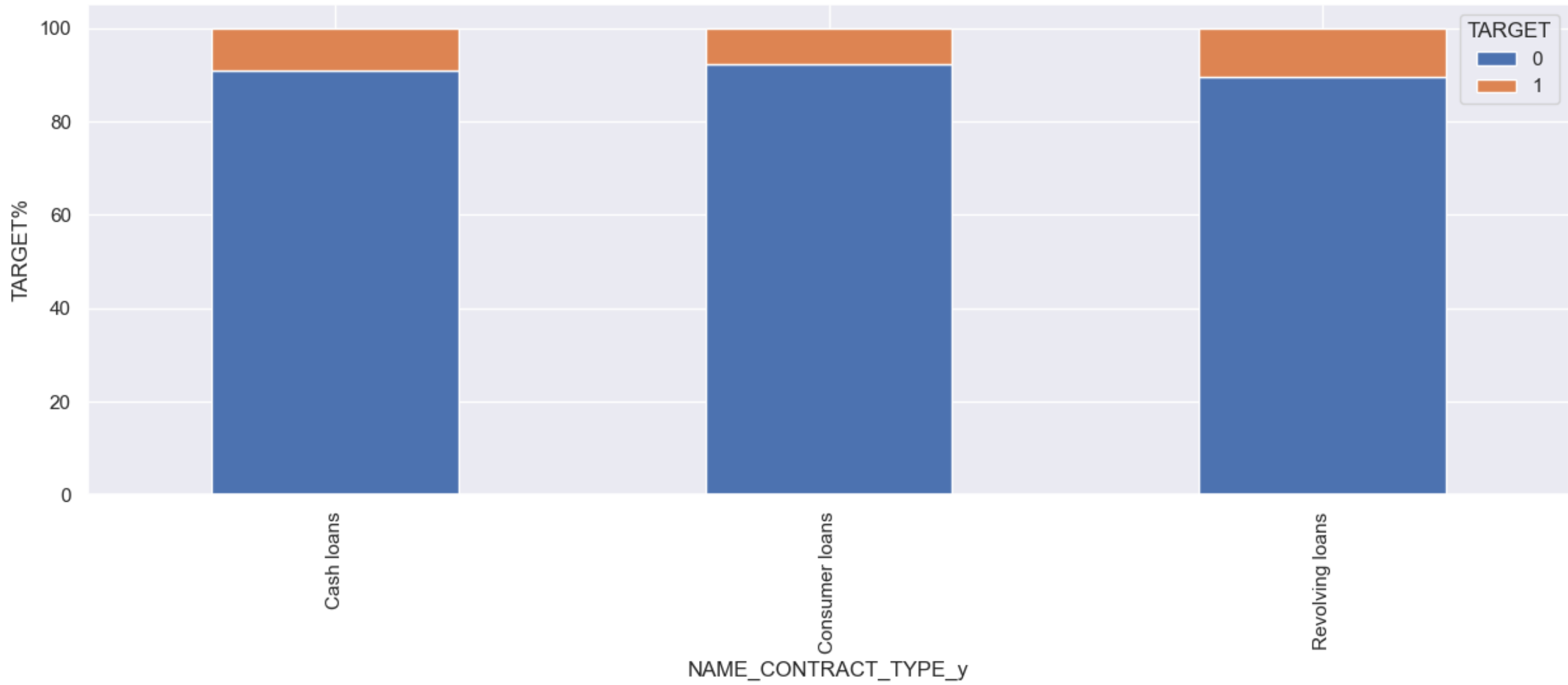


- plot-3: NAME_CONTRACT_STATUS - 62% of clients had their previous applications approved
- plot-4: NAME_PAYMENT_TYPE is also too unbalanced and can be dropped

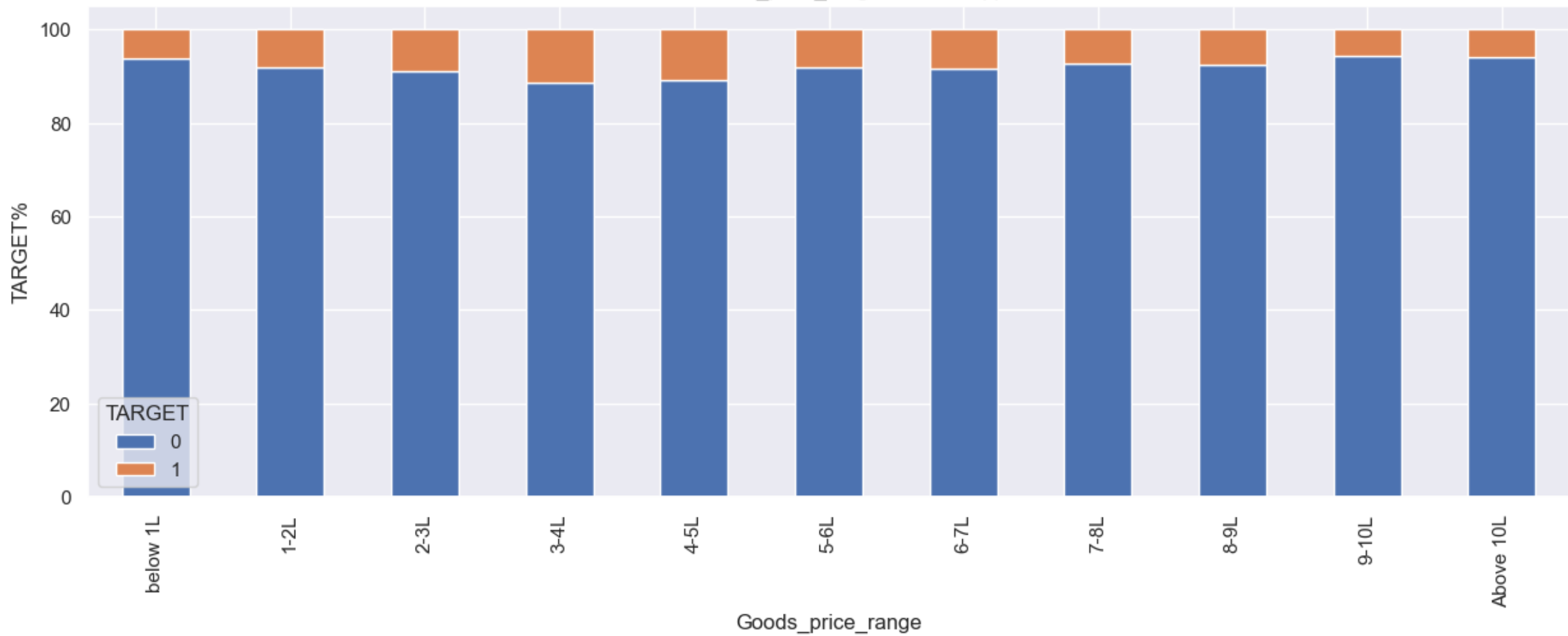
Effect Of NAME_CONTRACT_TYPE_x on Loan Approval



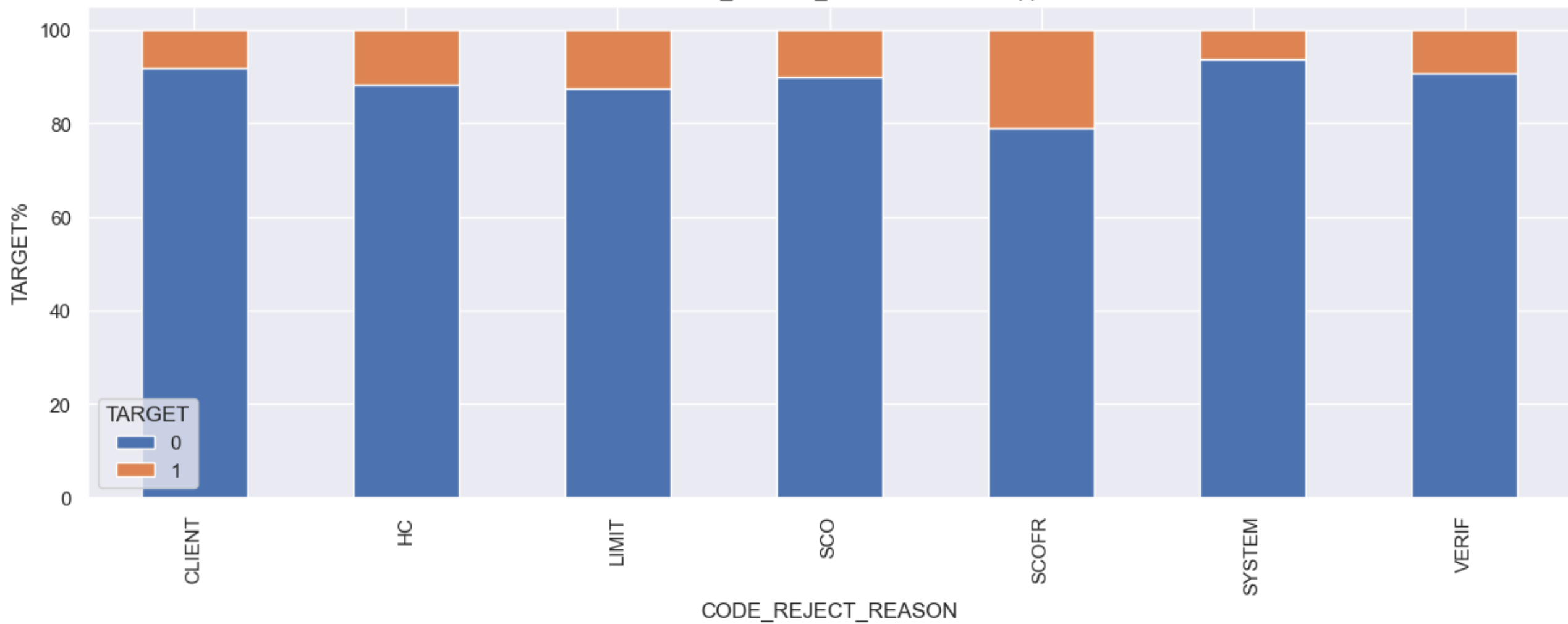
Effect Of NAME_CONTRACT_TYPE_y on Loan Approval



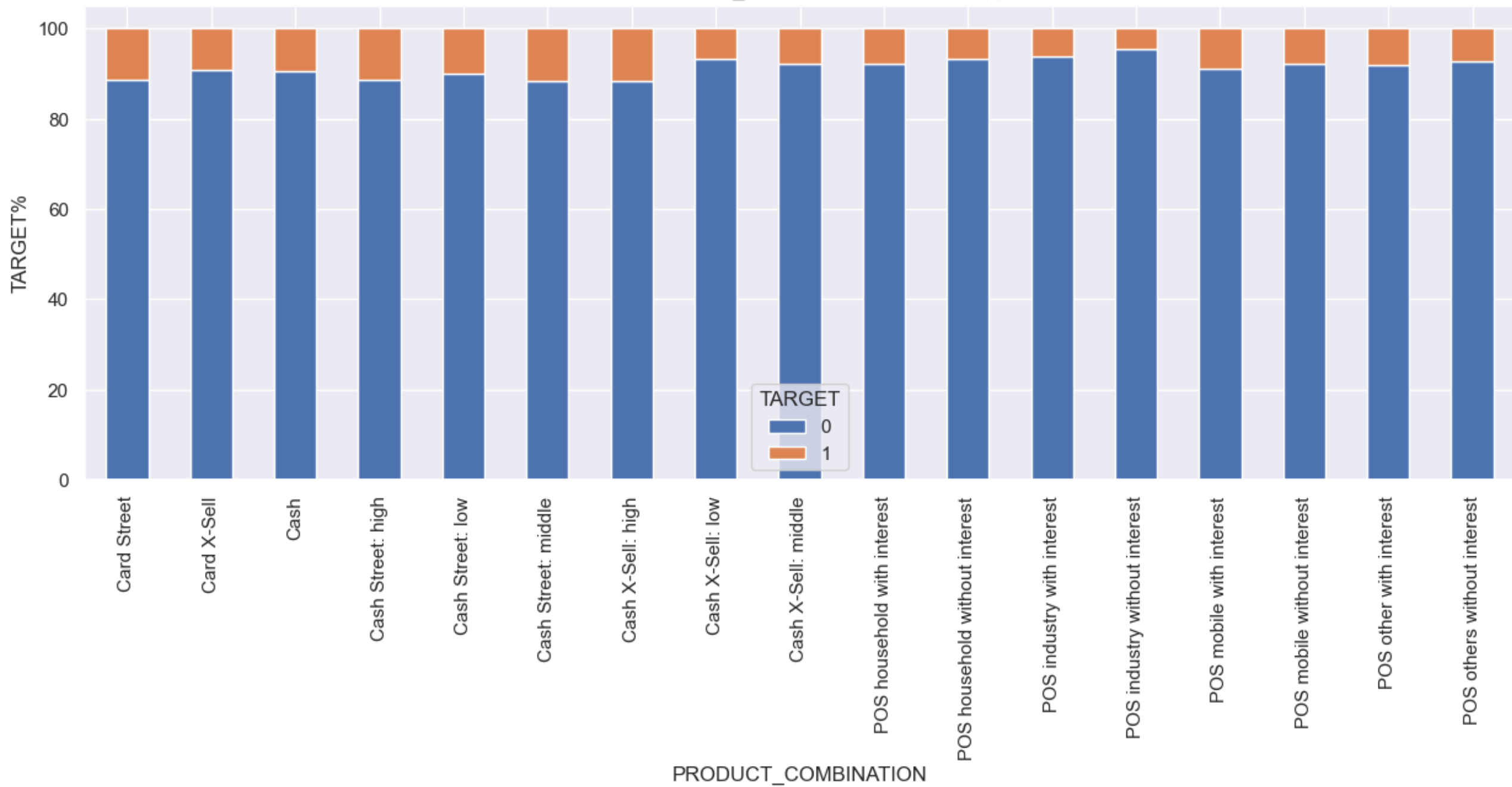
Effect Of Goods_price_range on Loan Approval



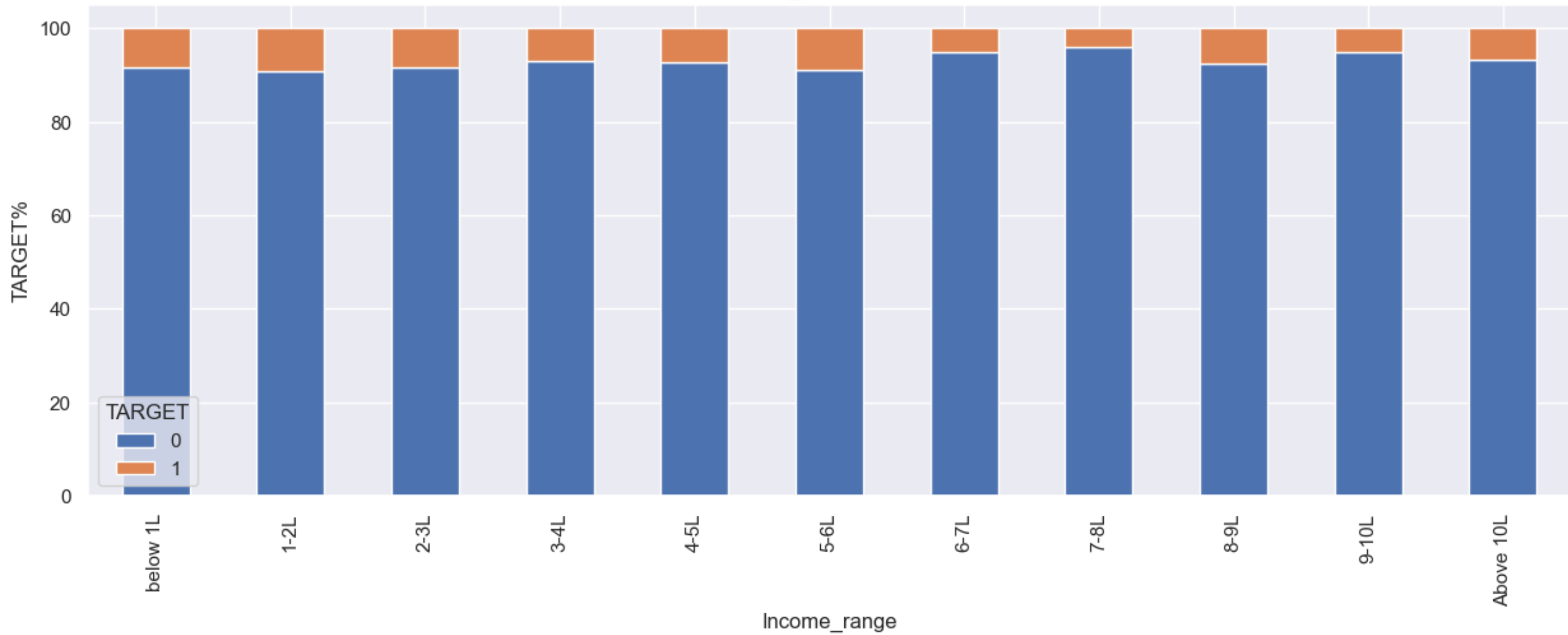
Effect Of CODE_REJECT_REASON on Loan Approval



Effect Of PRODUCT_COMBINATION on Loan Approval



Effect Of Income_range on Loan Approval



NOTE: as far as my understanding of the data set from the entire analysis is concerned, the factors that effect the loan approval for any applicant is driven by following attributes that greatly effect the target variable

- NAME_CONTRACT_TYPE_x
- NAME_CONTRACT_TYPE_y
- Income_range of client
- Goods_price_range for which the loan is being taken
- Reject reason for the previous loan
- Product combination of the client

THANK YOU