# EDA Case Study

**Group Members:**

- Mamta Mittal

- Sneha Boora

# Abstract

Business Understanding:
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.

Business Object:
This case study aims to identify patterns which indicate if a client has difficulty paying their installments 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.

Goals of Data Analysis:
To understand the  driving factors for loan default, identifying patterns which indicate if the client has difficulty paying installments using EDA.

# Data Understanding

1. *'application_data.csv'* contains all the information of the client at the time of application.
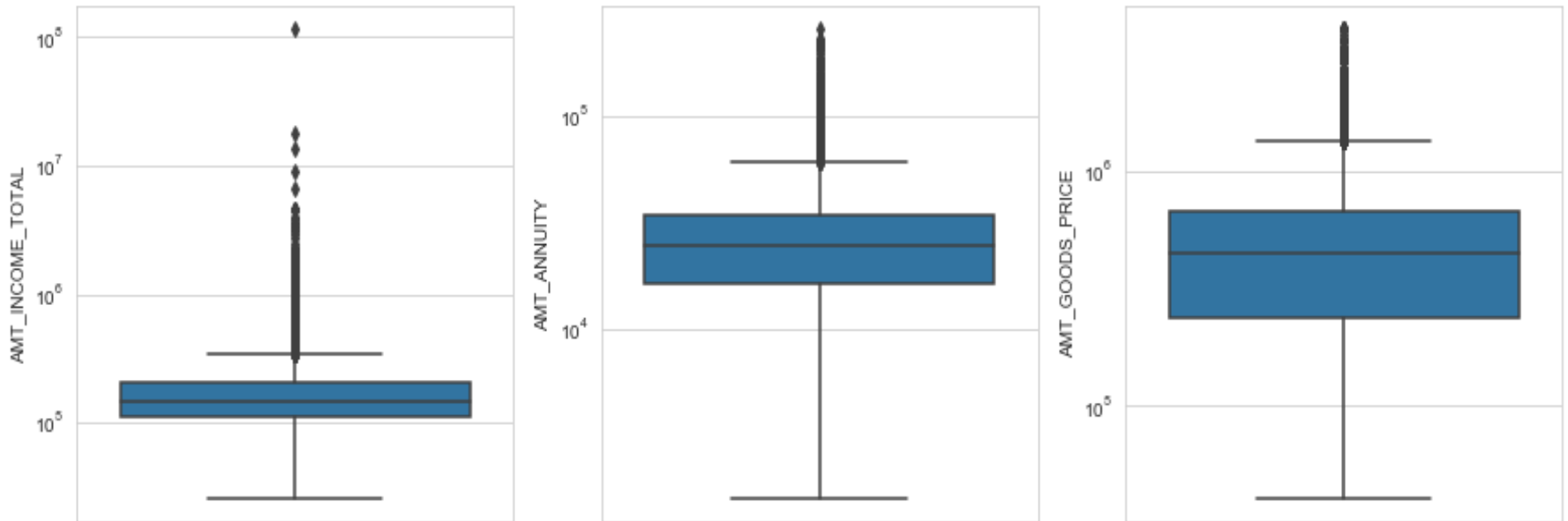The data is about whether a **client has payment difficulties.**

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.**

# Approach

1. Read the dataset
2. Clean the dataset(Drop unwanted rows and columns, check data type,)
3. Identify the outliers
4. Analysis the DF_APP_DATA and get the insight from the analysis of categorical variable and numeric variable by plotting required plots(univariate and Bivariate)
5. Analysis the DF_PREV_APP and get the insight from the analysis of categorical variable and numeric variable by plotting required plots(univariate and Bivariate)
6. Merge both the data set and analyse the impact of previous application data on the Current application by checking it against TARGET variable.
7. Find the major factors for Bank which can help in Credit risk analysis.
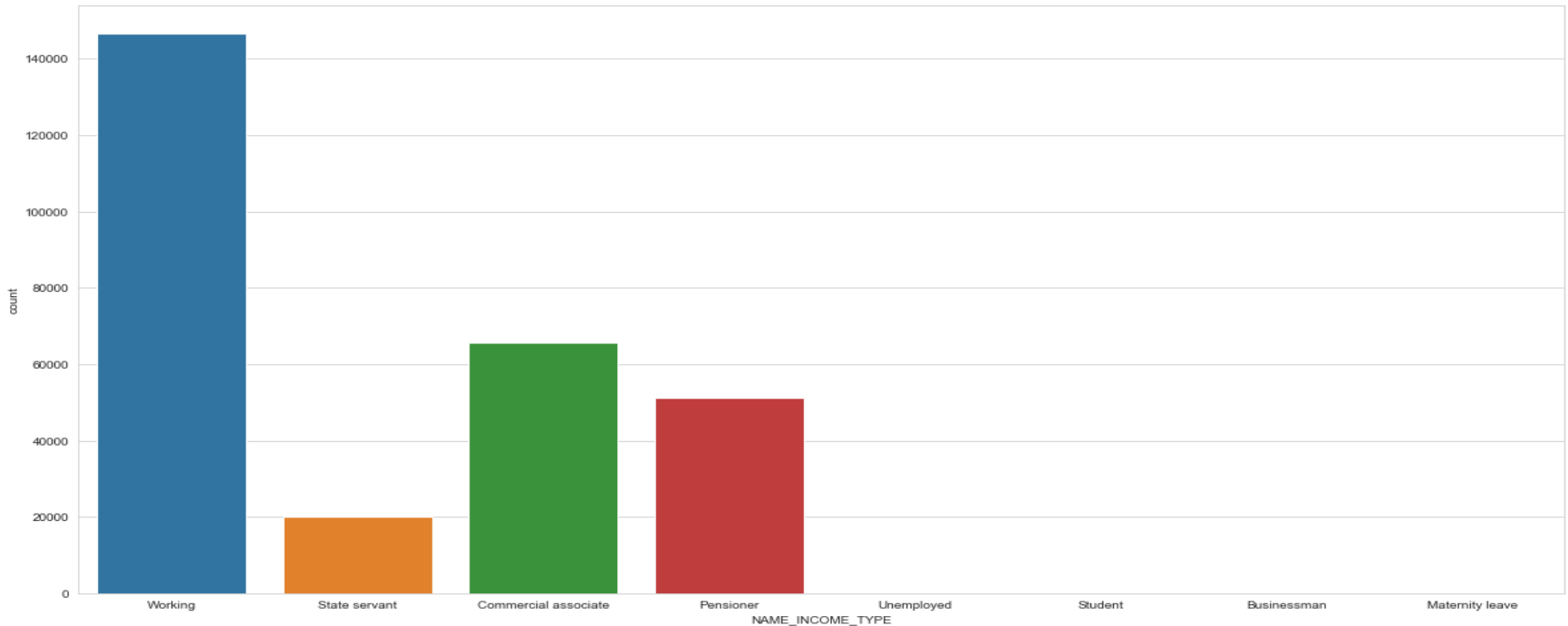
# Outliers



**Observation:**

•From the above box plot for AMT_INCOME_TOTAL/AMT_ANNUITY/AMT_GOODS_PRICE we can observe that there are lots of outliers in the data.

•Outliers are the extremely High and low values. Due to which data is skewed.

•Using DF_APP_DATA_NUM_COL.describe() metric we saw which columns have outliers by checking difference in mean, median values.
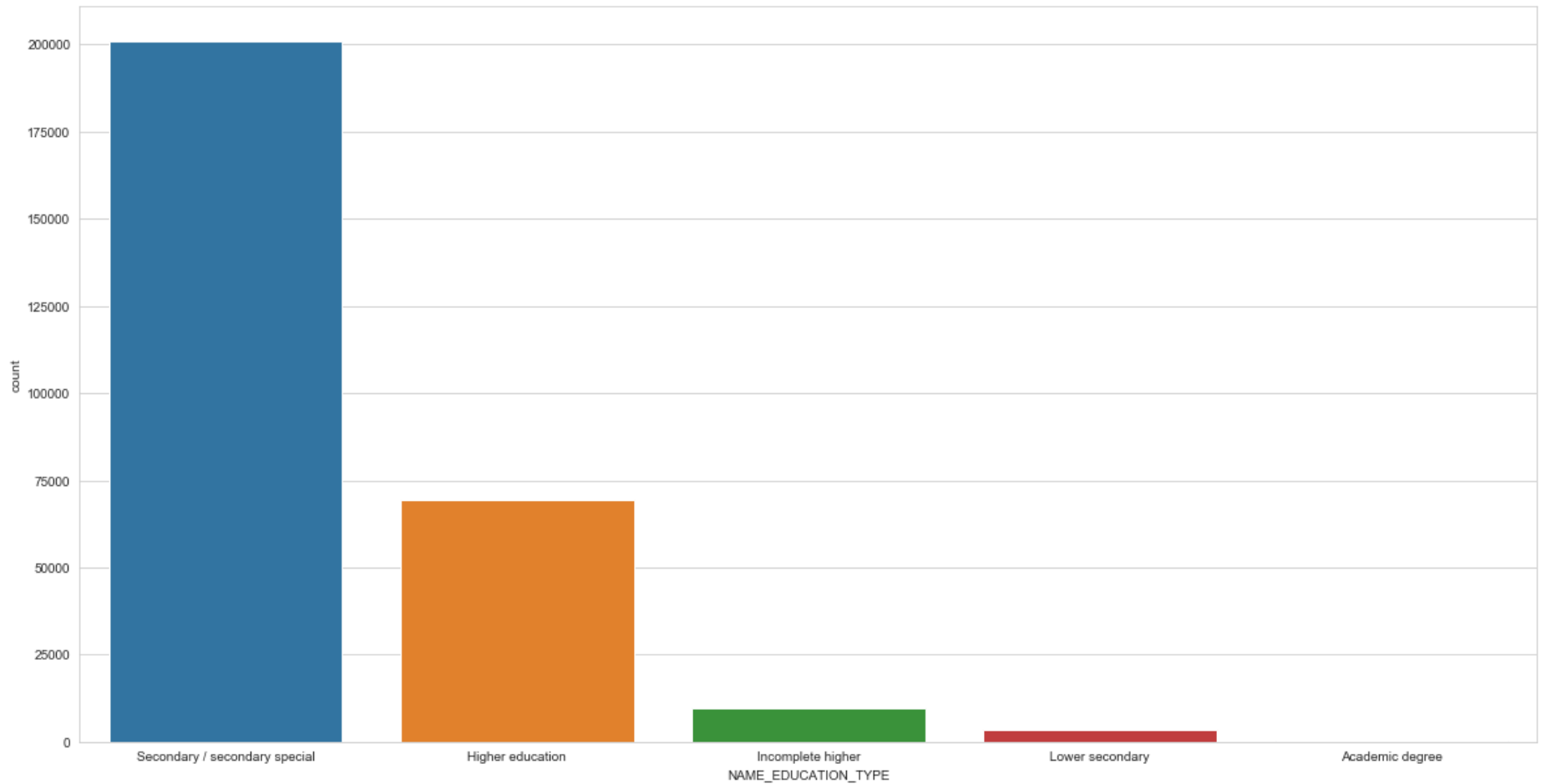
# Data Imbalance of 'TARGET' variable
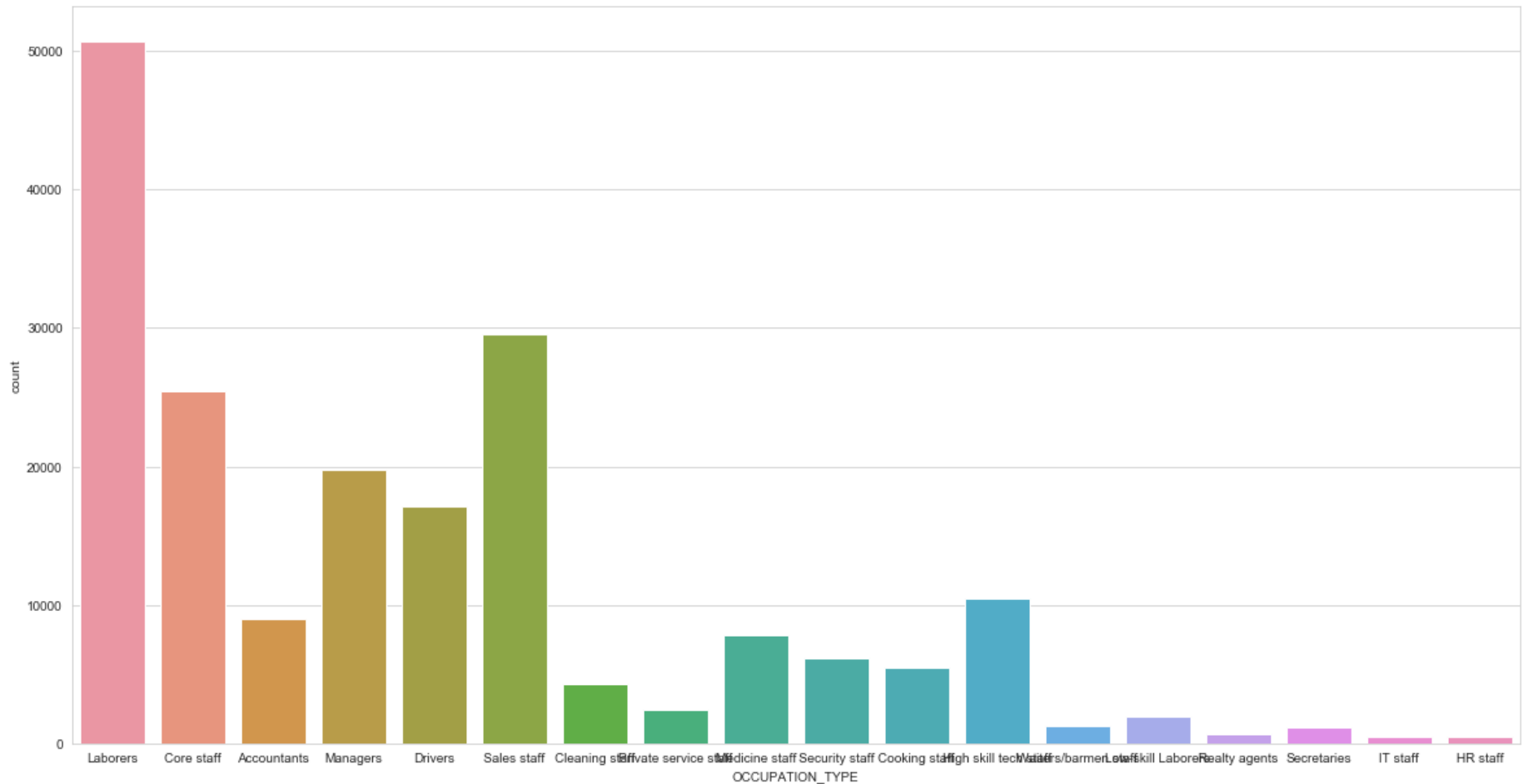


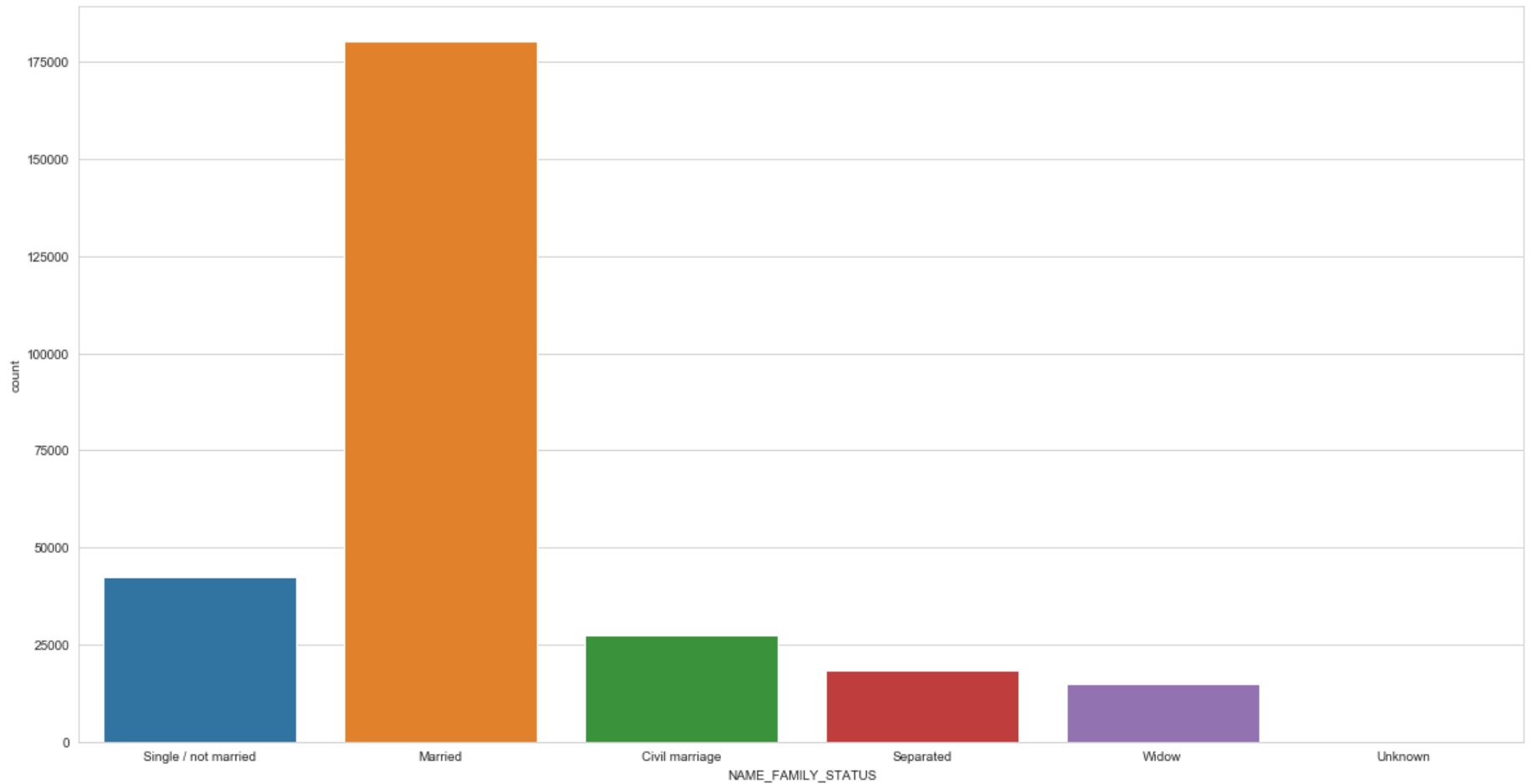Ratio of imbalance is found out to be 11.375961202376791

# Uni-variate Analysis



**Observation (NAME_INCOME_TYPE):** No. of people opting for Loan for income categories Working/ State Servent /Commercial Associate/Pensioner are significant as compare to other categories
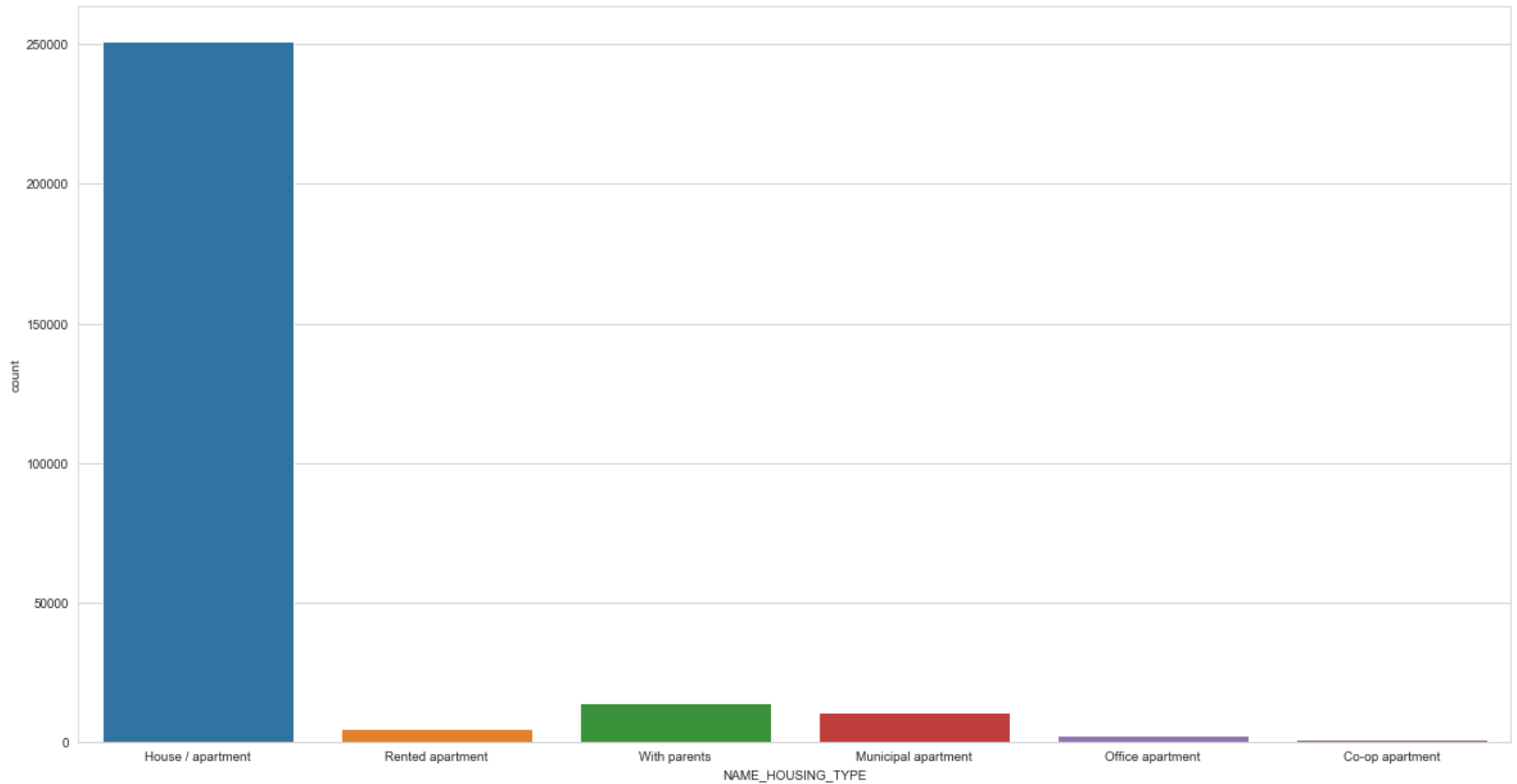
**Observation(NAME_EDUCATION_TYPE) :** No. of people opting for Loan for Education type categories Seconary/Higher are significant as compare to other categories.
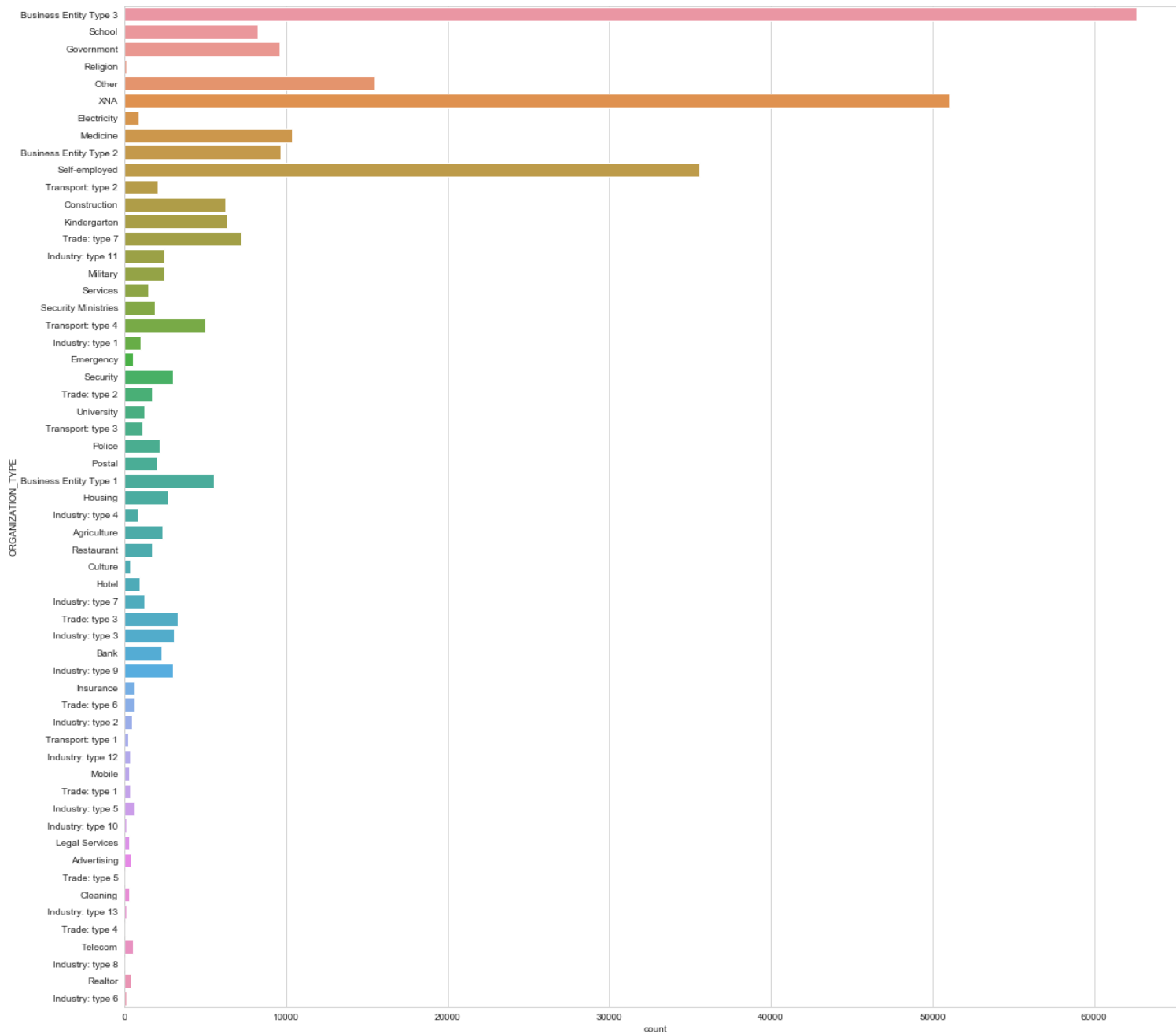
**Observation(OCCUPATION_TYPE) :** No. of people opting for Loan for occupation categories Laborers/Core Staff/Sales Staff are significant as compare to other categories.

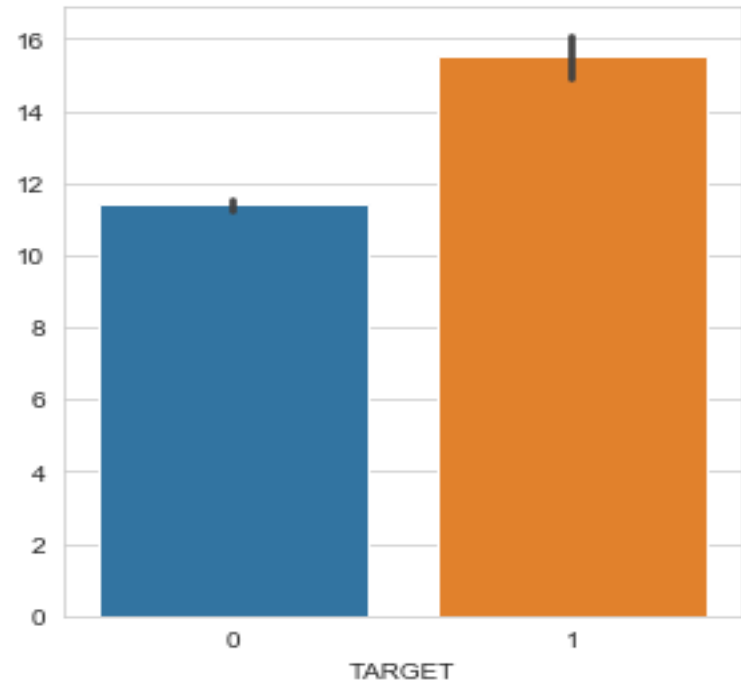**Observation(NAME_FAMILY_STATUS) :** No. of Married people applying for loan is much higher then other categories.
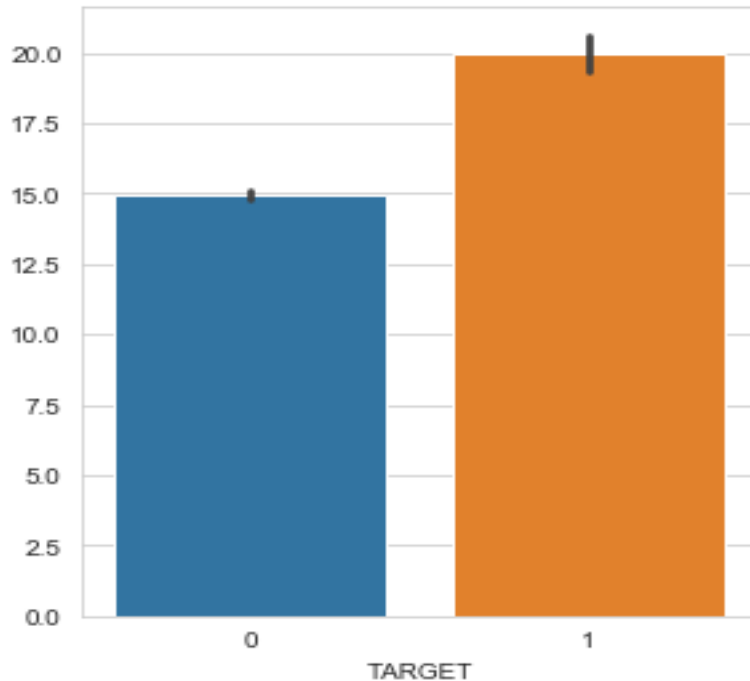
**Observation(NAME_HOUSING_TYPE) :** No. of people staying in Apartment are applying for loan is much higher then other categories.

**Observation( ORGANIZATION_TYPE) :** people in this category "Self Employed and Business entity type 3 " Apply for more loan
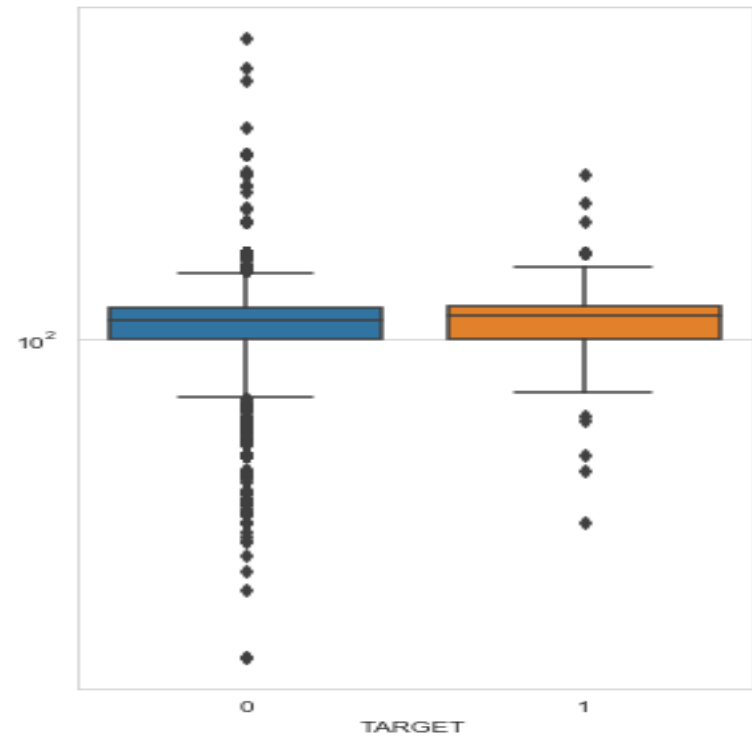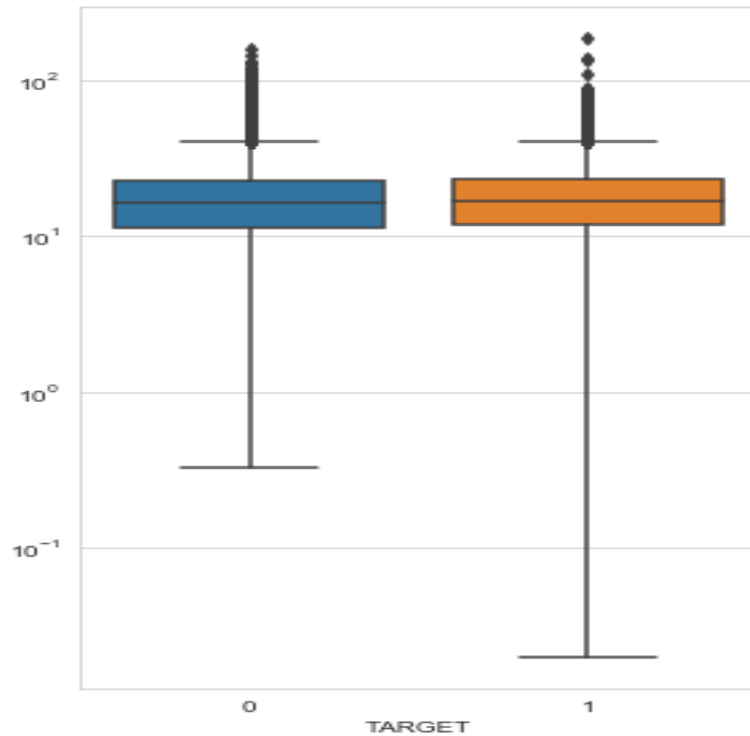
# Bi-variate Analysis



**Observation (OBS_30_CNT_SOCIAL_CIRCLE/OBS_60_CNT_SOCIAL_CIRCLE):**
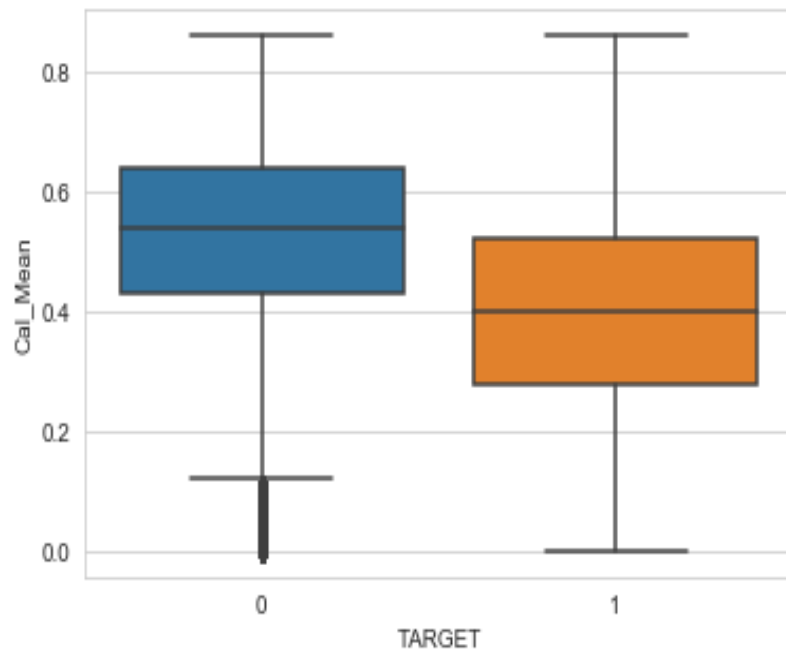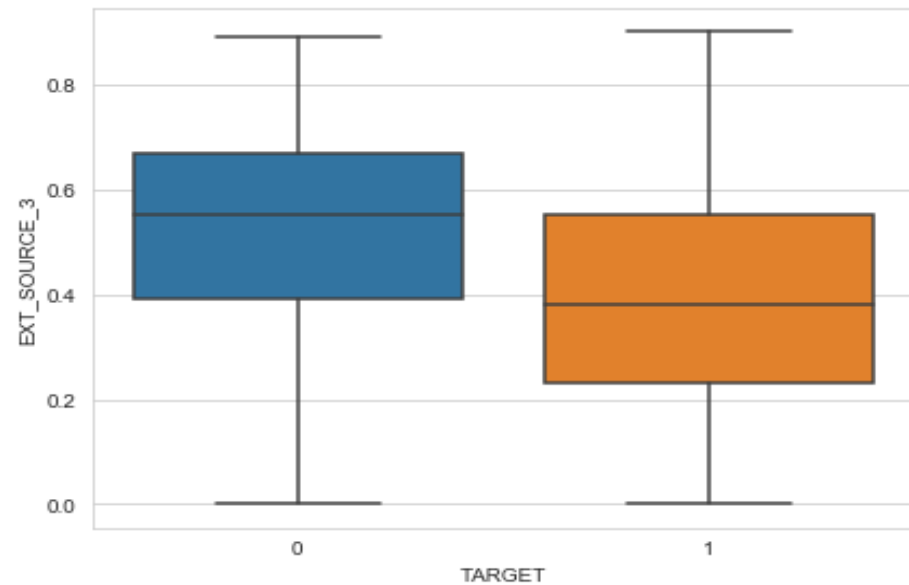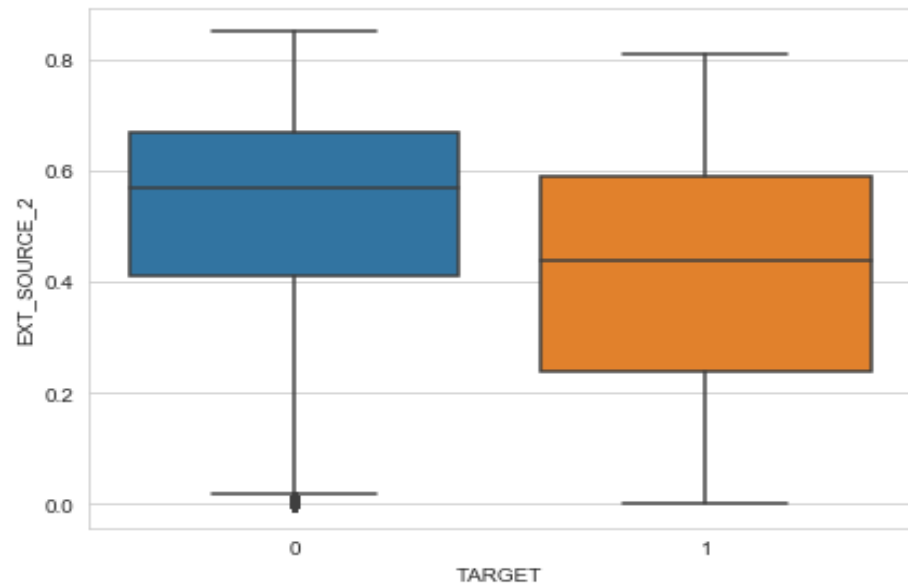As there are many people in the Social circle who are defaulter , Chance of defaulting a person increase because social circle has an impact on persons habits.
Hence this is one of the factor which influence the default case
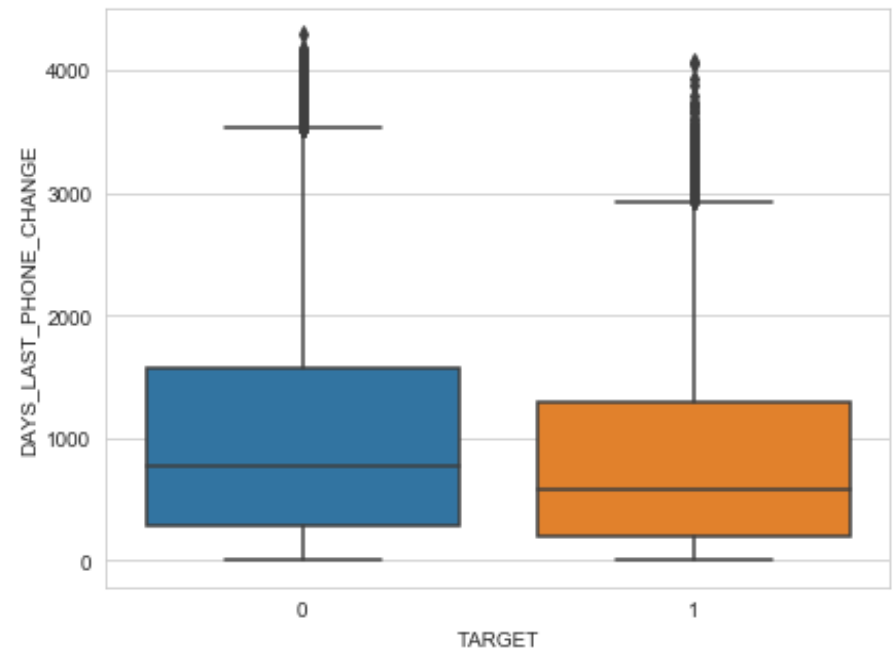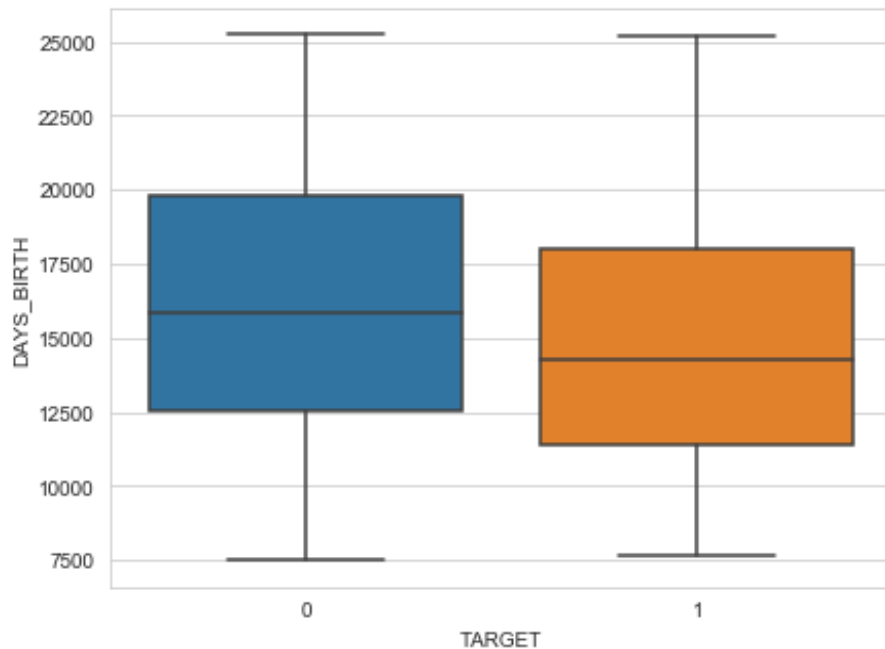
**Observation:**
*(AMT_ANNUITY/AMT_INCOME_TOTAL/AMT_CREDIT/AMT_GOODS_PRICE) : All These 4 columns does not seem to have an Impact on default cases*

**Observation:**

From All the above plots we can observe that EXT_SOURCE_1/2/3 has a huge impact on defaulter.

As EXT_SOURCE is the credit rating provided by the third party. Based on the Credit rating Bank may decide if the loan can be given to the customer. Customer with lower credit rating are more likely to default.
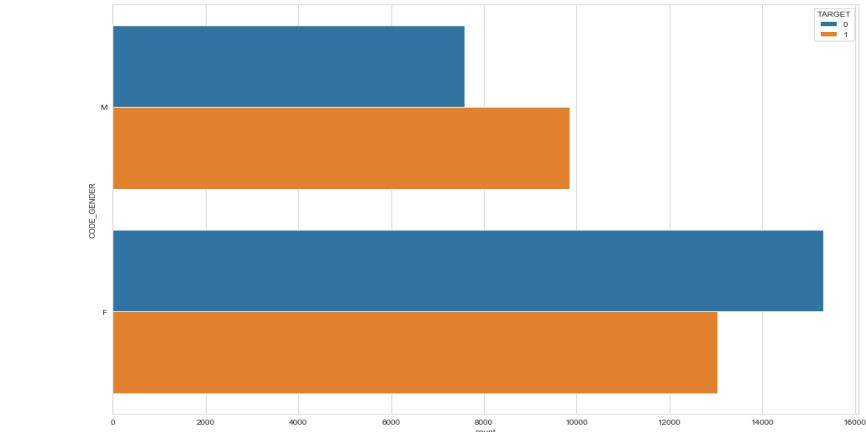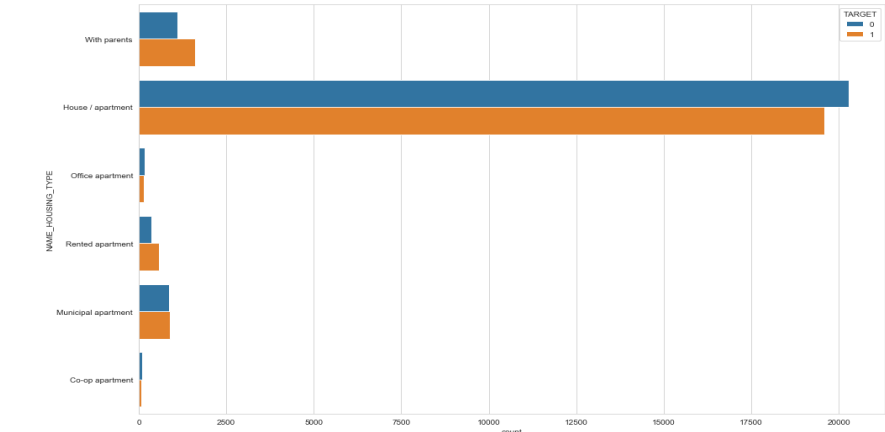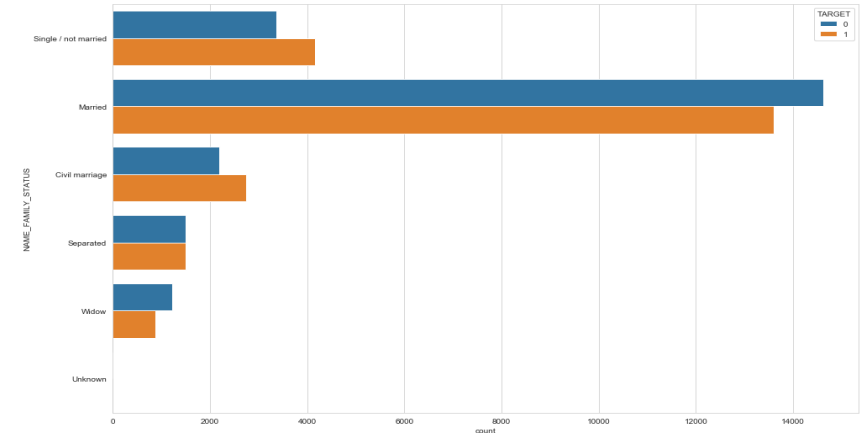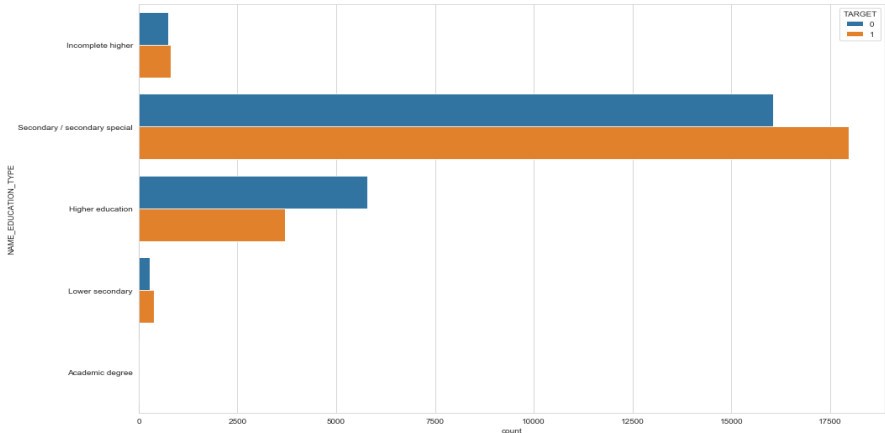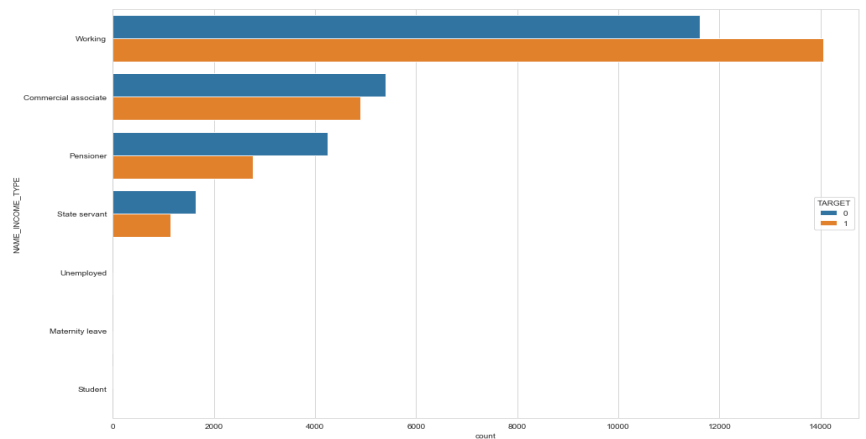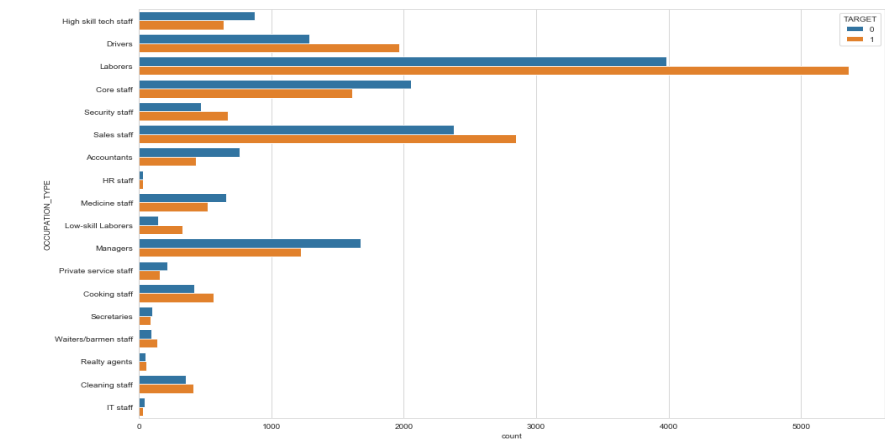
**Observations(**DAYS_BIRTH / DAYS_LAST_PHONE_CHANGE**):**
From All the above plots we can observe that DAYS_BIRTH : Higher the age(may be Old age), more likely the person will default
DAYS_LAST_PHONE_CHANGE: if the phone number has been changed recently, chances of a person to default increases
We may consider these 2 columns as itmay have impact on defaulter.

**Observation:**

**Following fields are significant**

- CODE_GENDER
  - 'FeMale' having much more % population with NO payment difficulty
- NAME_INCOME_TYPE
  - with values "State servant" and "Pensioner" having much more % population with no payment difficulty
- NAME_EDUCATION_TYPE
  - "Higher education" having much more % population with NO payment difficulty
  - "Lower secondary" and "Secondary" having much more % population with payment difficulty
- NAME_HOUSING_TYPE
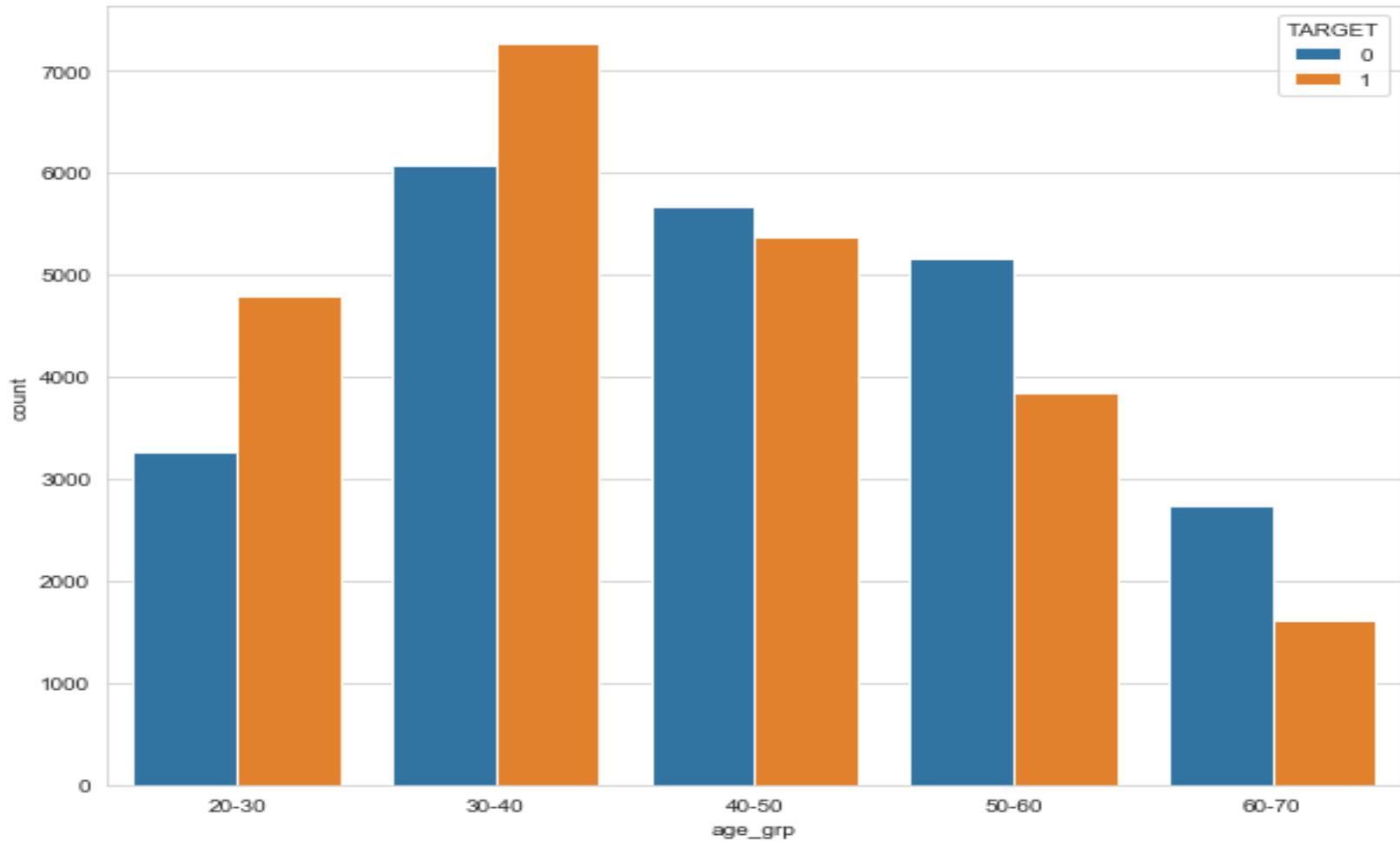  - "With parents" and "Rented apartment" having much more % population with payment difficulty
- OCCUPATION_TYPE
  - "Security staff", "Laborers", "Cooking staff", "Drivers" having much more % in population with payment difficulty
  - "Core staff","High skill tech staff","Accountants" having much more % in population with NO payment difficulty
- ORGANIZATION_TYPE
  - "Construction", "Industry: type 3" having much more % in population with payment difficulty
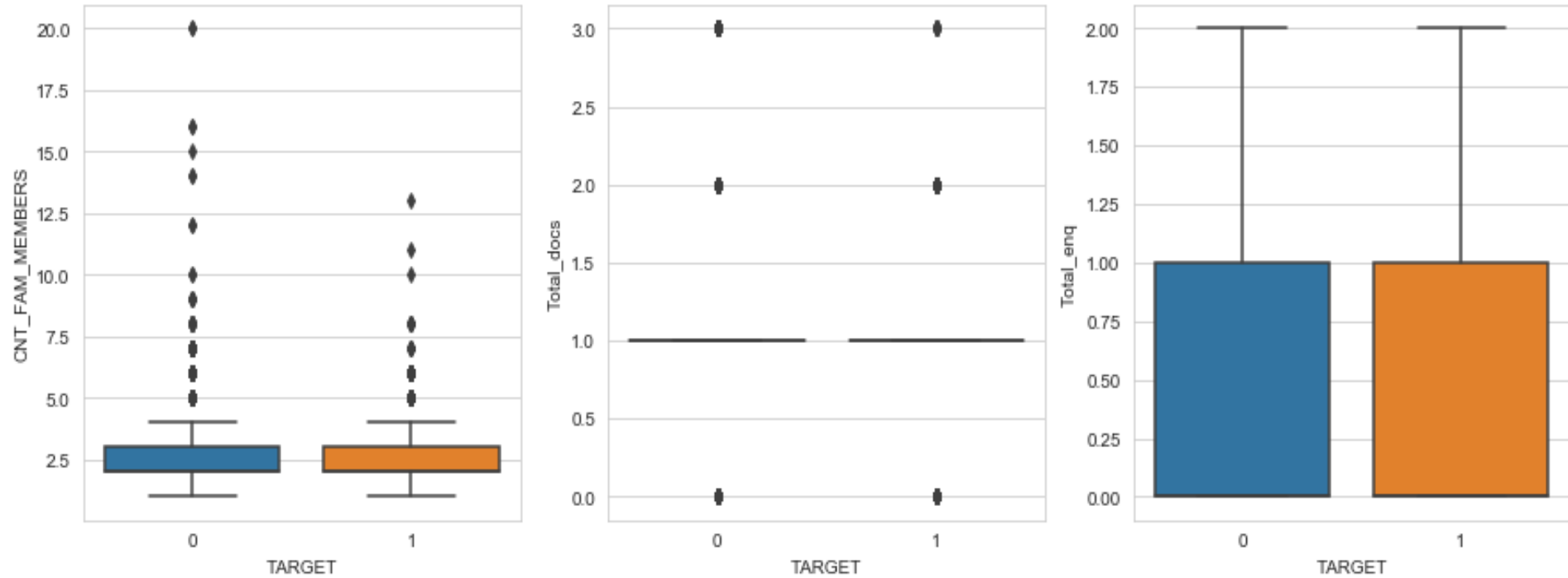  - "School" having much more % in population with NO payment difficulty

# Analysis of 'age_grp' variable which is derived by binning DAYS_OF_BIRTH



**Observation:**
age_grp: 20-30/30-40 having much more % population with payment difficulty.
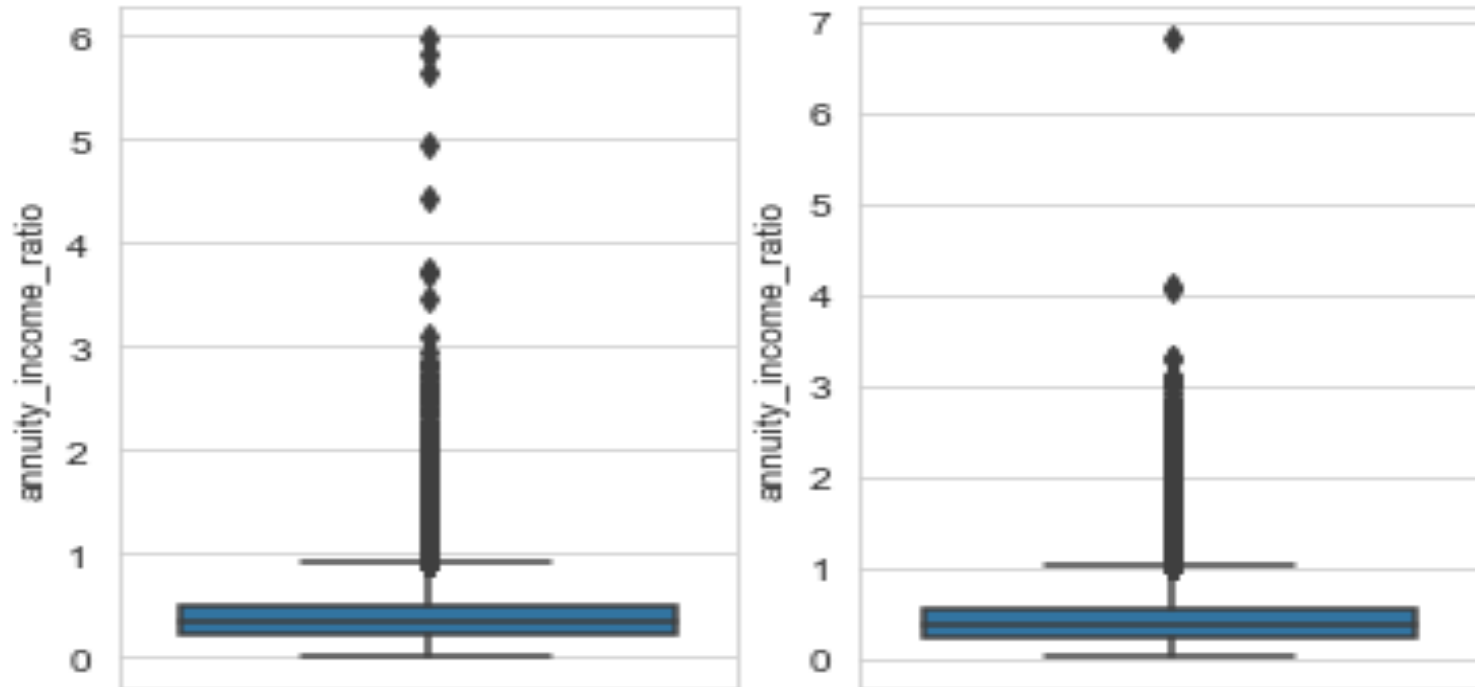
Few other variables:



**Observation: NO Impact**
CNT_FAM_MEMBERS/Total_docs/Total_enq/DAYS_REGISTRATION/DAYS_EMPLOYED/REG_REGION_NOT_LIVE_REGION
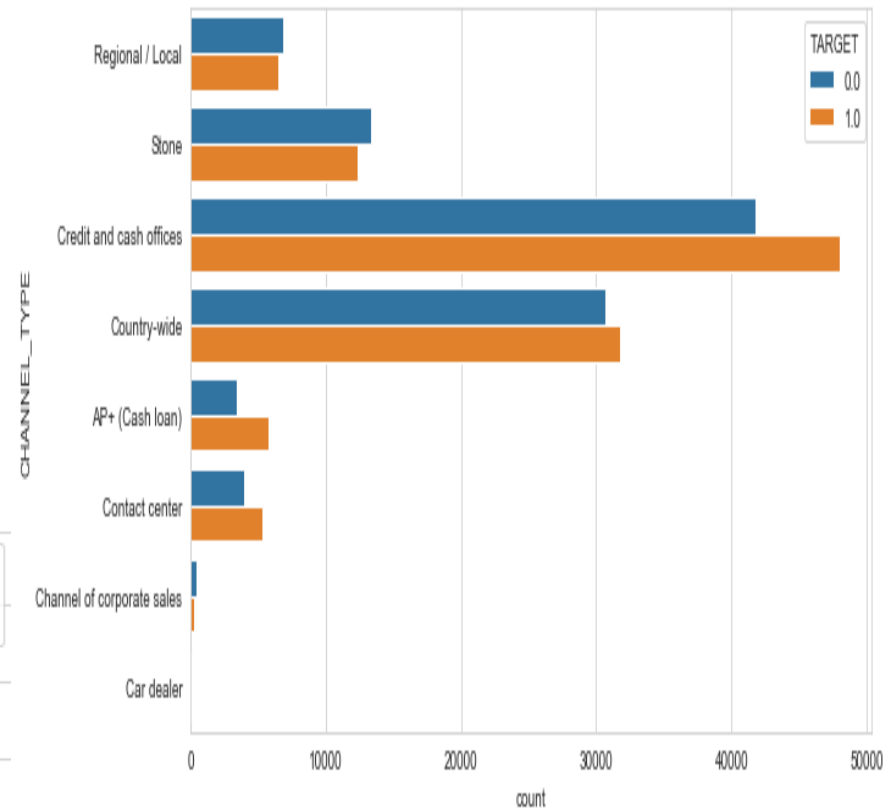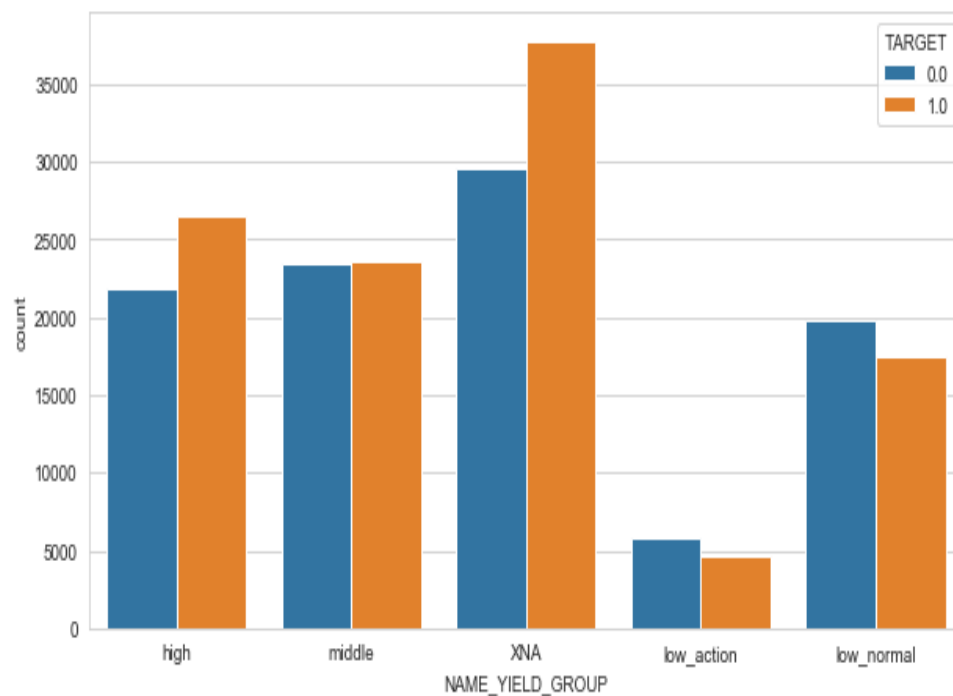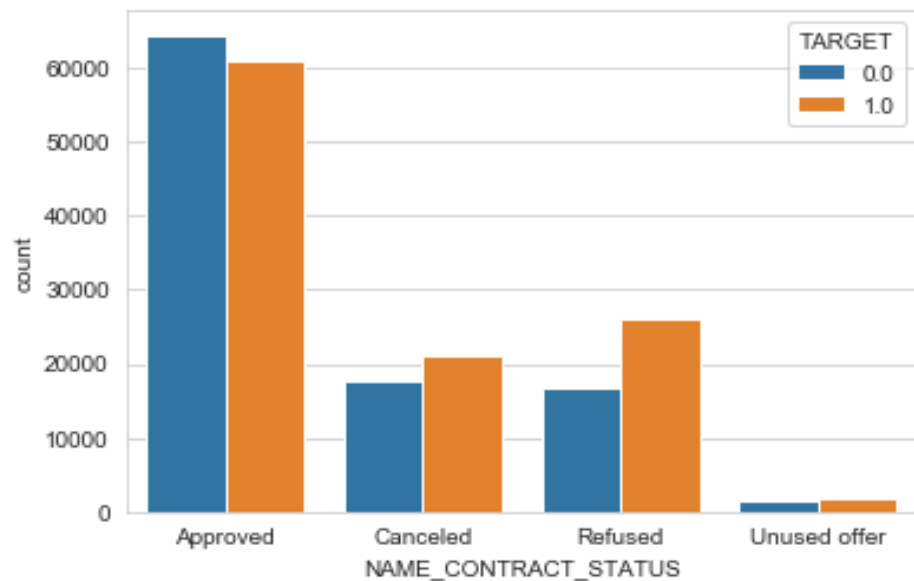
# Previous application - Approach

•Sum the annuity for all previous applications where the payment is still due
•filter where DAYS_TERMINATION is still +ve i.e. loan is still active

•Add the outstanding annuity to current application annuity

•Find the annuity to income ratio

•Compare between Target 0 and target 1 population and see if any significant difference

# Further Analysis



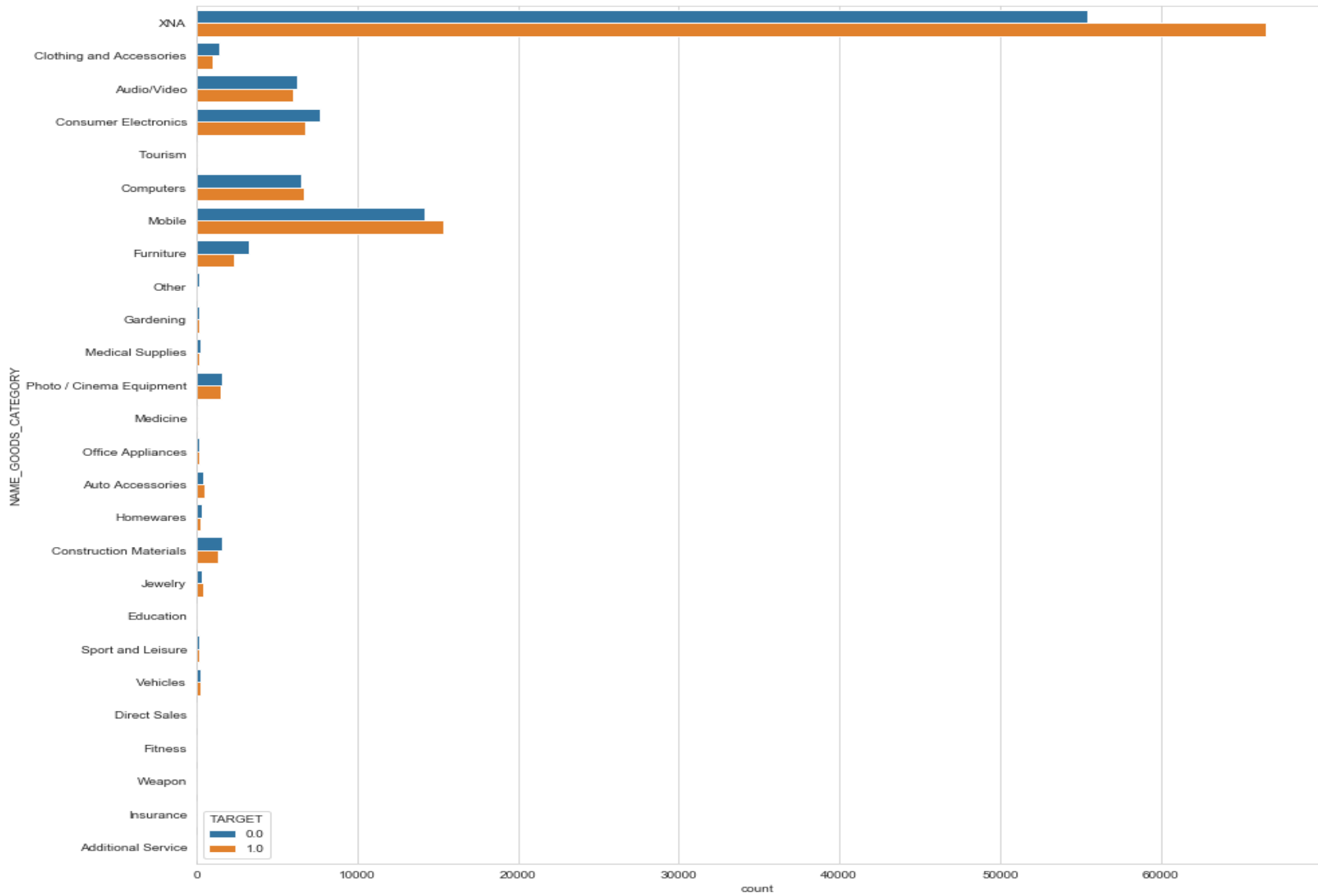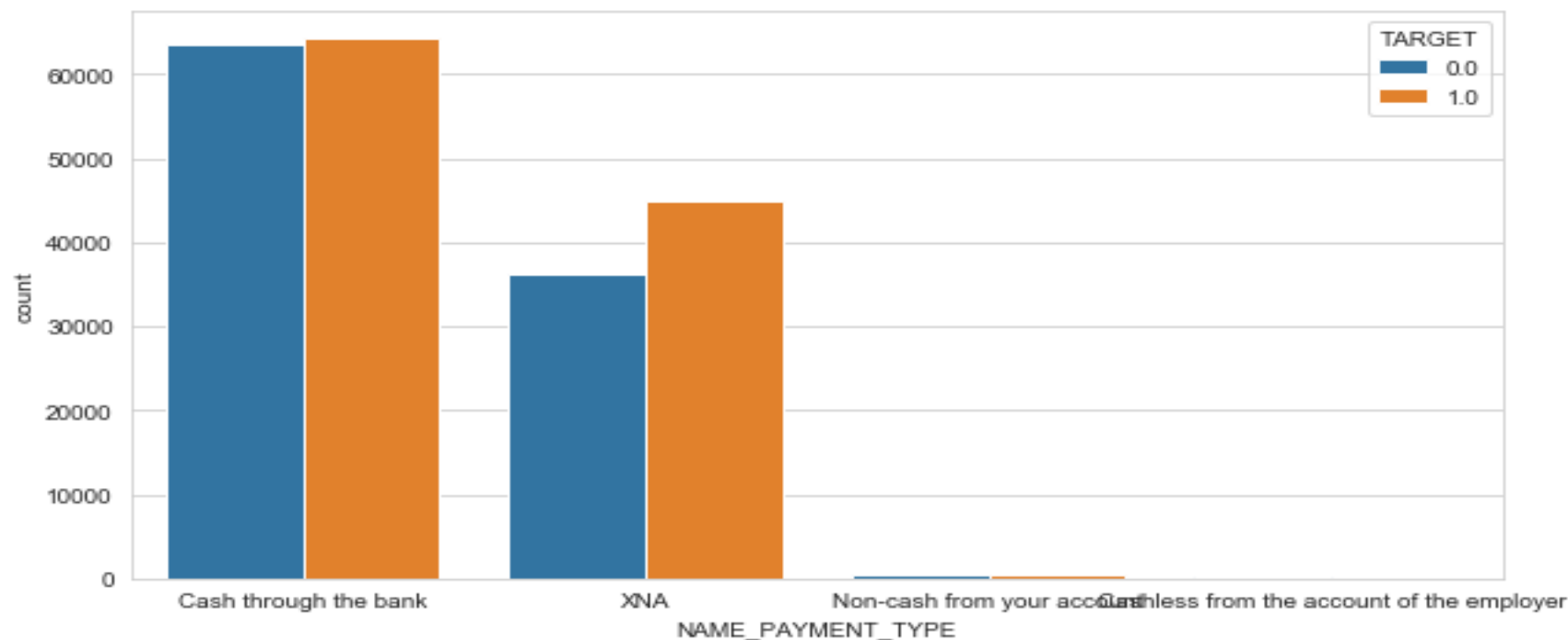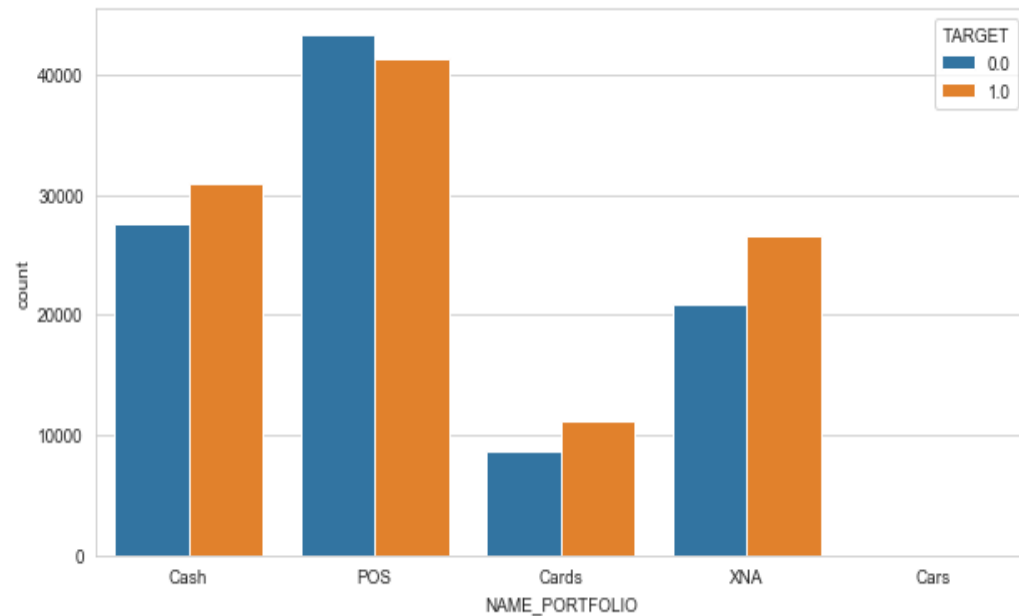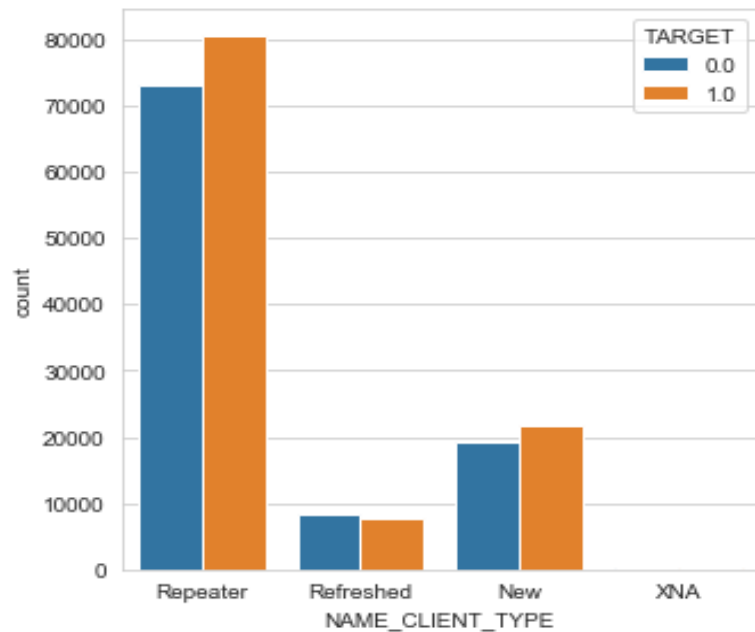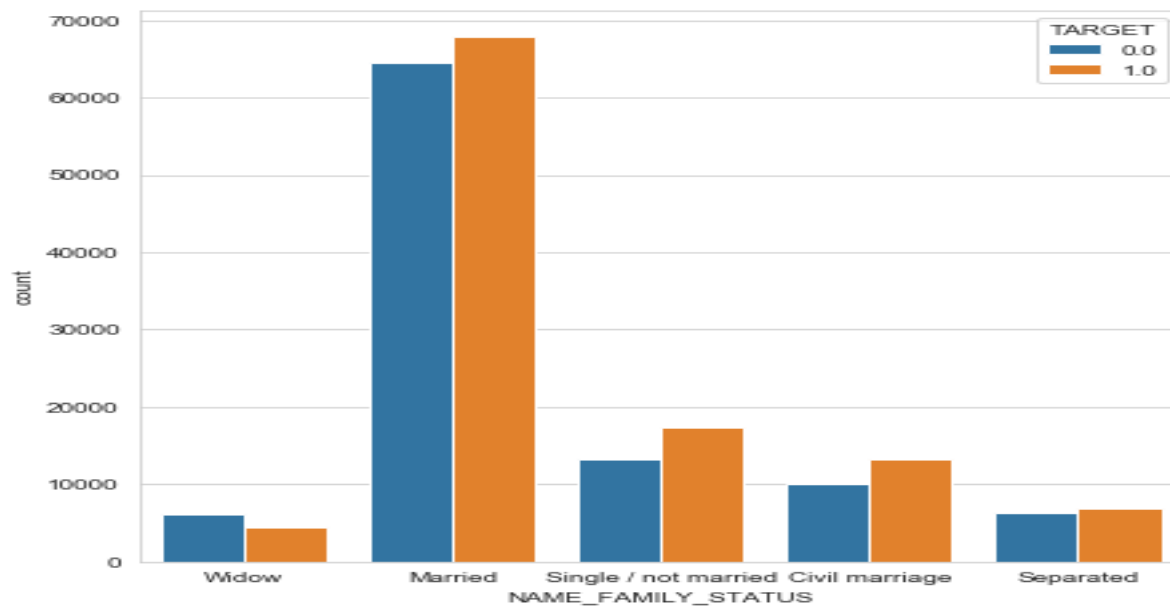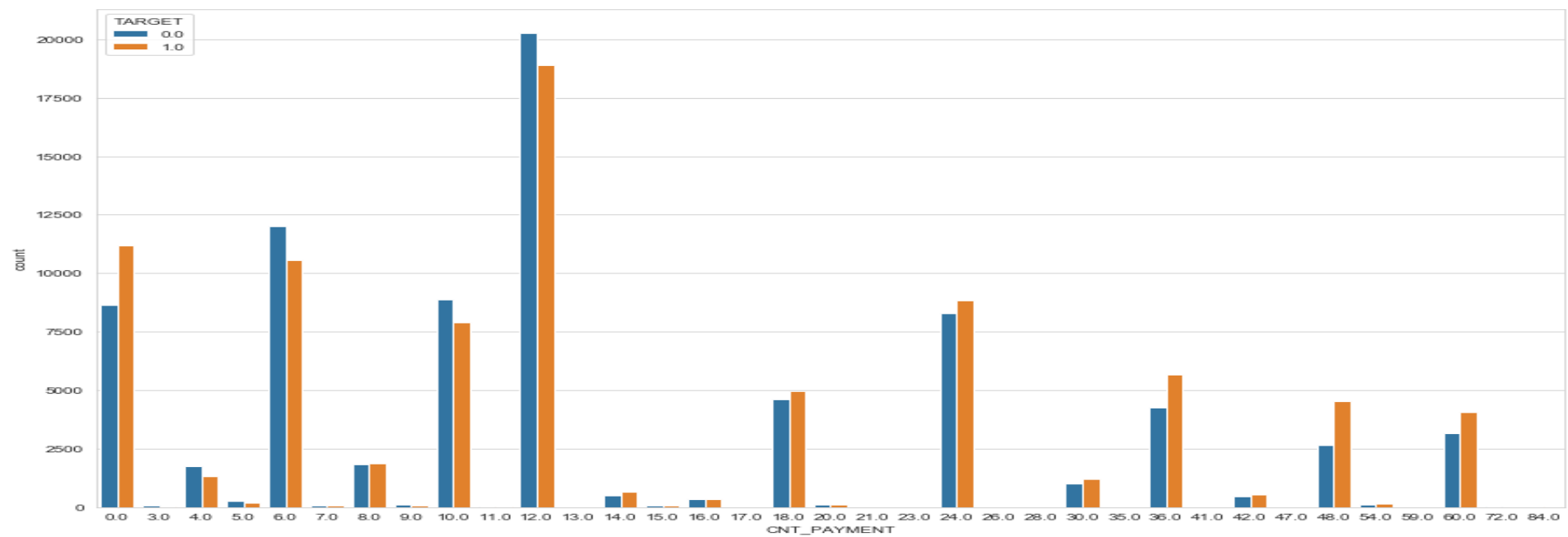**Conclusion : No Significance**
As we can see the mean, median and inter-quartile ranges are almost similar between Target 0 and Target1. So it can be concluded that Annuity and Income does not have any significant difference for Target 0 and Target 1 population

**Observation:**

**Following fields are significant**

NAME_CONTRACT_STATUS

'Refused' having much more % population with payment difficulty. Which is correct. Loan should be refused to customer with payment difficulties

NAME_YIELD_GROUP

'High interest rate' having much more % population with payment difficulty

CHANNEL_TYPE

with values "Credit and cash Office/contact center/ AP+(cash loan)" having much more % population with payment difficulty

NAME_SELLER_INDUSTRY

"Connectivity" having much more % population with payment difficulty

PRODUCT_COMBINATION

"CASH/Card Street/POS Mobile with interest" having much more % population with payment difficulty

NAME_GOODS_CATEGORY

"Mobile" having much more % population with payment difficulty

NAME_CLIENT_TYPE

"Repeater" having much more % in population with payment difficulty

NAME_PAYMENT_TYPE

No Significant Impact

NAME_PORTFOLIO

"Cash/Cards" having much more % in population with payment difficulty

CNT_PAYMENT

'long Term loan on previous application' having much more % in population with payment difficulty

NAME_FAMILY_STATUS

'Single/Not Married/Civil_Marriage' having much more % in population with payment difficulty

**Below are the Top 10 variables which are correlated.**

**Observation: df_merged_target0**

    AMT_CREDIT and AMT_GOODS_PRICE (.99)

    AMT_CREDIT and AMT_APPLICATION (.97)

    REGION_RATING_CLIENT and REGION_RATING_CLIENT_W_CITY(.95)

    DAYS_LAST_DUE and DAYS_TERMINATION (.93)

    CNT_FAM_MEMBERS and CNT_CHILDREN(.89)

    REG_REGION_NOT_LIVE_REGION and LIVE_REGION_NOT_WORK_REGION(.88)

    DEF_30_CNT_SOCIAL_CIRCLE and DEF_60_CNT_SOCIAL_CIRCLE(.87)

    REG_CITY_NOT_WORK_CITY and REG_CITY_NOT_LIVE_CITY(.84)

    AMT_ANNUITY and AMT_GOODS_PRICE(.82)

    AMT_ANNUITY and AMT_CREDIT(.82)

    AMT_ANNUITY and AMT_APPLICATION(.81)

**Observation: df_merged_target1**

    AMT_CREDIT and AMT_GOODS_PRICE(.99)

    AMT_CREDIT and AMT_APPLICATION(.97)

    REGION_RATING_CLIENT and REGION_RATING_CLIENT_W_CITY(.96)

    DAYS_LAST_DUE and DAYS_TERMINATION(.94)

    CNT_FAM_MEMBERS and CNT_CHILDREN(.89)

    REG_REGION_NOT_LIVE_REGION and LIVE_REGION_NOT_WORK_REGION(.88)

    DEF_30_CNT_SOCIAL_CIRCLE and DEF_60_CNT_SOCIAL_CIRCLE(.86)

    AMT_ANNUITY and AMT_GOODS_PRICE(.84)

    AMT_ANNUITY and AMT_CREDIT(.84)

    AMT_ANNUITY and AMT_APPLICATION(.82)

    REG_CITY_NOT_WORK_CITY and REG_CITY_NOT_LIVE_CITY(.79)

# Summary

- Social Circle, Credit rating(EXT_SOURCE), Occupation type, Education type, Family Status, Income type, Tenure of the loan are the factors to be considered during Credit risk assessment

- As we observed correlation among variables is almost SAME for Target 0 and Target 1