



Tanzeel Iqbal

Credit EDA Assignment

Problem Statement - I

Business Understanding

The loan providing companies find it hard to give loans to the people due to their insufficient or non-existent credit history. Because of that, some consumers use it to their advantage by becoming a defaulter. Suppose you work for a consumer finance company which specialises in lending various types of loans to urban customers. You have to use EDA to analyse the patterns present in the data. This will ensure that the applicants capable of repaying the loan are not rejected.

When the company receives a loan application, the company has to decide for loan approval based on the applicant's profile. Two types of risks are 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.

Problem Statement - II

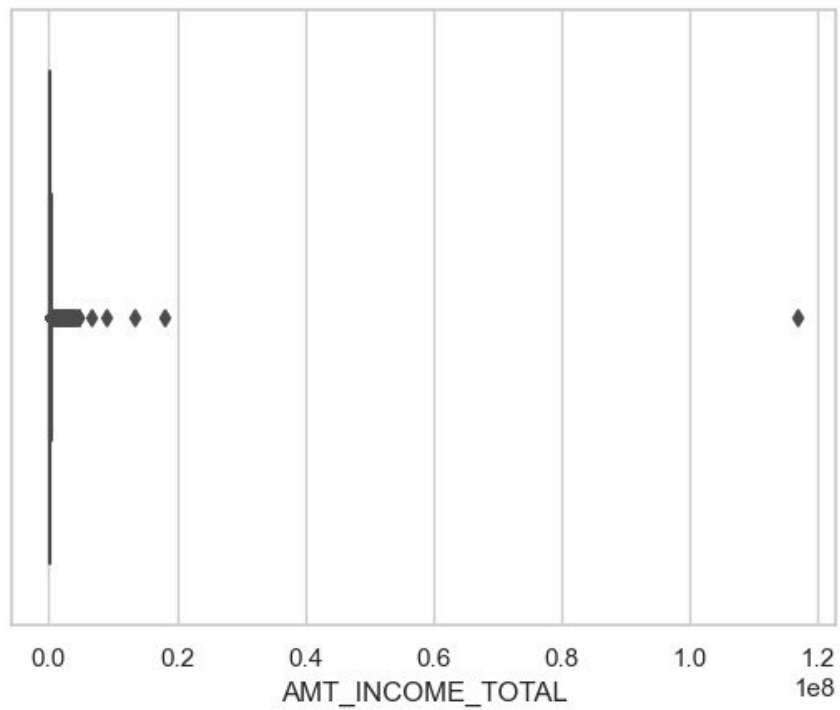
Present the overall approach of the analysis in a presentation. Mention the problem statement and the analysis approach briefly.

Identify the missing data and use appropriate method to deal with it. (Remove columns/or replace it with an appropriate value)

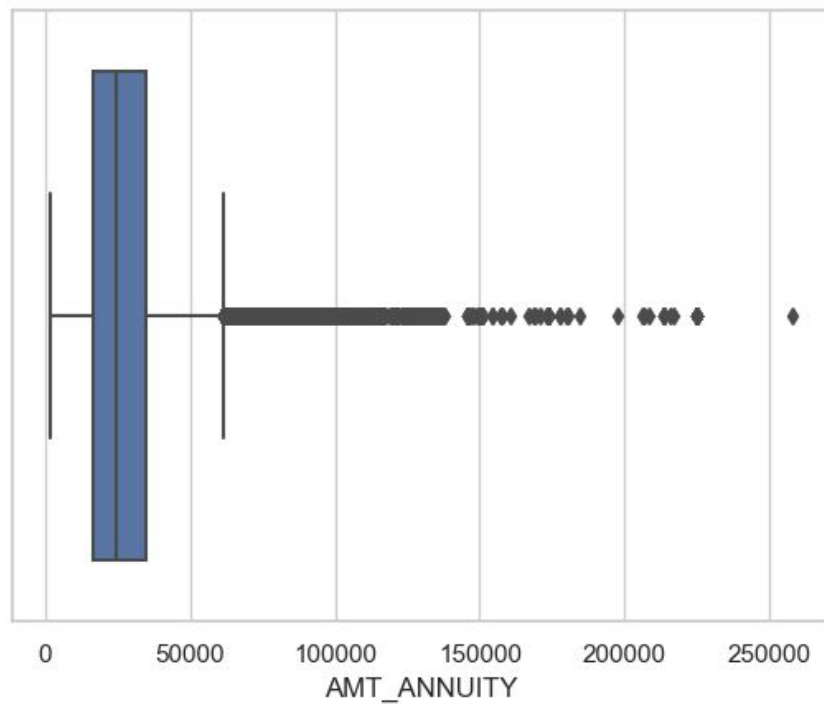
Outlier Analysis

These are highly deviated values from our dense region of data. These need to be analysed and dropped if needed.

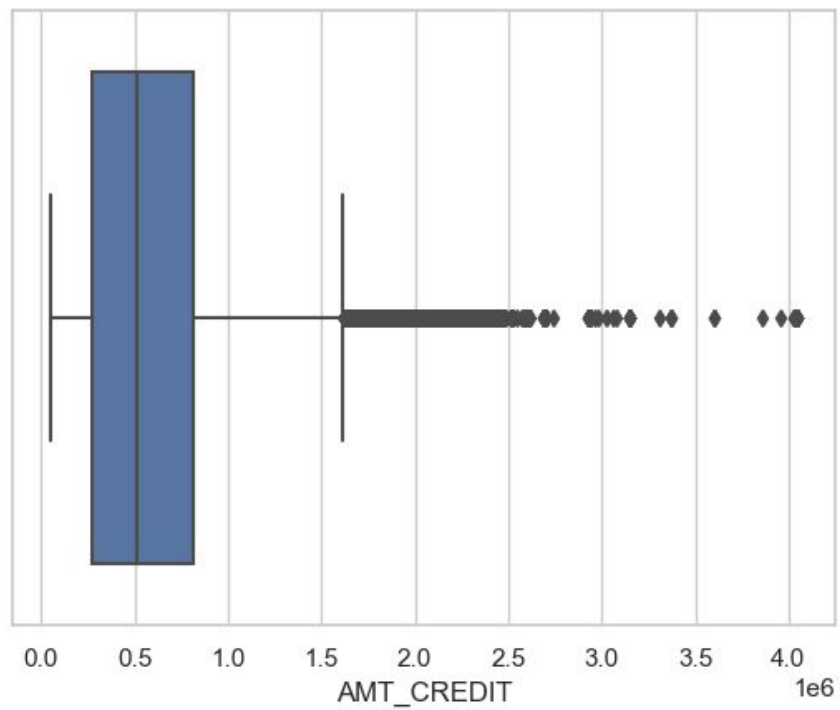
#PS Everything above 1M is an outlier



#PS values above 69898 seems to be outliers



#PS values above 1843771 seems to be outliers



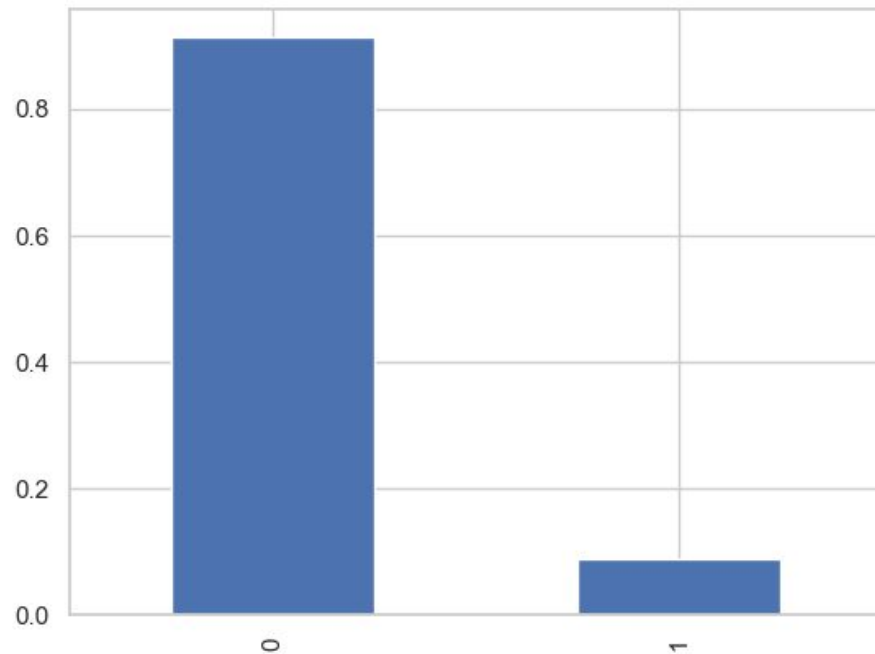
Target Analysis

Target = 1 are clients who are having difficulty repaying

Target = 0 are the clients who have been paying on time

```
#PS less than 10% people have difficulty paying back
```

```
]: <AxesSubplot:>
```

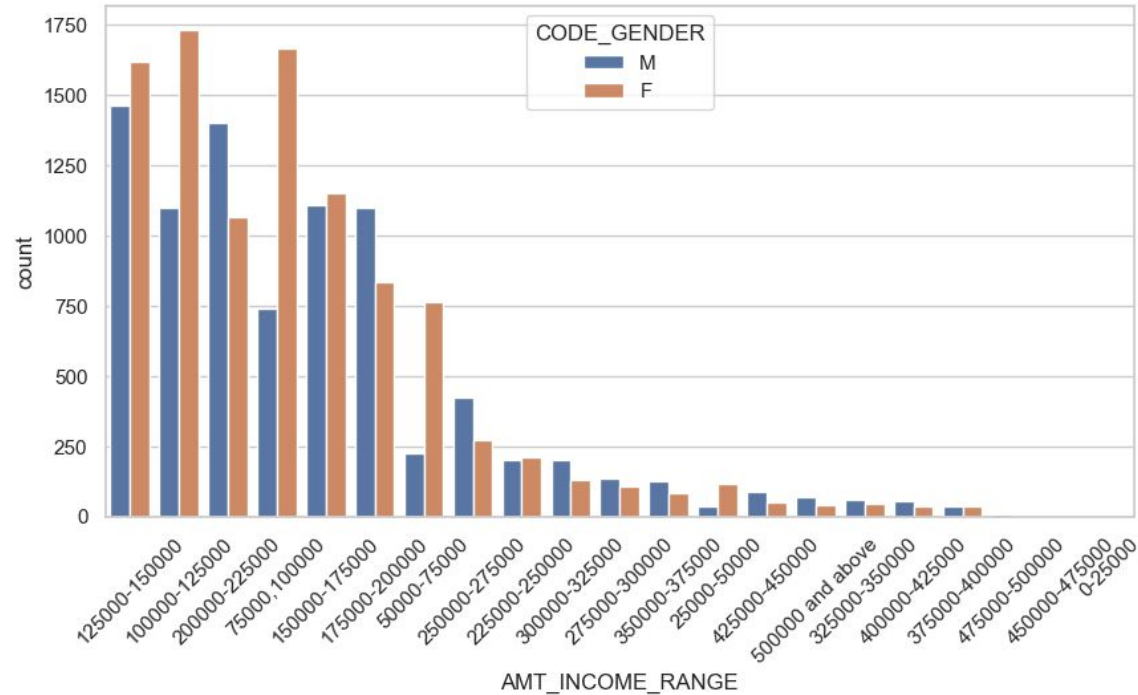


Deep dive into Target = 1

Analyse data and figure out
patterns for our target clients

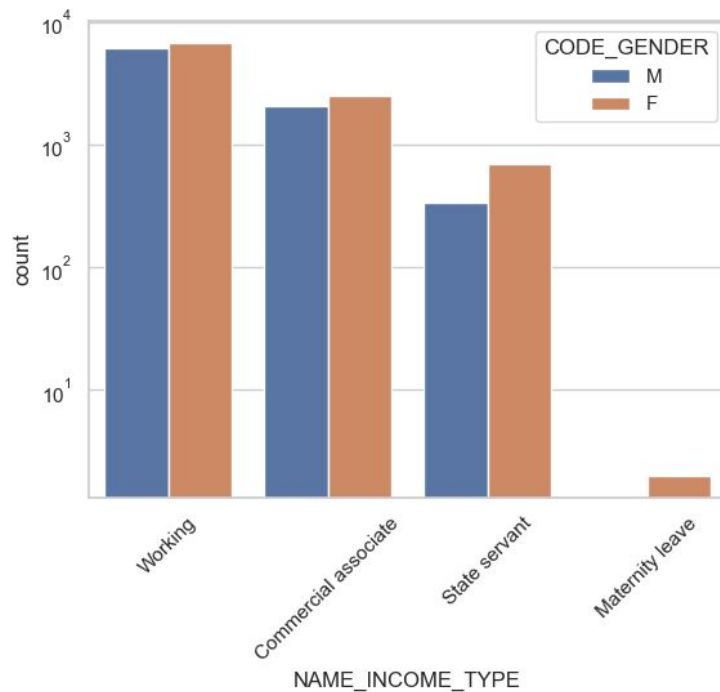
#Plotting AMT_INCOME_RANGE with CODE_GENDER

#The following plot shows that female count of default is higher than male



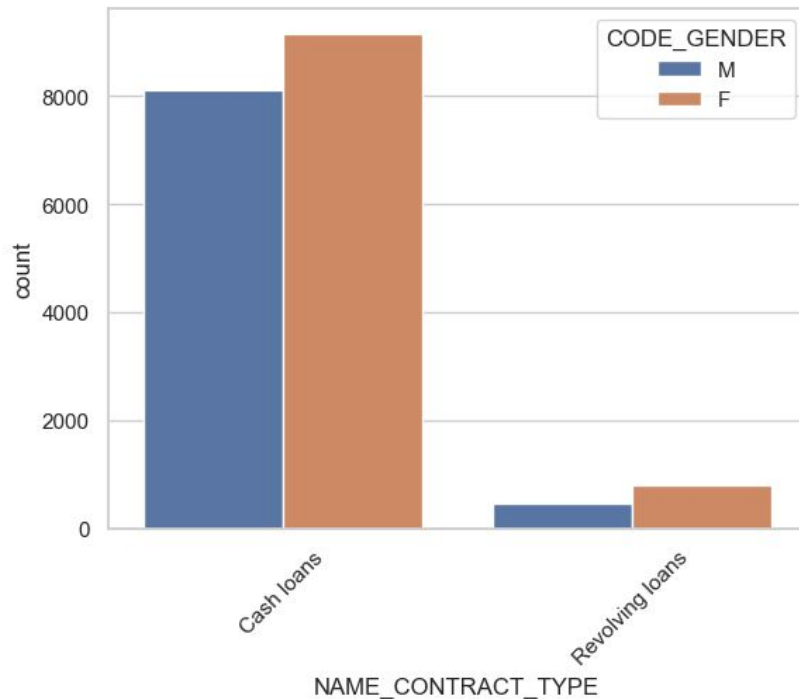
#Plotting NAME_INCOME_TYPE with CODE_GENDER

```
#The following plot shows that working people are most of the clients  
#There are more female clients in the top 3 categories of income type
```

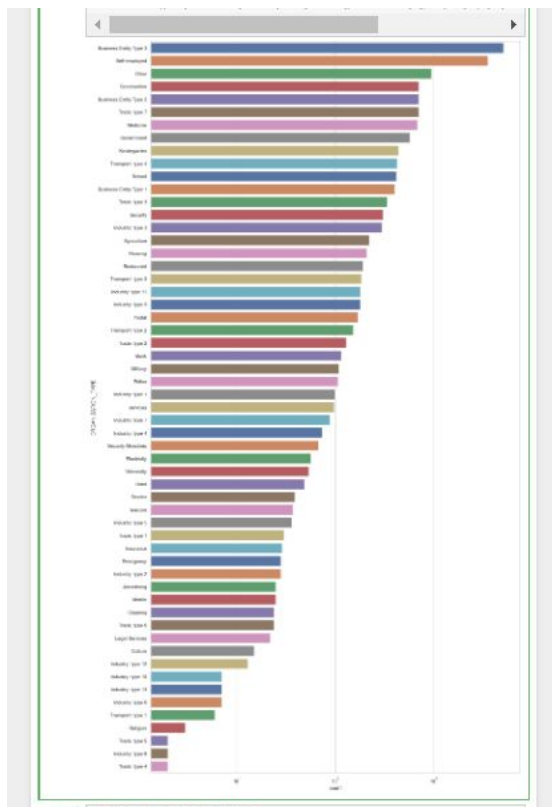


#Plotting NAME_CONTRACT_TYPE with CODE_GENDER

```
#Cash loans clearly has more clients per credit compared to revolving loans. In both cases,
```

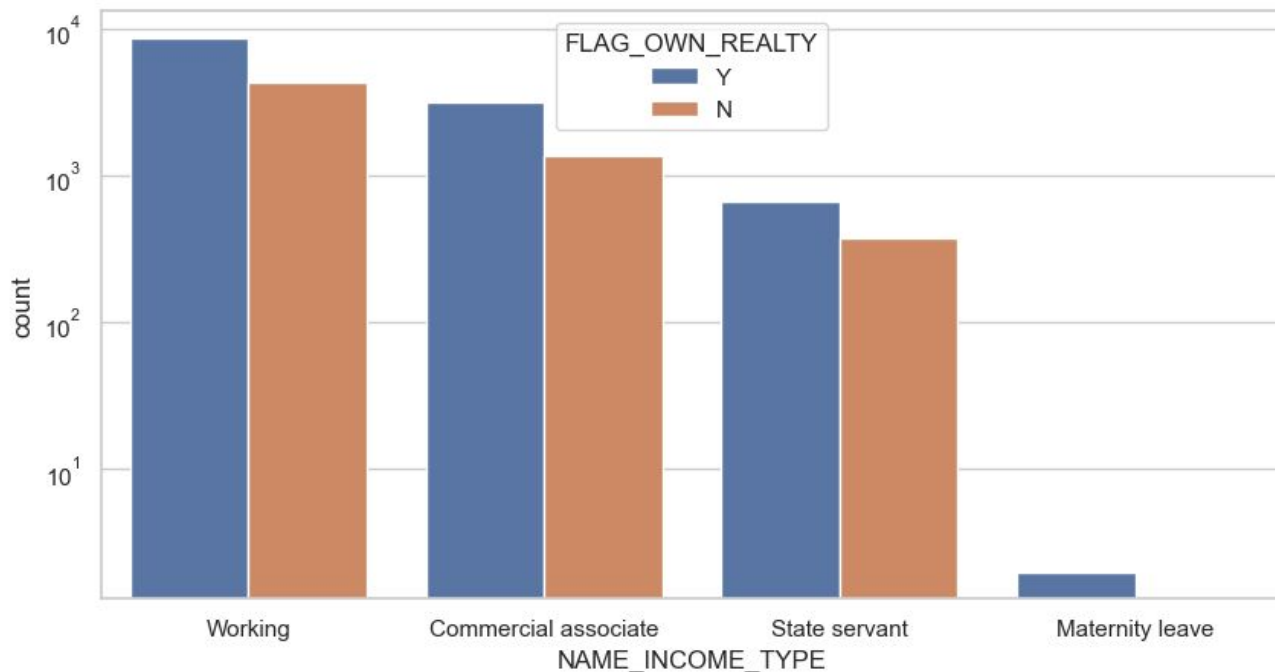


Plotting for ORGANIZATION_TYPE



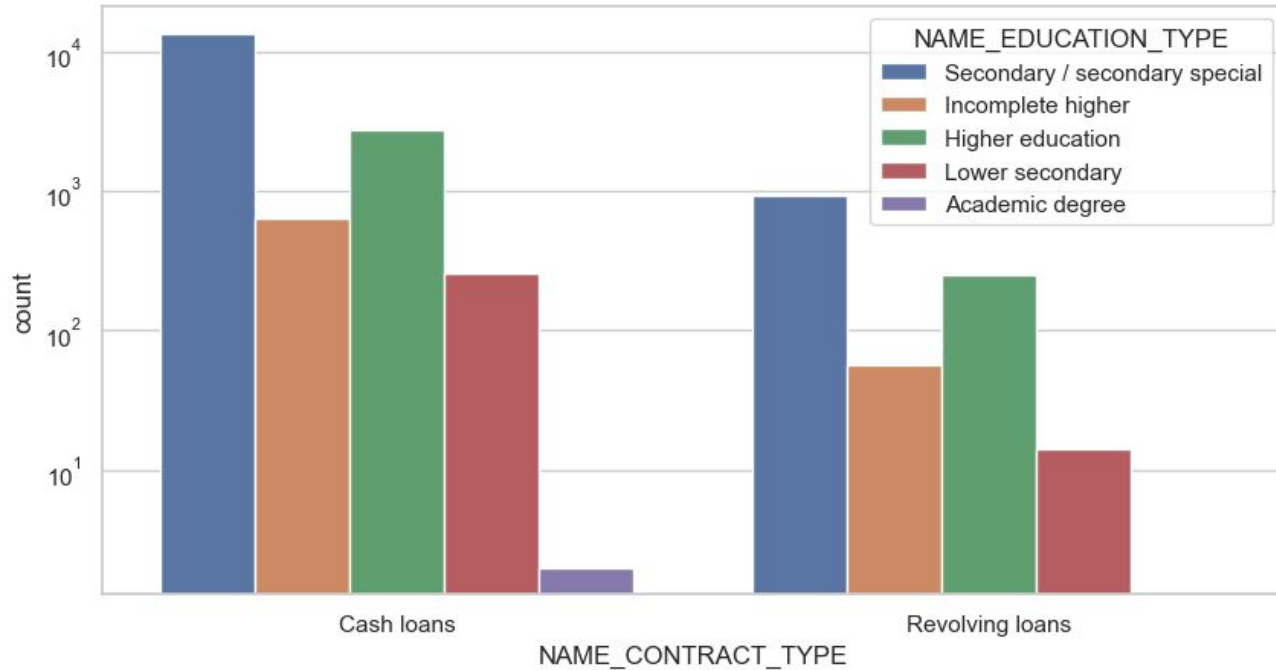
#Plotting NAME_INCOME_TYPE with FLAG_OWN_REALTY

```
#Working customers, obviously, have a higher count.  
#As we can see, most customers do have their own property (house or a flat) but a large number of customers can be stat
```



#Plotting for NAME_CONTRACT_TYPE with NAME_EDUCATION_TYPE

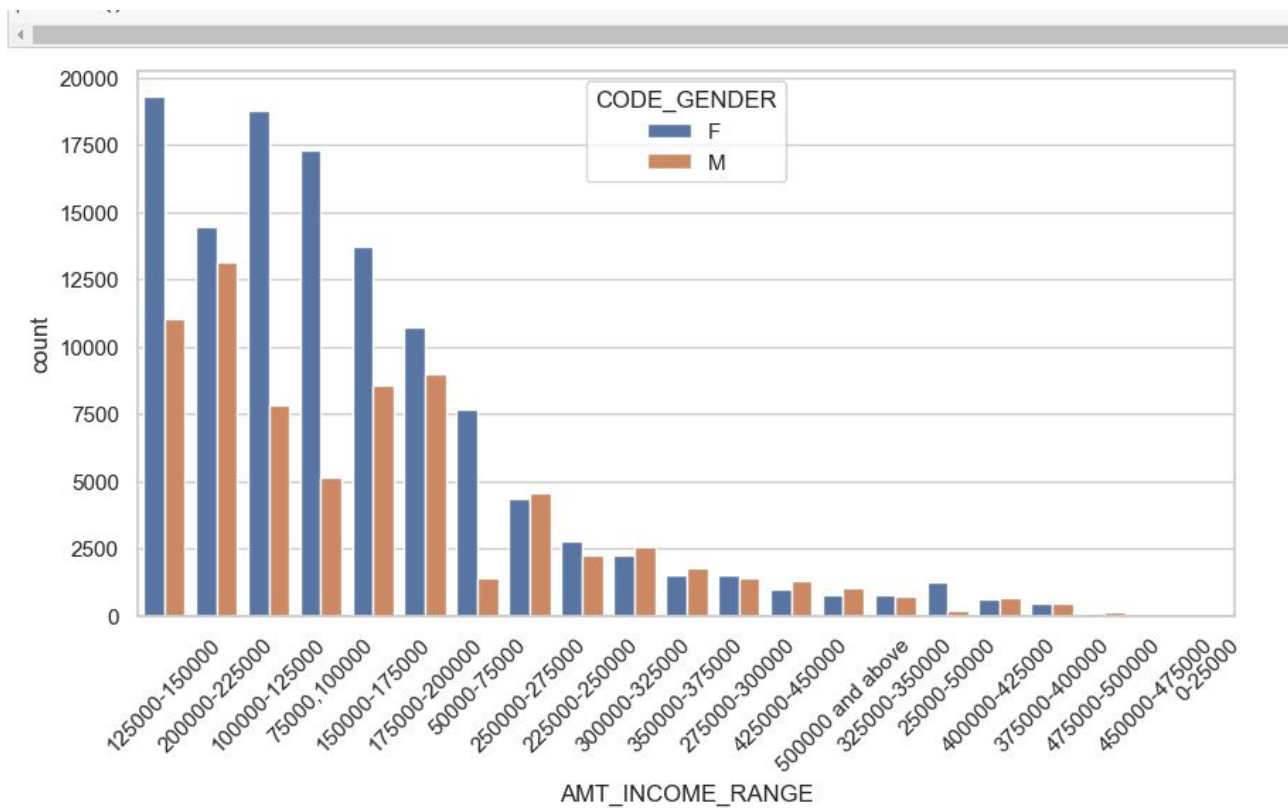
*#Cash loans, as we can see, are preferred by clients of all education backgrounds with a majority.
#People with only an academic degree do not prefer revolving loans at all.*



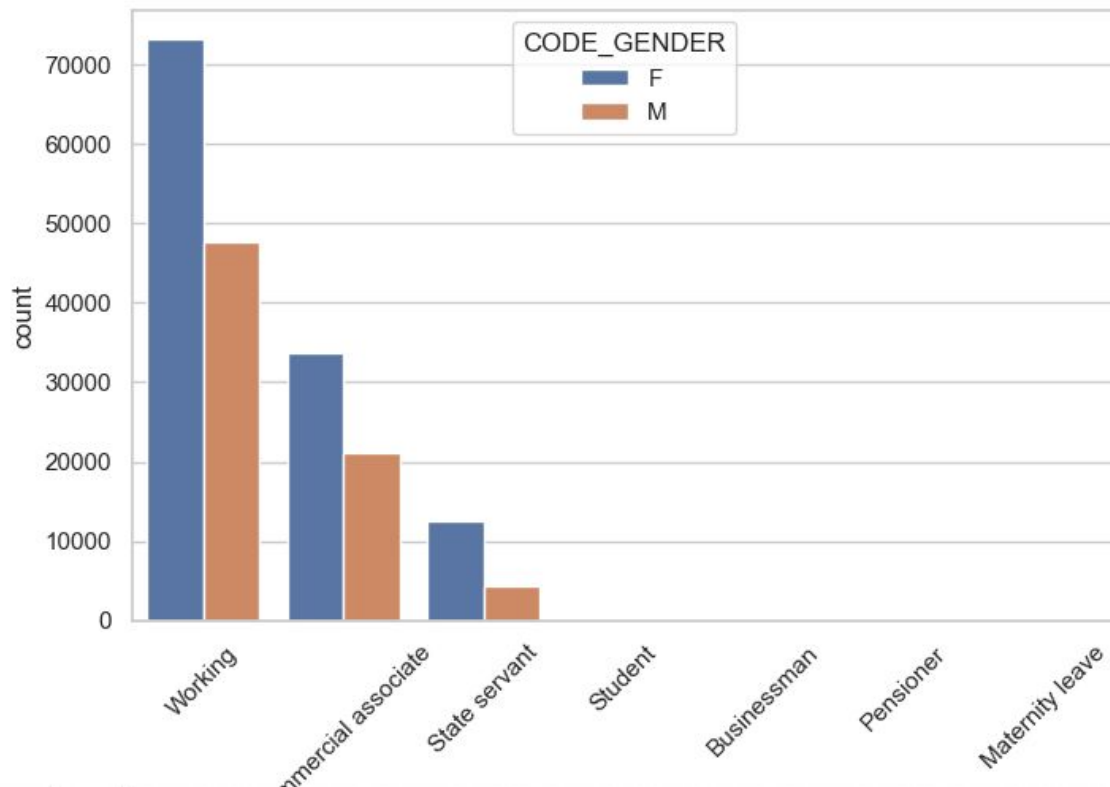
Deep dive into Target = 0

Analyse data and figure out
patterns for our clients who have
been paying on time

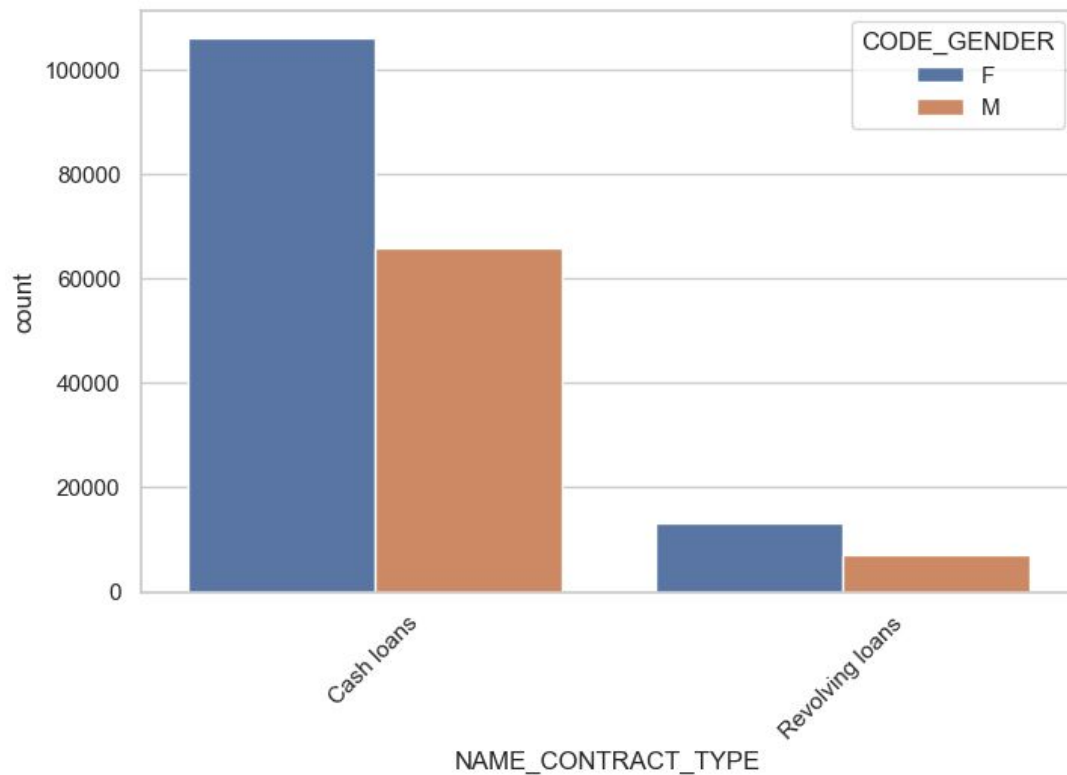
#Plotting AMT_INCOME_RANGE with gender



#Plotting NAME_INCOME_TYPE with gender



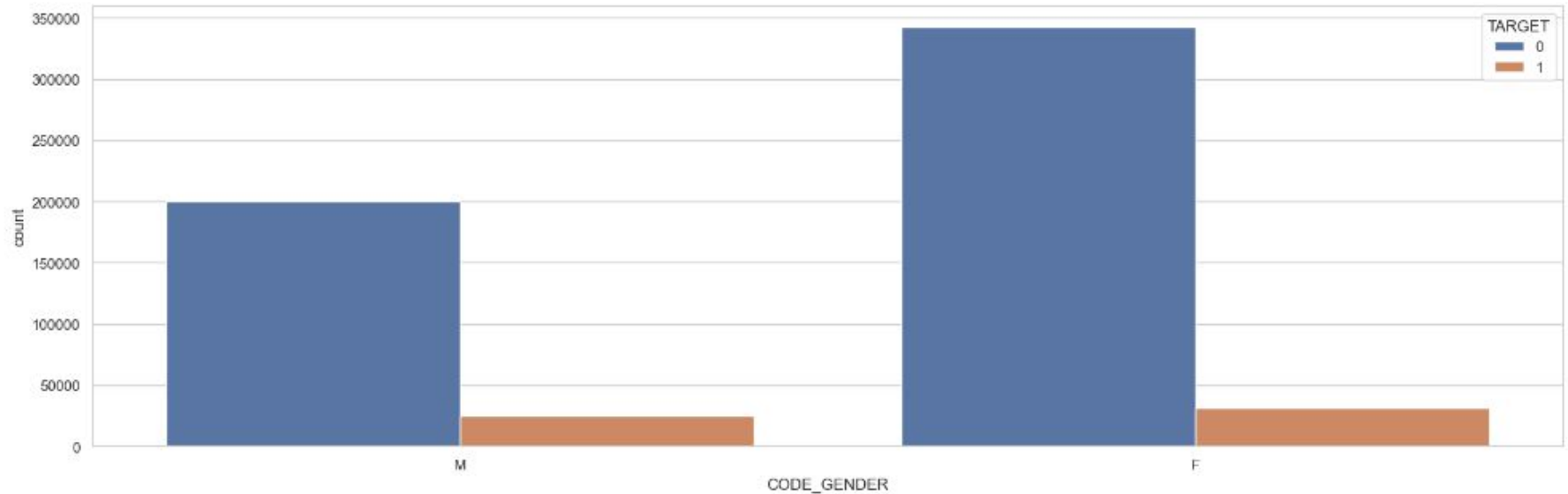
#Plotting NAME_CONTRACT_TYPE with gender



Merging the two data set

Merging previous application data to application data for further analysis

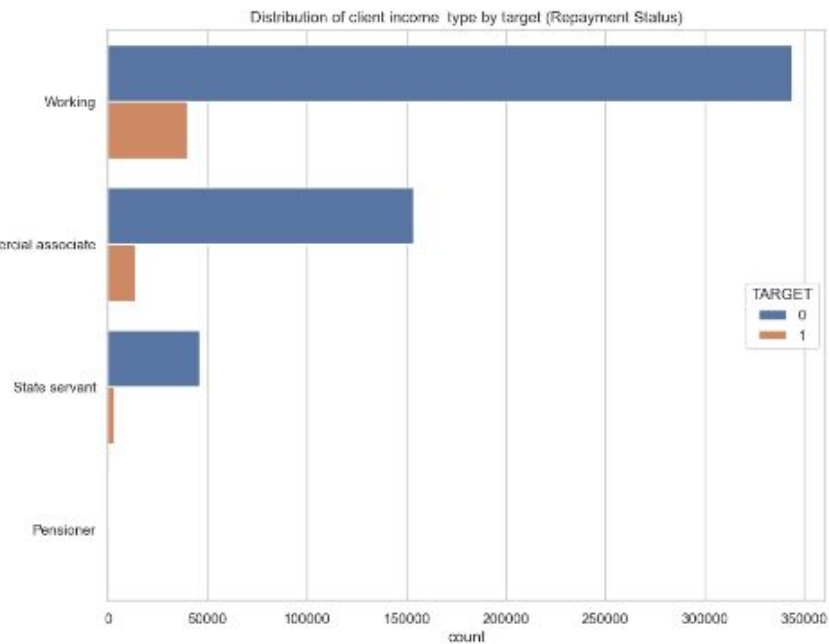
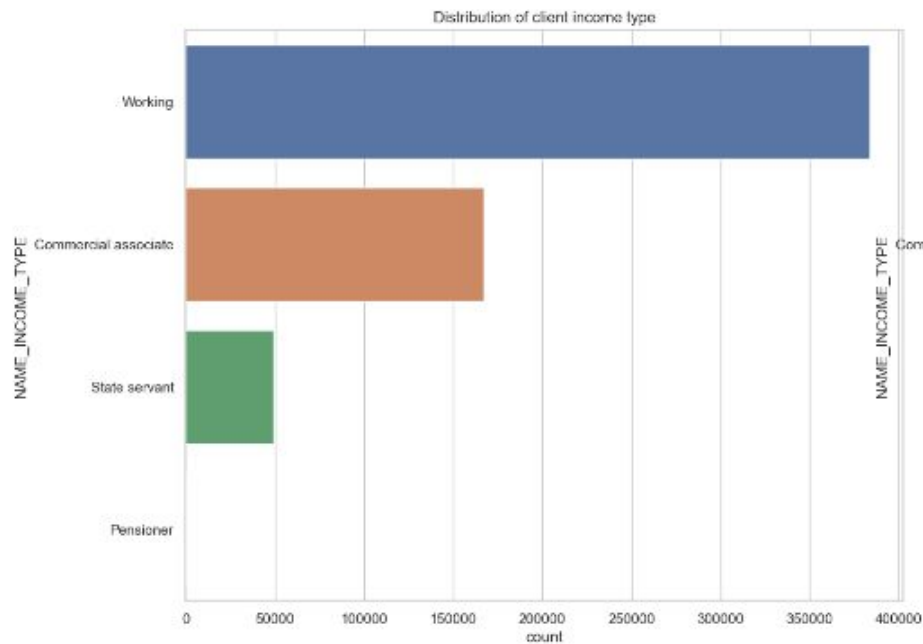
#Distribution of CODE_GENDER by TARGET for repayment status



Clearly, female clients are the best repayers of their loan (almost double the amount of males).

Amount of defaulters in both genders are almost equally distributed.

#Distribution of client income



Most clients as per both cases of repayment status, are working

On the other end, the least amount of clients are pensioners

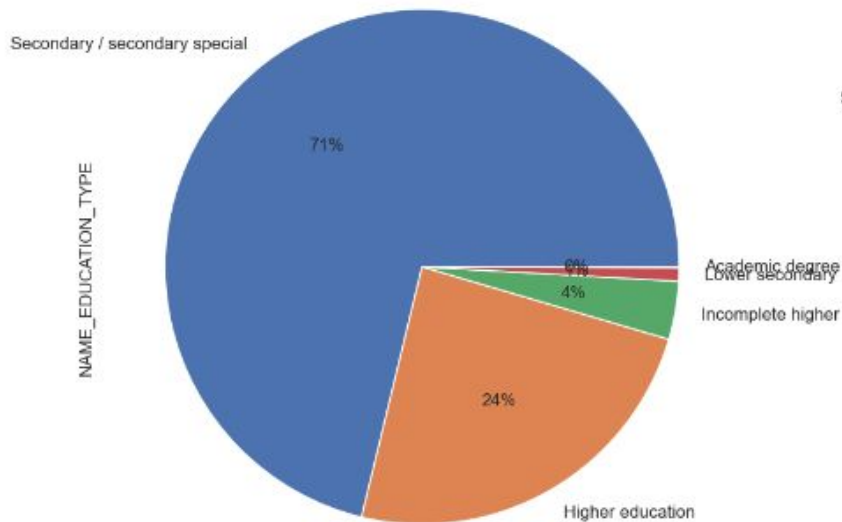
#Distribution of education type by repayment status

#For clients with secondary edu, default is proportionally 9% higher compared to clients who do not default

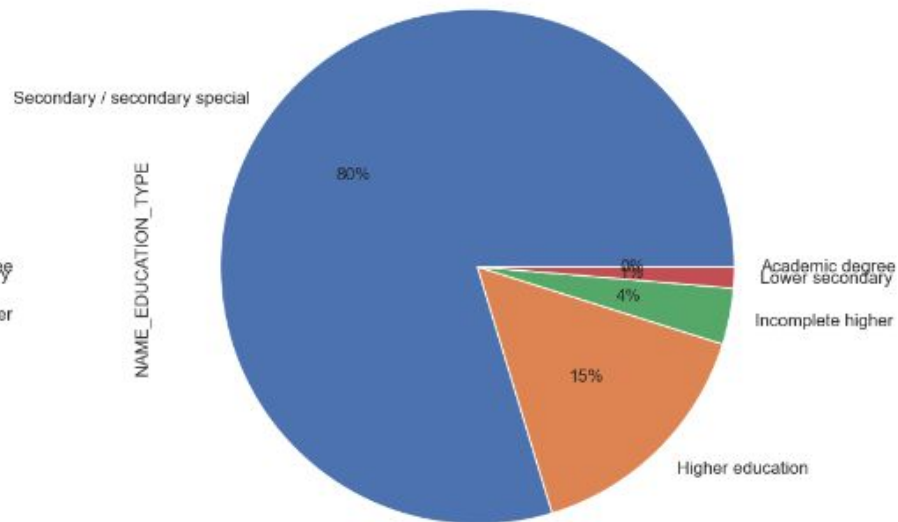
#In the higher education category, clients who default are 8% fewer.

#In both cases of repayment status, lower secondary and academic degree categories are the minority.

Distribution of Education for Repayers



Distribution of Education for Defaulters

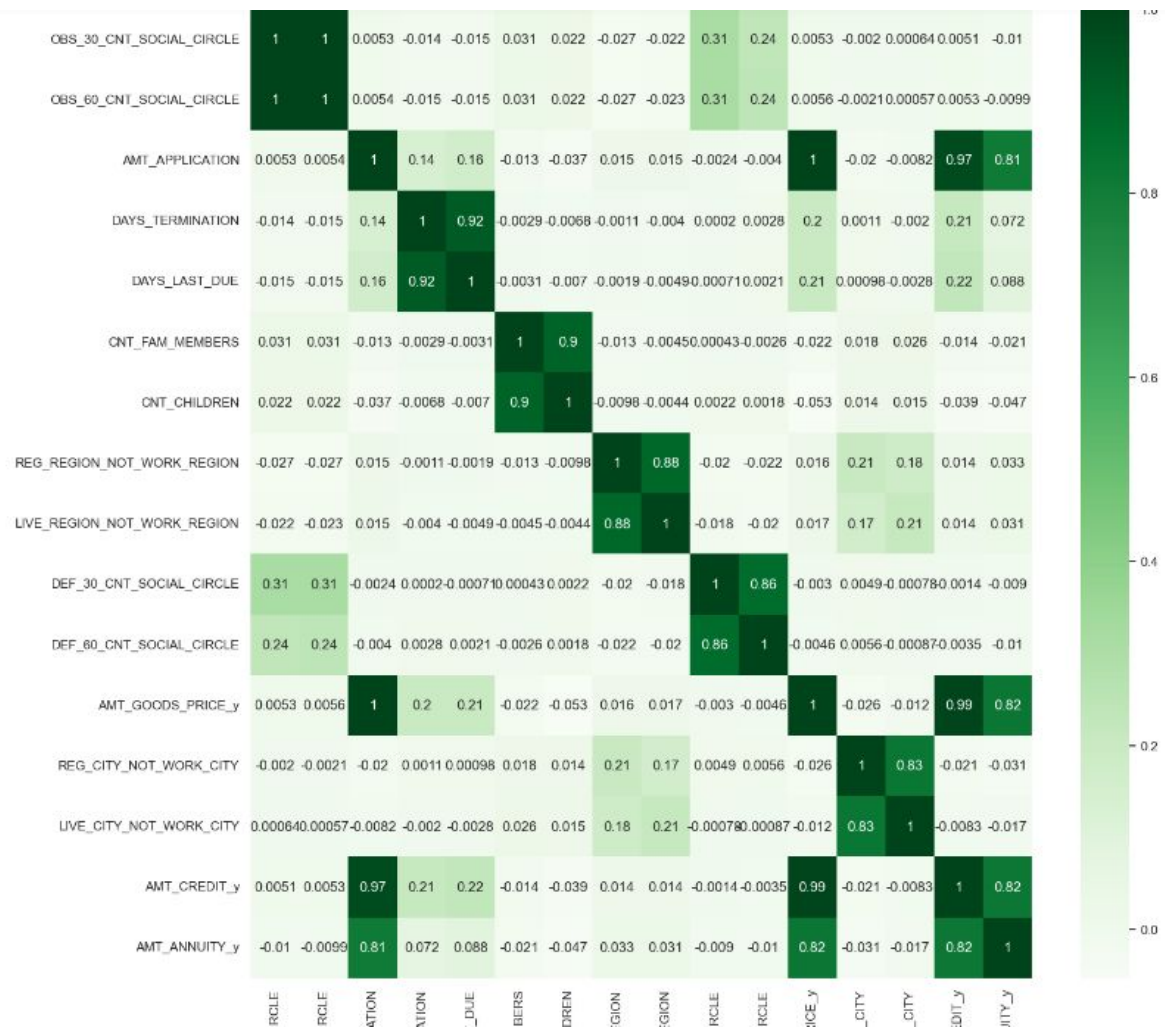


Finding Correlation

Correlations between numerical data can be used to infer

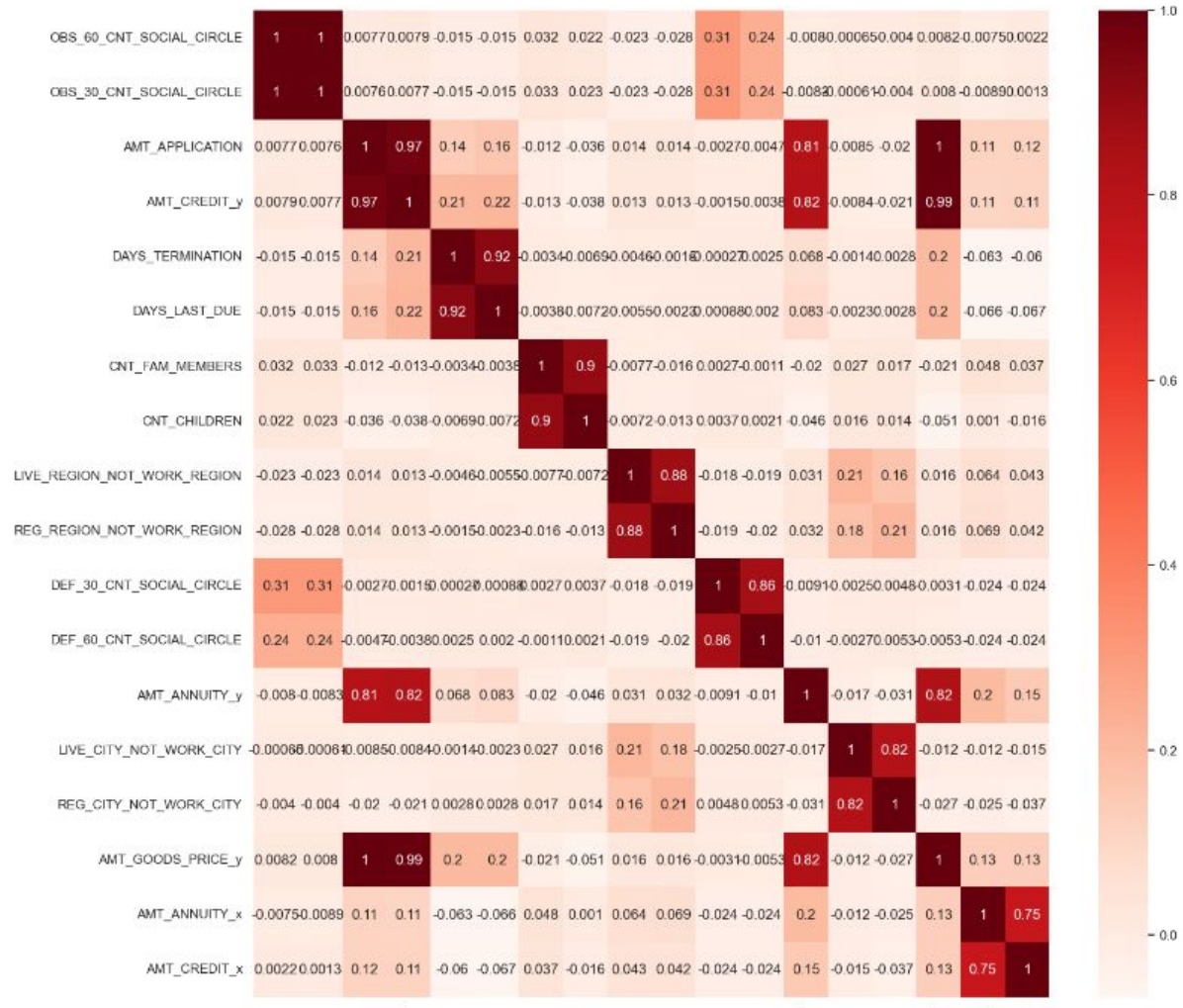
Top 10 correlated columns for repayers

- OBS_30_CNT_SOCIAL_CIRCLE OBS_60_CNT_SOCIAL_CIRCLE 1.00
- AMT_CREDIT_y AMT_APPLICATION 0.97
- DAYS_TERMINATION DAYS_LAST_DUE 0.93
- CNT_FAM_MEMBERS CNT_CHILDREN 0.90
- REG_REGION_NOT_WORK_REGION LIVE_REGION_NOT_WORK_REGION 0.88
- DEF_30_CNT_SOCIAL_CIRCLE DEF_60_CNT_SOCIAL_CIRCLE 0.87
- AMT_GOODS_PRICE_y AMT_CREDIT_y 0.86
- AMT_APPLICATION AMT_GOODS_PRICE_y 0.85
- REG_CITY_NOT_WORK_CITY LIVE_CITY_NOT_WORK_CITY 0.83
- AMT_CREDIT_y AMT_ANNUITY_y 0.81



Top 10 correlations for Defaulters

- OBS_60_CNT_SOCIAL_CIRCLE OBS_30_CNT_SOCIAL_CIRCLE 1.00
- AMT_APPLICATION AMT_CREDIT_y 0.97
- DAYS_TERMINATION DAYS_LAST_DUE 0.95
- CNT_FAM_MEMBERS CNT_CHILDREN 0.90
- LIVE_REGION_NOT_WORK_REGION REG_REGION_NOT_WORK_REGION 0.87
- DEF_30_CNT_SOCIAL_CIRCLE DEF_60_CNT_SOCIAL_CIRCLE 0.86
- AMT_CREDIT_y AMT_ANNUITY_y 0.83
- LIVE_CITY_NOT_WORK_CITY REG_CITY_NOT_WORK_CITY 0.78
- AMT_ANNUITY_y AMT_GOODS_PRICE_y 0.76
- AMT_ANNUITY_x AMT_CREDIT_x 0.74




CONCLUSION:

Target

1. Students, Pensioners and Commercial Associates with a housing type such as office/co-op/municipal apartments
2. Clients living with parents
3. Female clients, as they have a high rate of repayment
4. Should target client who own car

Not to target

1. Based on education, as there is no sufficient evidence available
2. Females on maternity leave
3. Clients 'working', as they have highest defaulters
4. 'Repairs' purpose of loans have the most amount of defaulters, loan disbursement should be of low risk



Thank You. It was a pleasure
brainstorming and finishing this
assignment.