# Credit EDA Case Study

DYUTIMAYA DAS

SHRUTI SRIVASTAVA

# Table of Contents

# Business Objectives

- Identify Patterns indicating clients facing difficulty paying their loans. (this will help in denying the loan, reducing the loan amount, increasing the interest rate)

- Consumers capable of repaying the loan is NOT REJECTED

- Consumers incapable of repaying the loan are REJECTED

- Driving factors/variable behind the loan default.

- Identify and handle missing data

- Identify outliers

- Identify the data imbalance

# Dataset :

➢ **ApplicationData** - contains all the information of the client at the time of application. The data is about whether a **client has payment difficulties.**

➢ **Previous_application** 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.**

➢ **Columns_description** is data dictionary which describes the meaning of the variables.

# Loading the Application Data using the read_csv command in pandas

```
In [3]:  #Load and read the CSV file
         df_application=pd.read_csv('application_data.csv')
         pd.set_option('display.max_columns', 122)

In [4]:  # Show first five rows
         df_application.head()
```

| | SK_ID_CURR | TARGET | NAME_CONTRACT_TYPE | CODE_GENDER | FLAG_OWN_CAR | FLAG_OWN_REALTY | CNT_CHILDREN | AMT_INCOME_TOTAL | AMT_CREDIT | AMT_ANNUITY | AMT_GOODS_PRICE |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 100002 | 1 | Cash loans | M | N | Y | 0 | 202500.00000 | 406597.50000 | 24700.50000 | 351000.00000 |
| 1 | 100003 | 0 | Cash loans | F | N | N | 0 | 270000.00000 | 1293502.50000 | 35698.50000 | 1129500.00000 |
| 2 | 100004 | 0 | Revolving loans | M | Y | Y | 0 | 67500.00000 | 135000.00000 | 6750.00000 | 135000.00000 |
| 3 | 100006 | 0 | Cash loans | F | N | Y | 0 | 135000.00000 | 312682.50000 | 29686.50000 | 297000.00000 |
| 4 | 100007 | 0 | Cash loans | M | N | Y | 0 | 121500.00000 | 513000.00000 | 21865.50000 | 513000.00000 |

# Application Data Inspection

Using functions like shape, describe to get a view of the data along with some statistical information.

```
df_application.shape
```

```
(307511, 122)
```

```
df_application.describe()
```

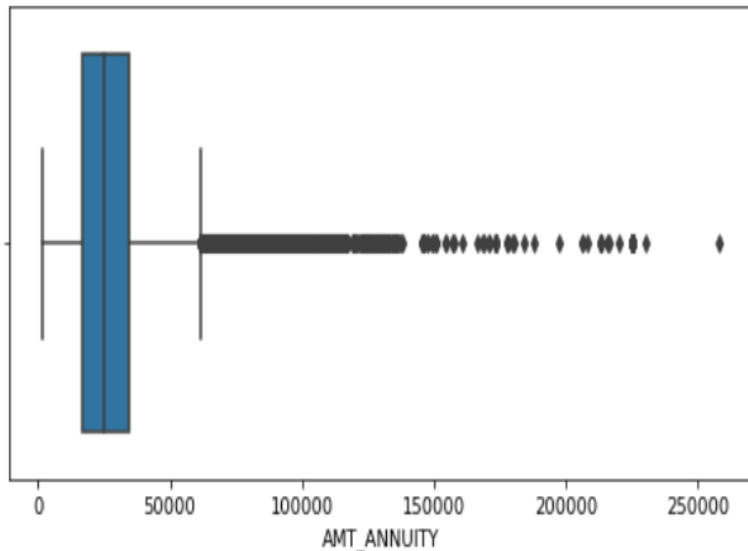|  | SK_ID_CURR | TARGET | CNT_CHILDREN | AMT_INCOME_TOTAL |
|---|---|---|---|---|
| count | 307511.00000 | 307511.00000 | 307511.00000 | 307511.00000 |
| mean | 278180.51858 | 0.08073 | 0.41705 | 168797.91930 |
| std | 102790.17535 | 0.27242 | 0.72212 | 237123.14628 |
| min | 100002.00000 | 0.00000 | 0.00000 | 25650.00000 |
| 25% | 189145.50000 | 0.00000 | 0.00000 | 112500.00000 |
| 50% | 278202.00000 | 0.00000 | 0.00000 | 147150.00000 |
| 75% | 367142.50000 | 0.00000 | 1.00000 | 202500.00000 |
| max | 456255.00000 | 1.00000 | 19.00000 | 117000000.00000 |

# Data Cleaning

- **Handling Missing values** – Check the null values using the isnull python function and drop the columns with more that 50% of missing data.

- **Checking for Outliers** - Use the describe and box plot function to check the presence of outliers.

- **Removing Columns irrelevant to the analysis** – Based on the understanding of the data columns description remove the columns which are not important for our analysis.

- **Binning Data –** Use the qcut function to put data into bins for easier analysis (AGE, INCOME_RANGE)

Missing Data, outliers can reduce the statistical power of a study and can produce biased estimates, leading to invalid conclusions.

# Data Cleaning

## Handling Outliers

```
# AMT_ANNUITY Column contains so less missing values lets checck
plt.figure(figsize=[8,4])
sns.boxplot(appdata['AMT_ANNUITY'])
plt.show()
```



## Binning Data

```
#Checking new Binned Variable"INCOME_RANGE"

df_application['INCOME_RANGE'].value_counts()
```

```
High          83013
Low           71797
Medium        67665
Very_high     20332
Name: INCOME_RANGE, dtype: int64
```

## • Percentage of missing values

```
# Checking the number of null values in the var
df_missingData = df_application.isnull().sum()

# Checking the columns with missing values grea
df_missingData = df_missingData[df_missingData.
df_missingData
```

```
AMT_GOODS_PRICE              278
NAME_TYPE_SUITE            1292
OWN_CAR_AGE             202929
OCCUPATION_TYPE          96391
EXT_SOURCE_1            173378
```

```
In [22]:  # Listing out the unwanted columns and dropping the columns to analyse the data better
not_required=['FLAG_MOBIL', 'FLAG_EMP_PHONE', 'FLAG_WORK_PHONE', 'FLAG_CONT_MOBILE','CNT_FAM_MEMBERS',
        'FLAG_PHONE', 'FLAG_EMAIL','REGION_RATING_CLIENT','REGION_RATING_CLIENT_W_CITY','FLAG_EMAIL','REGION_RATING_CLIENT',
        'REGION_RATING_CLIENT_W_CITY','DAYS_LAST_PHONE_CHANGE', 'FLAG_DOCUMENT_2', 'FLAG_DOCUMENT_3','FLAG_DOCUMENT_4', 'FLAG_DOCU
        'FLAG_DOCUMENT_7', 'FLAG_DOCUMENT_8', 'FLAG_DOCUMENT_9','FLAG_DOCUMENT_10', 'FLAG_DOCUMENT_11', 'FLAG_DOCUMENT_12',
        'FLAG_DOCUMENT_13', 'FLAG_DOCUMENT_14', 'FLAG_DOCUMENT_15','FLAG_DOCUMENT_16', 'FLAG_DOCUMENT_17', 'FLAG_DOCUMENT_18',
        'FLAG_DOCUMENT_19', 'FLAG_DOCUMENT_20', 'FLAG_DOCUMENT_21','AMT_REQ_CREDIT_BUREAU_HOUR','AMT_REQ_CREDIT_BUREAU_DAY','AMT_F
df_application=df_application.drop(labels=not_required,axis=1)
```

# Data Imbalance

- Splitting the **application data** set into two based on the TARGET variable, we notice that there is high imbalance between the **defaulters** and **non-defaulters**.
- The ratio of data imbalance is approx. 10.50



Non - default population 91.3%

8.7% client with payment difficulties

**Defaulters:** Client with Payment Difficulties

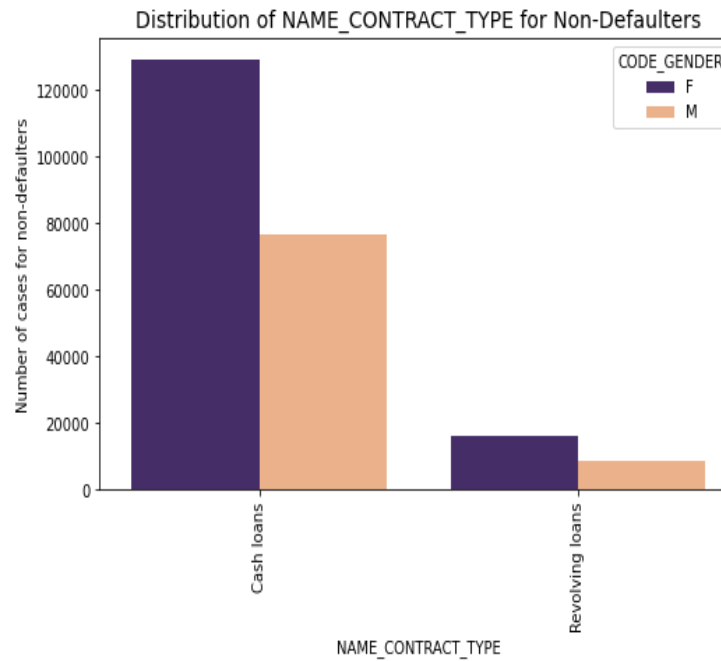**Non-Defaulters** – all others in the population.

# Univariate Analysis – Unordered Categorical Variables

**NAME_CONTRACT_TYPE :**

- Revolving loans are lesser in the defaulted population.

- Comparing both the plots it seems like the percentage of male for cash loans in defaulters are more than that of others

**NAME_INCOME_TYPE :**

- Most of the defaulters are in the working population, the reason may be this income_type apply for more loans

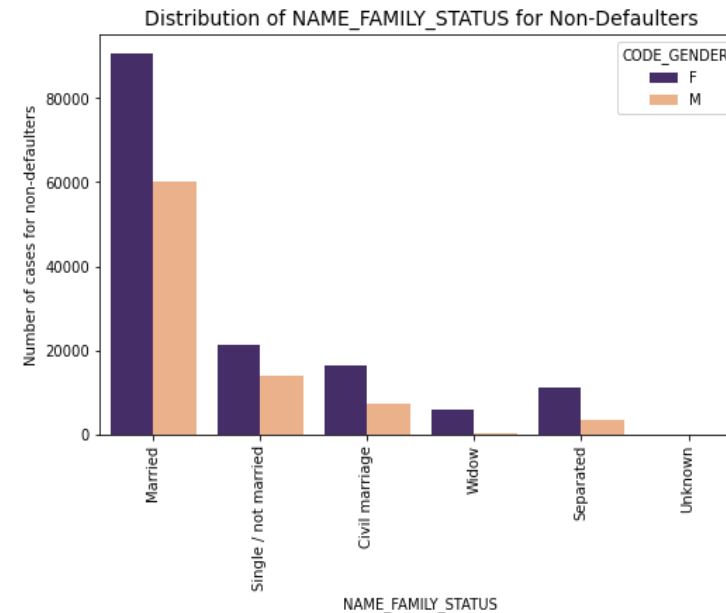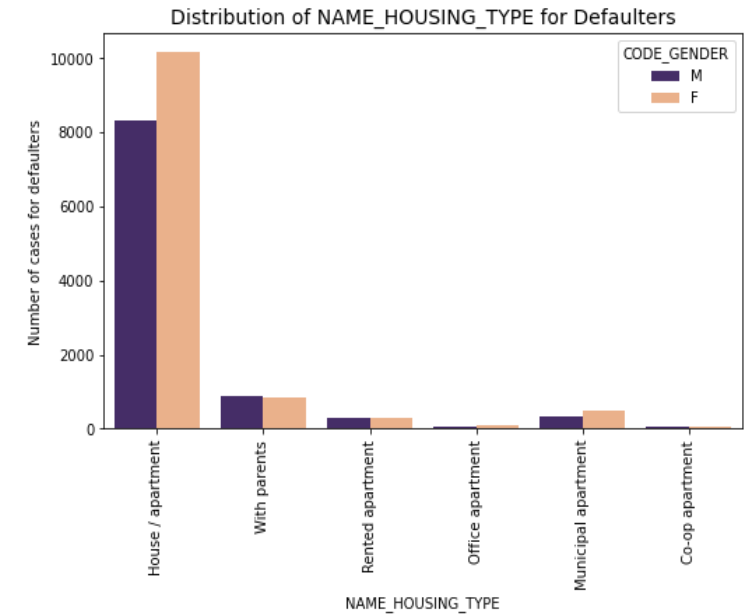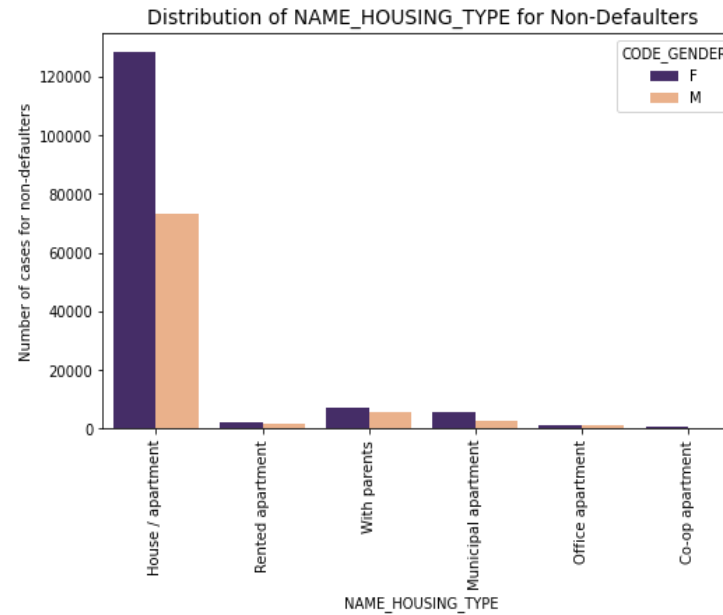- Most of the defaulters are female and belongs to working type.

# Univariate Analysis – Unordered Categorical Variables

**Name Housing  Type:**

- Population living in Rented apartments and those living with parents have higher default rate as they have higher proportion in the Defaulted population as compared to non defaulted population.

- Living in rental apartment means a cash outflow towards rent and thus less cash left for repayment of loan. Living with parents may suggest that the income is not too high and thus difficulty in repayment of loan.
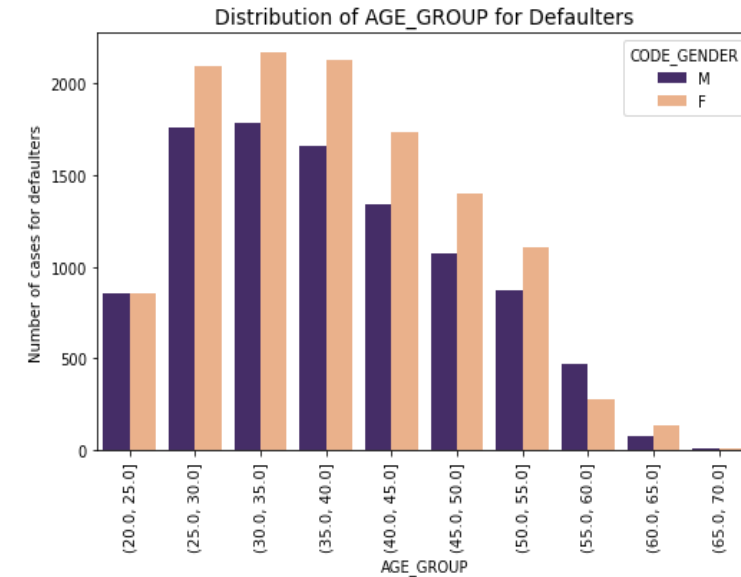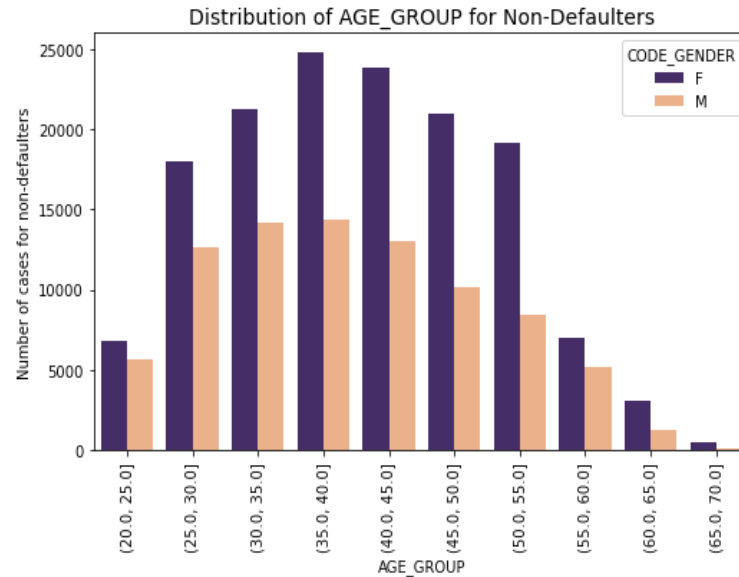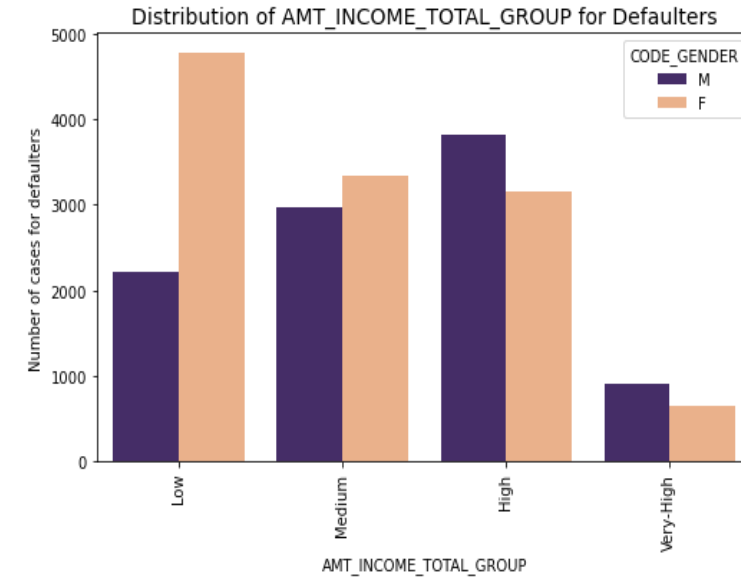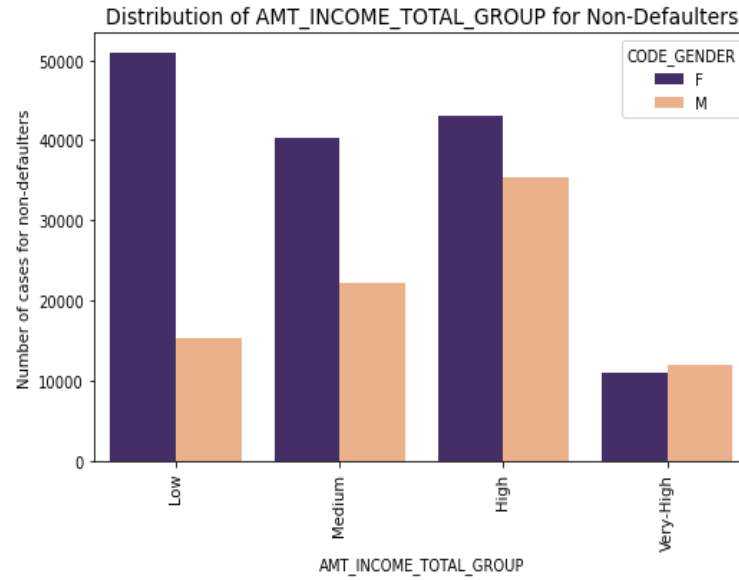
**Name Family Status:**

- Single male applicants have comparatively more defaulters than the female applicants

- Married male are more no. of defaulters in percentage than the married female

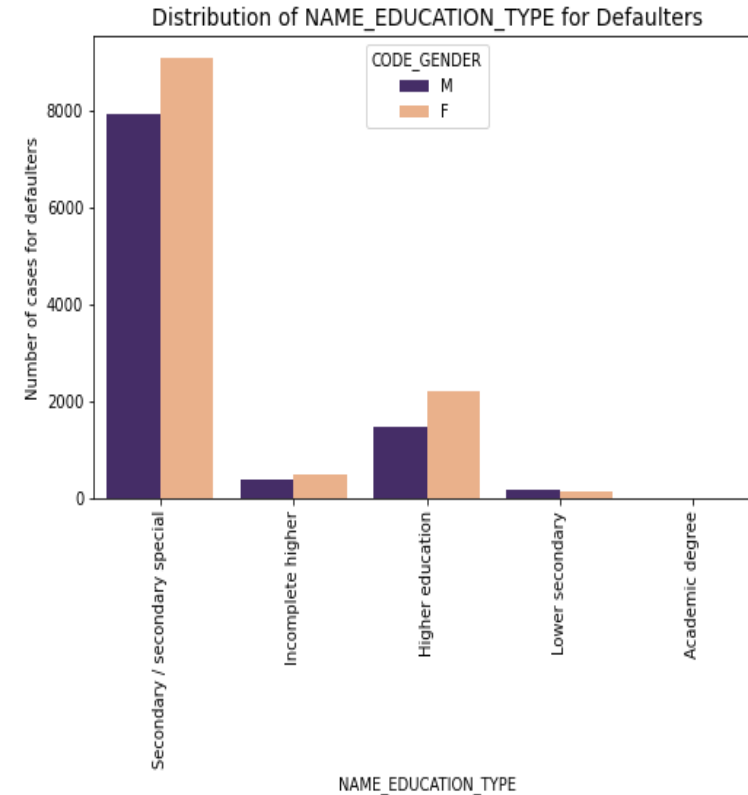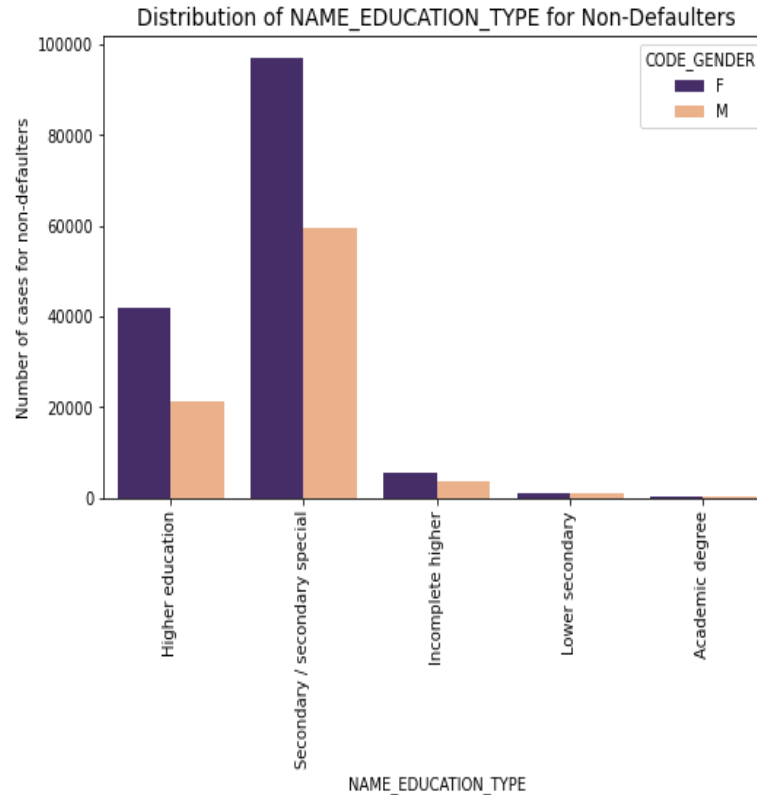# Univariate Analysis – Ordered Categorical and Continuous Variables

- **AMT_INCOME_TOTAL_GROUP** : Low income have higher defaults

- **AGE_GROUP** : Middle aged have higher number of defaulters

# Univariate Analysis – Ordered Categorical Variables

**Name Education Type**:

- Defaults from the clients with Academic degree are quite low.

- Approving application from clients having secondary/secondary special is a risk as they have more defaults. The no. of female defaulters is more than that of males.

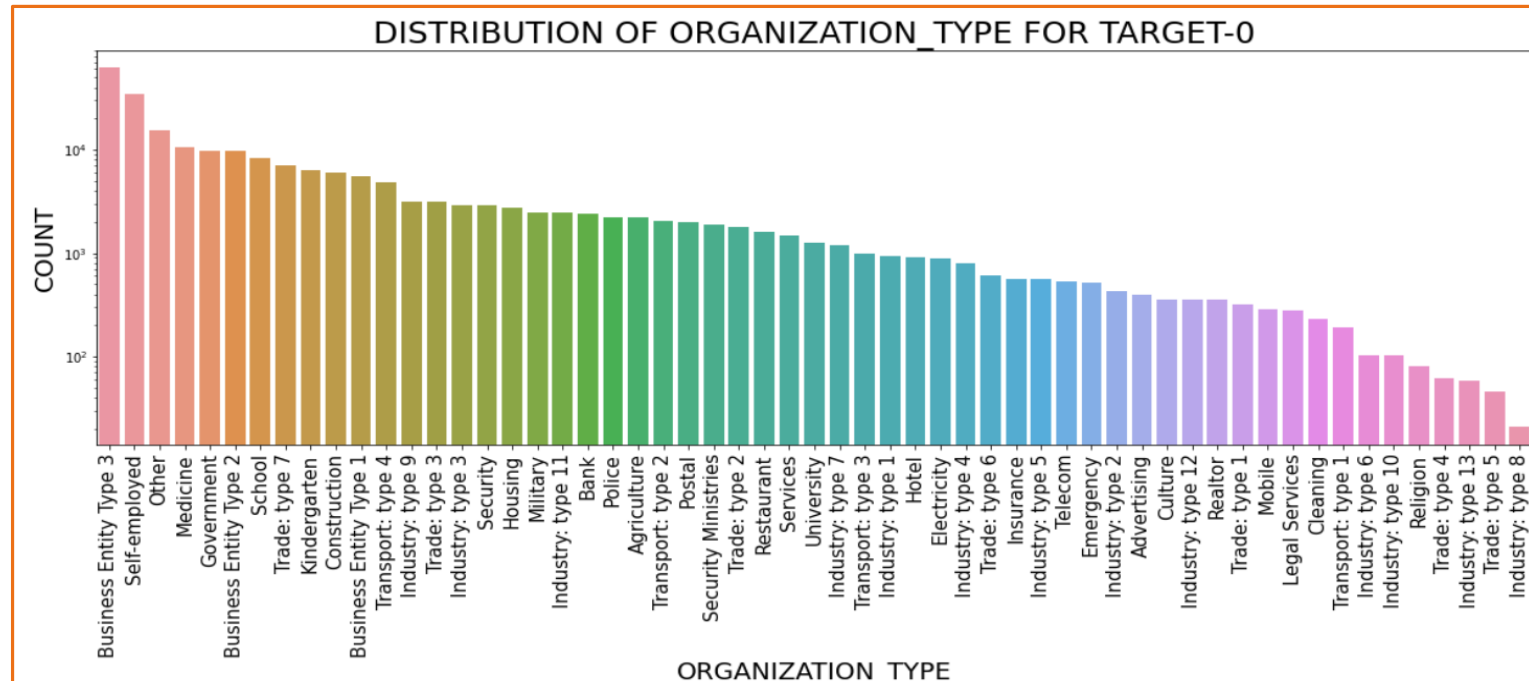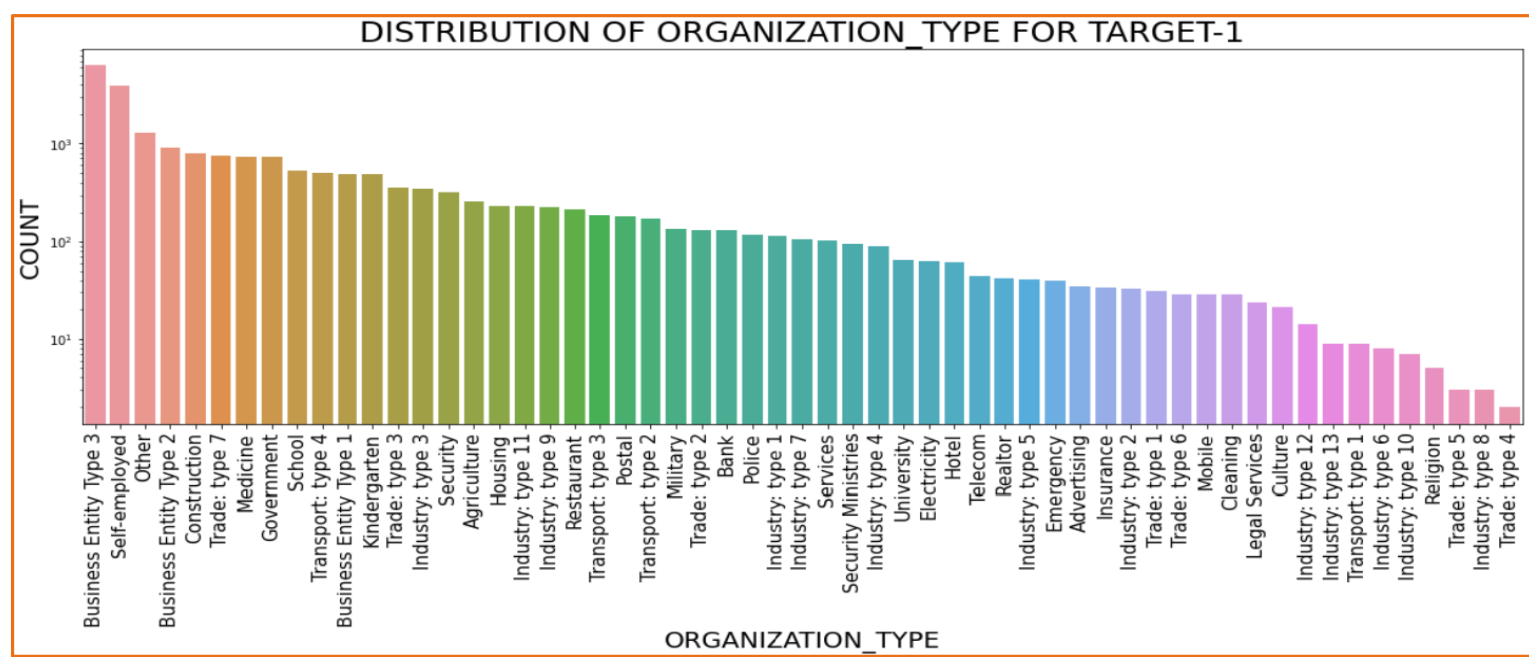- Clients having secondary/secondary special took more loan as compared to others.

# DISTRIBUTION PLOT OF ORGANIZATION_TYPE FOR TARGET 0 AND 1

**Defaulters (Target 1)**:

- Clients which have applied for credits are from most of the organization type 'Business entity Type 3', 'Self employed', 'Other', 'Medicine' and 'Government'.

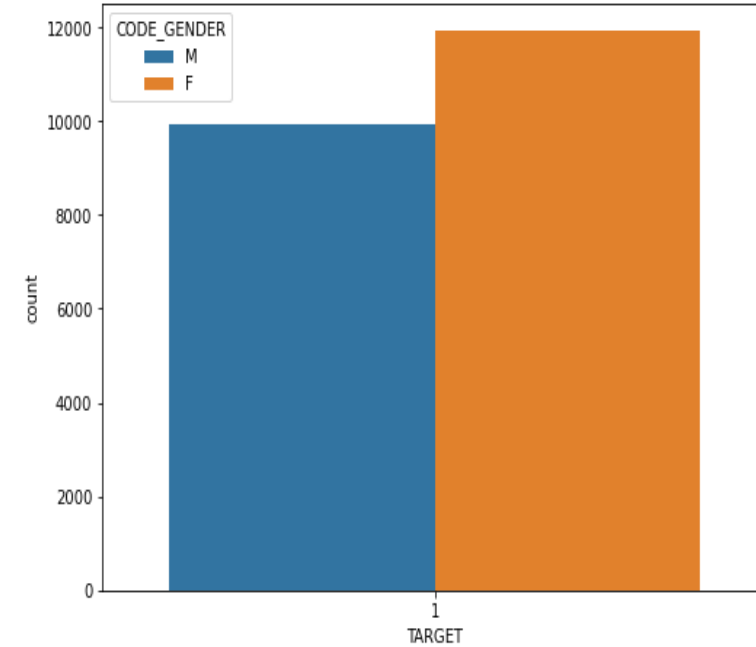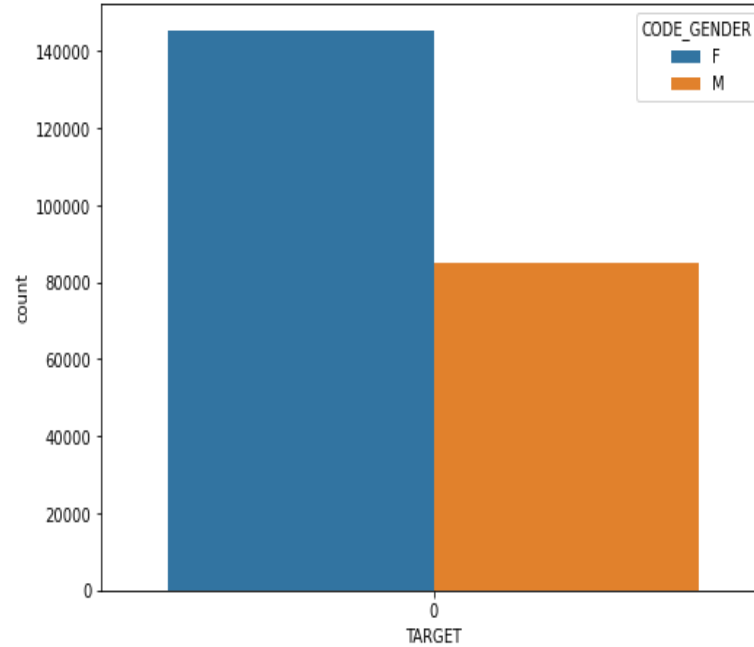- Less clients are from Industry type 8, type 6, type 10, religion and trade type 5, type 4.

**Non-Defaulters(Target 0)**:

- Most of the clients who have applied for credits are from "Business Entity Type 3", "Self-employed", "Other","Medicine", "Government", "Business Entity Type 2"

- Less number of clients are from "Industry:type 6", "Industry:type 10", "Religion", "Trade:type 4", "Industry:type 13", "Trade:type 5", "Industry:type 8"
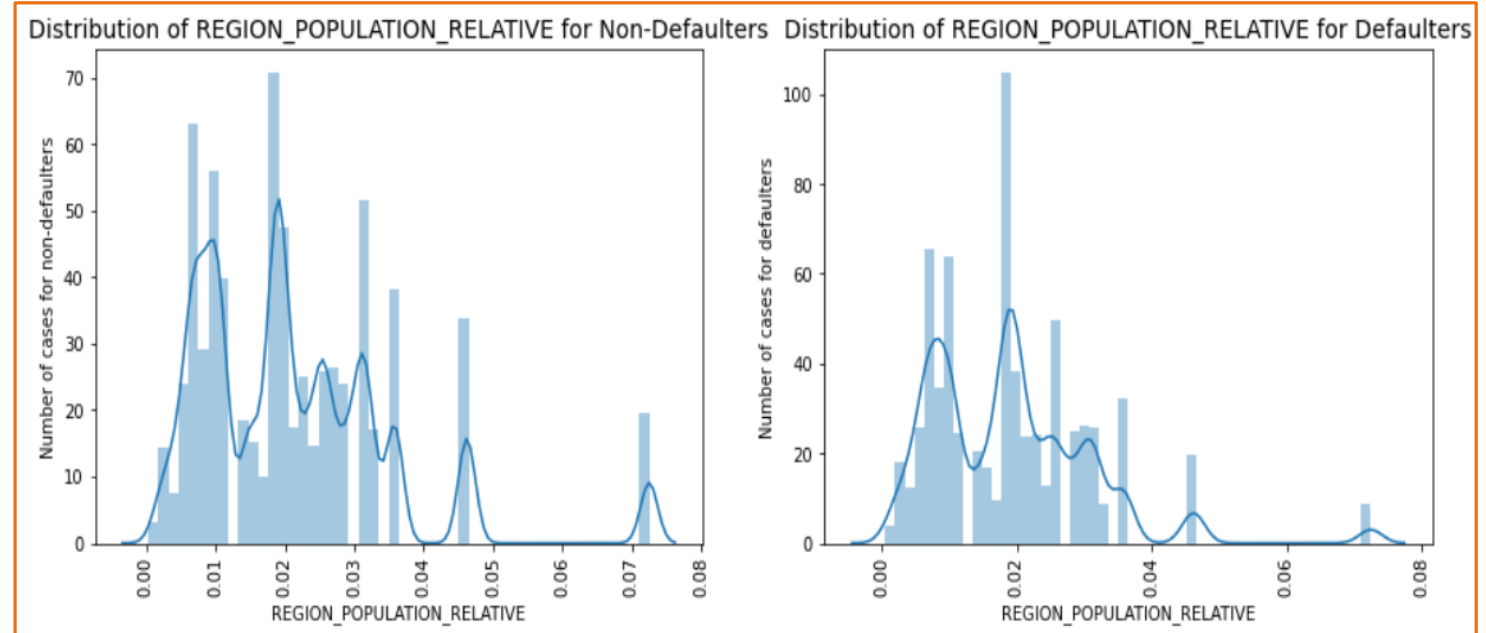
# Segmented Analysis

- The female count is more than that of male count in applicants

- The percentage of male defaulters are more than that of male non-defaulters

# Univariate Analysis – Continuous Variables

**REGION_POPULATION_RELATIVE** :

- The population relative more with respect to the no-defaulters

- The distribution is more for defaulters having region population 0.02

# Outliers Detection for Target-0
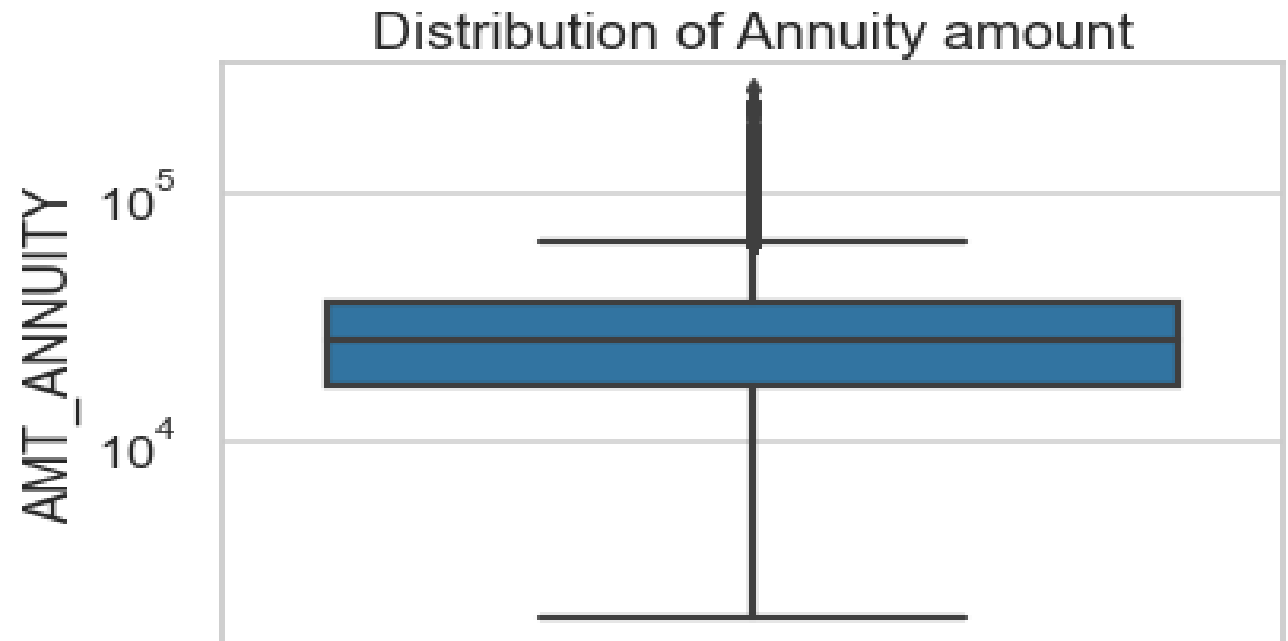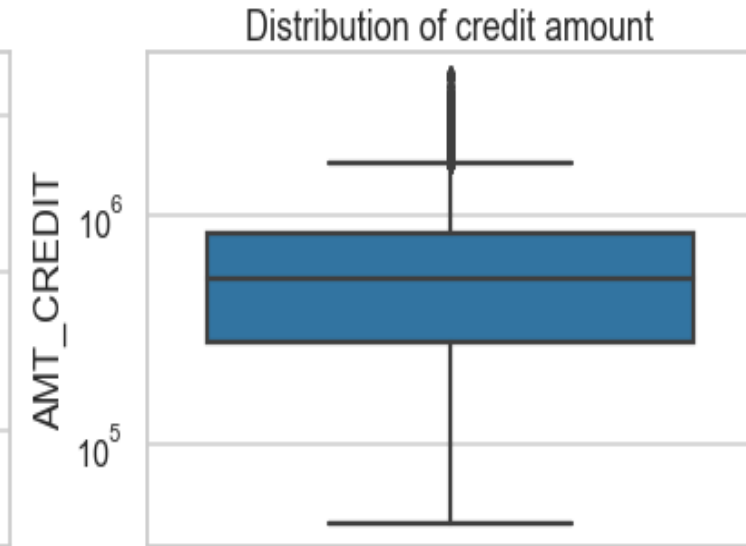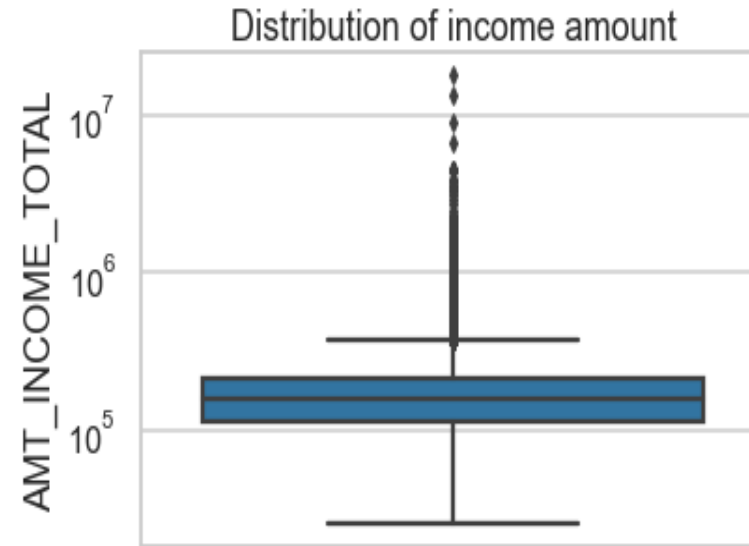
**Amt Income Total:**

- There are some outliers in the income amount column.

- The third quartiles is very slim for income amount.

**Amt Credit:**

- The column contains some outliers.

- The first quartile is bigger than third quartile for credit amount which means most of the credits of clients are present in the first quartile
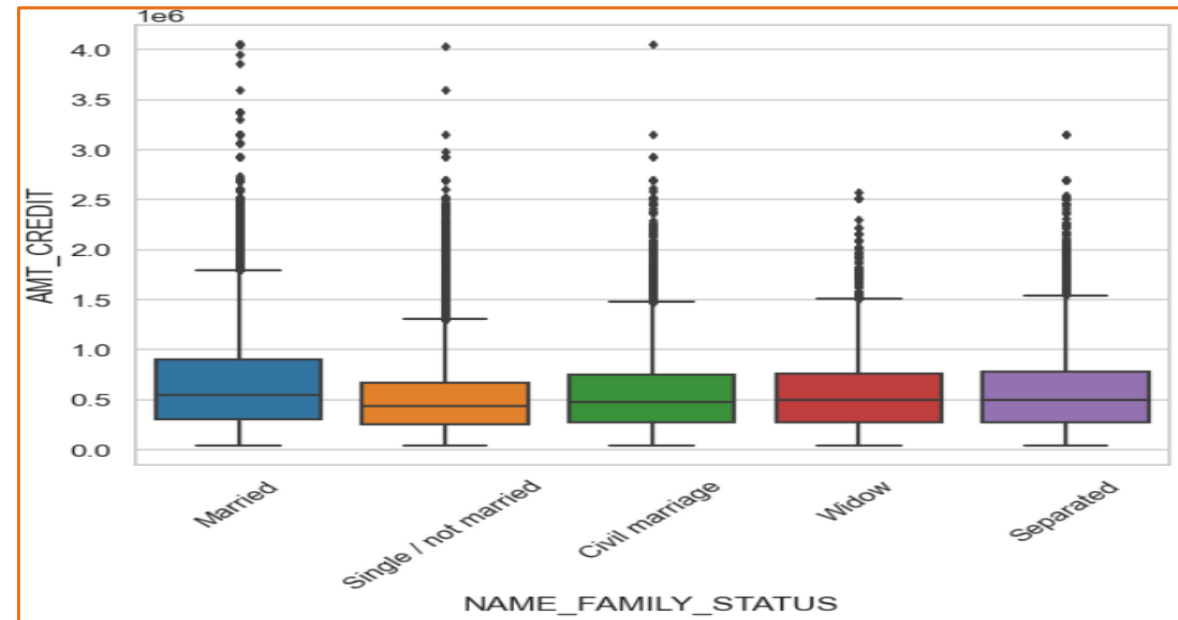
**Amt Annuity:**

- Some outliers are noticed in annuity amount.

- The first quartile is bigger than third quartile for annuity amount which means most of the annuity clients are from first quartile.



Distribution of income amount



Distribution of credit amount



Distribution of Annuity amount
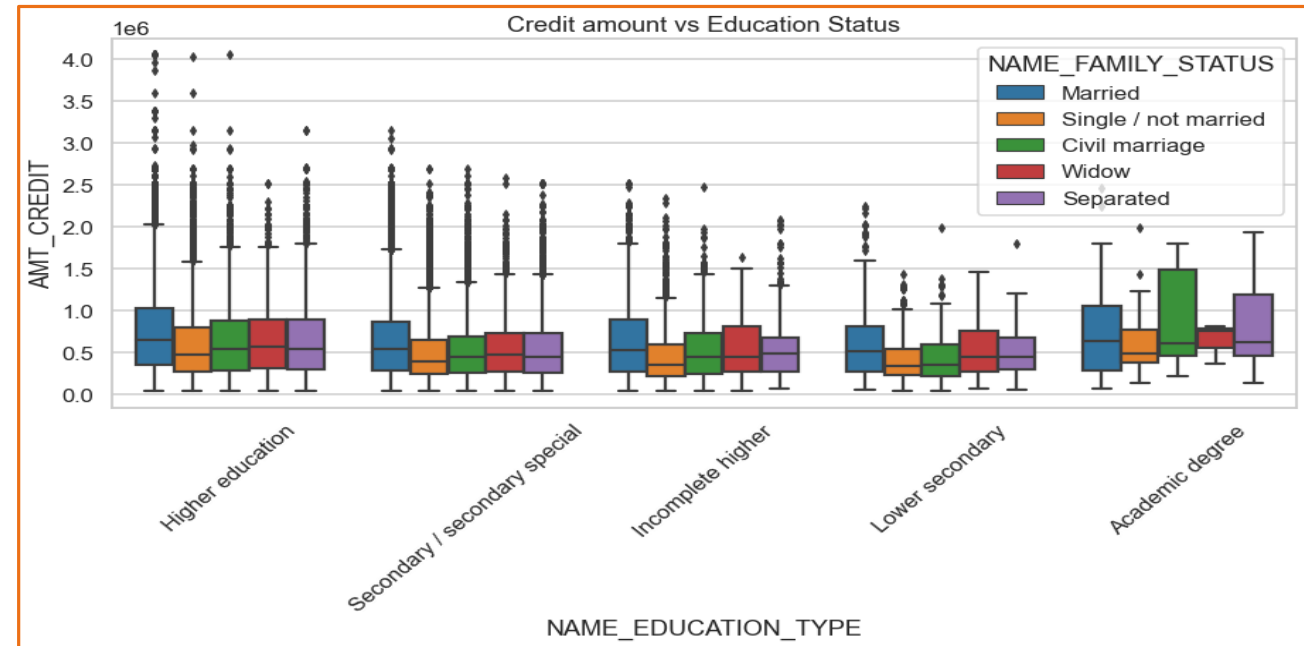
# Bi-variate Analysis for Target 0

## AMT_CREDIT & NAME_FAMILY_STATUS

- Applicants with Family Status as "Married" have higher credit amount of the loan.
- The second quartile for family status as "Marriage" is larger as compared to other family statuses.
- Outliers are present in all the statuses.

## AMT_CREDIT & NAME_EDUCATION_TYPE

- Family status of 'civil marriage', 'marriage' and 'separated' of Academic degree education are having higher number of credits than others.
- Also, higher education of family status of 'marriage', 'single' and 'civil marriage' are having more outliers.
- Civil marriage for Academic degree is having most of the credits in the third quartile.
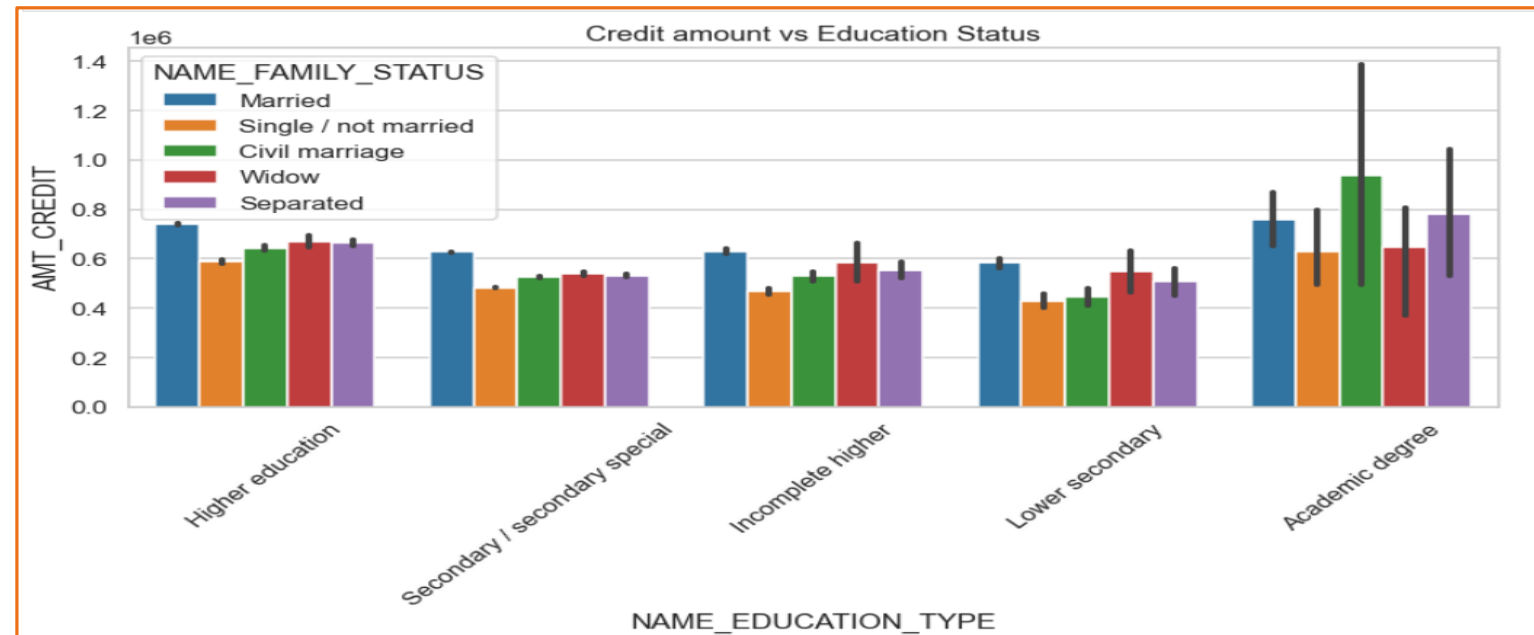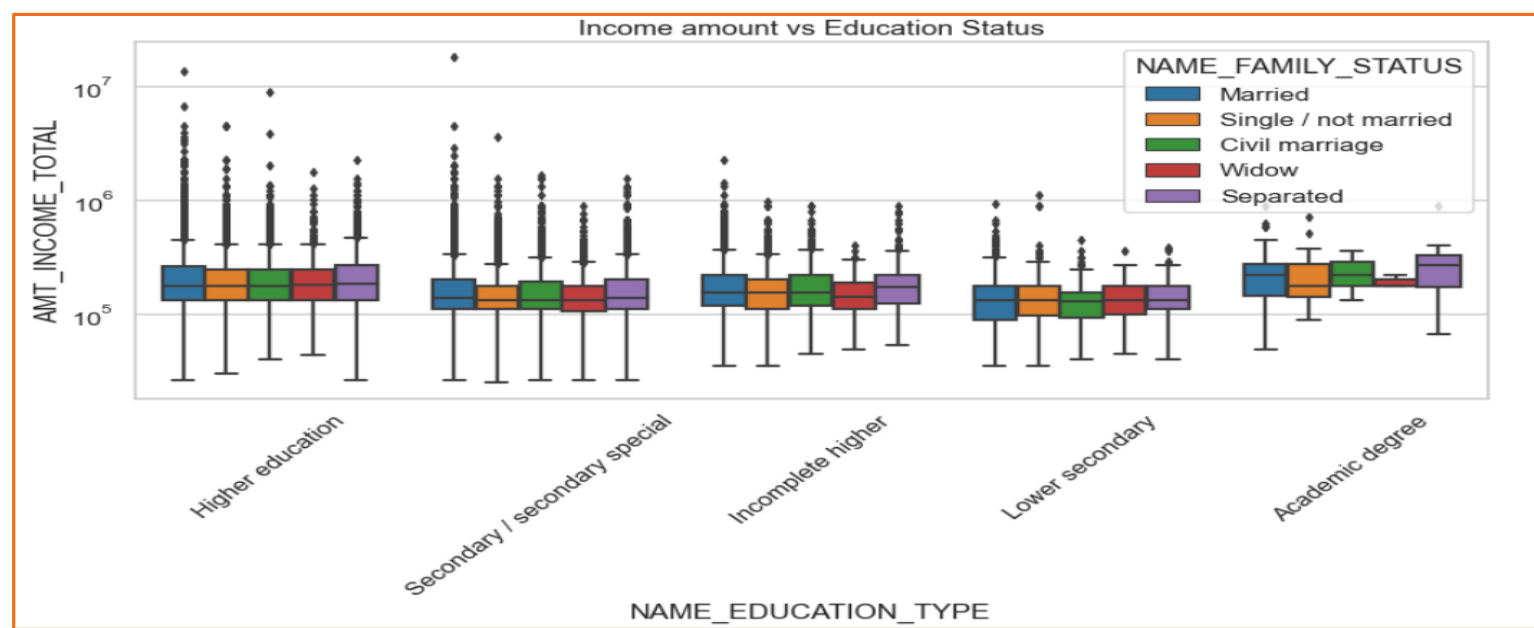
# Bi-variate Analysis for Target 0 (contd..)

**AMT_INCOME_TOTAL & NAME_EDUCATION_TYPE**

- For Education type 'Higher education' the income amount is almost equal with the family status. It contains outliers.
- Less outliers are there for applicants with Academic degree, but their income amount is little higher that Higher education.
- Lower secondary of civil marriage family status are having" less income amount than others.
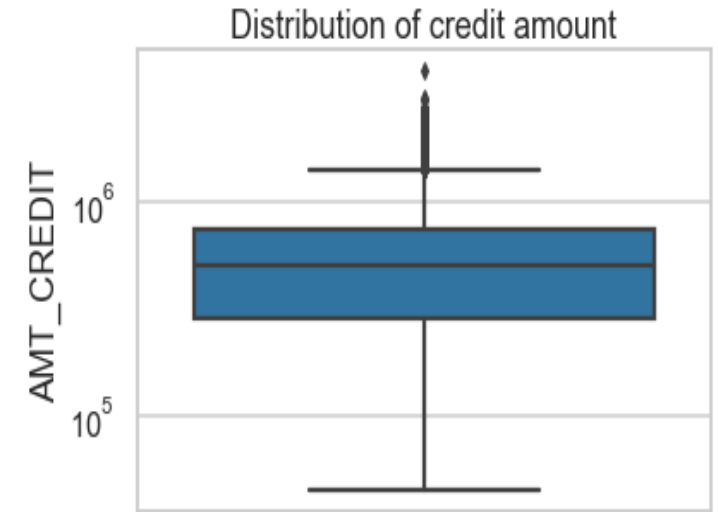
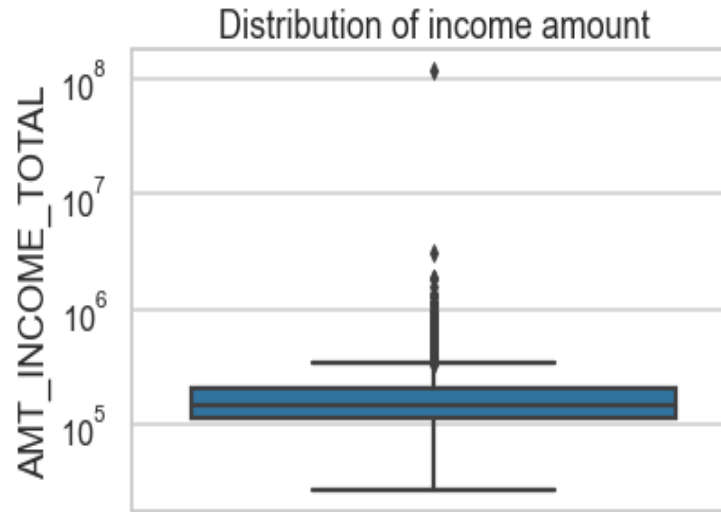**AMT_CREDIT & NAME_EDUCATION_TYPE**

- Amount of credit is high in case of Higher education and Academic degree populations married persons took more credit than others.
- Amount of credit is minimum for the populations having lower secondary education and all type of family status
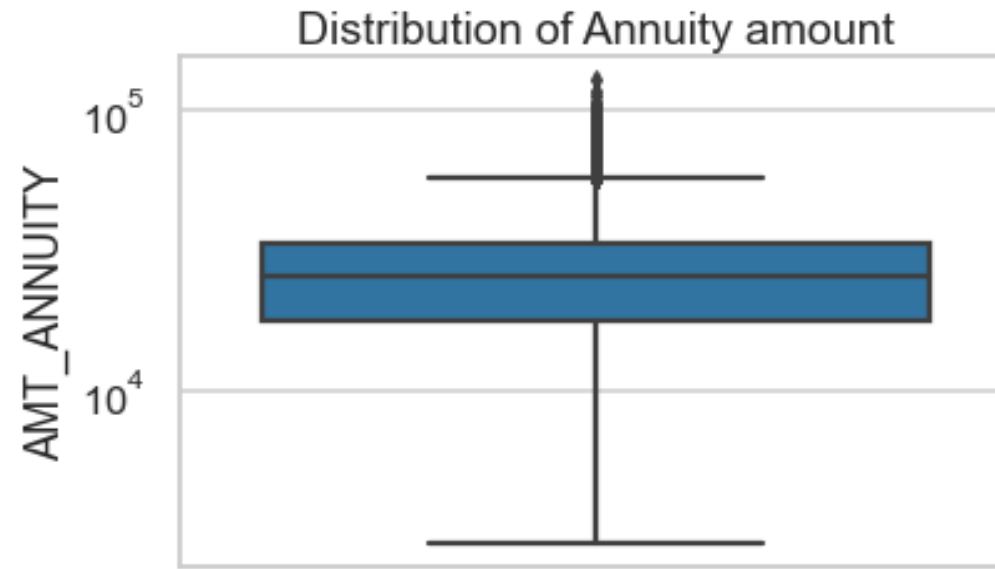
# Outliers for TARGET-1

**AMT INCOME TOTAL:**

- Some outliers are noticed in income amount.

- The third quartiles is very slim for income amount.

- Most of the clients of income are present in first quartile

**AMT CREDIT:**

- Some outliers are noticed in credit amount.

- The first quartile is bigger than third quartile for credit amount which means most of the credits of clients are present in the first quartile.

**AMT ANNUITY:**

- Some outliers are noticed in annuity amount.

- The first quartile is bigger than third quartile for annuity amount which means most of the annuity clients are from first quartile.
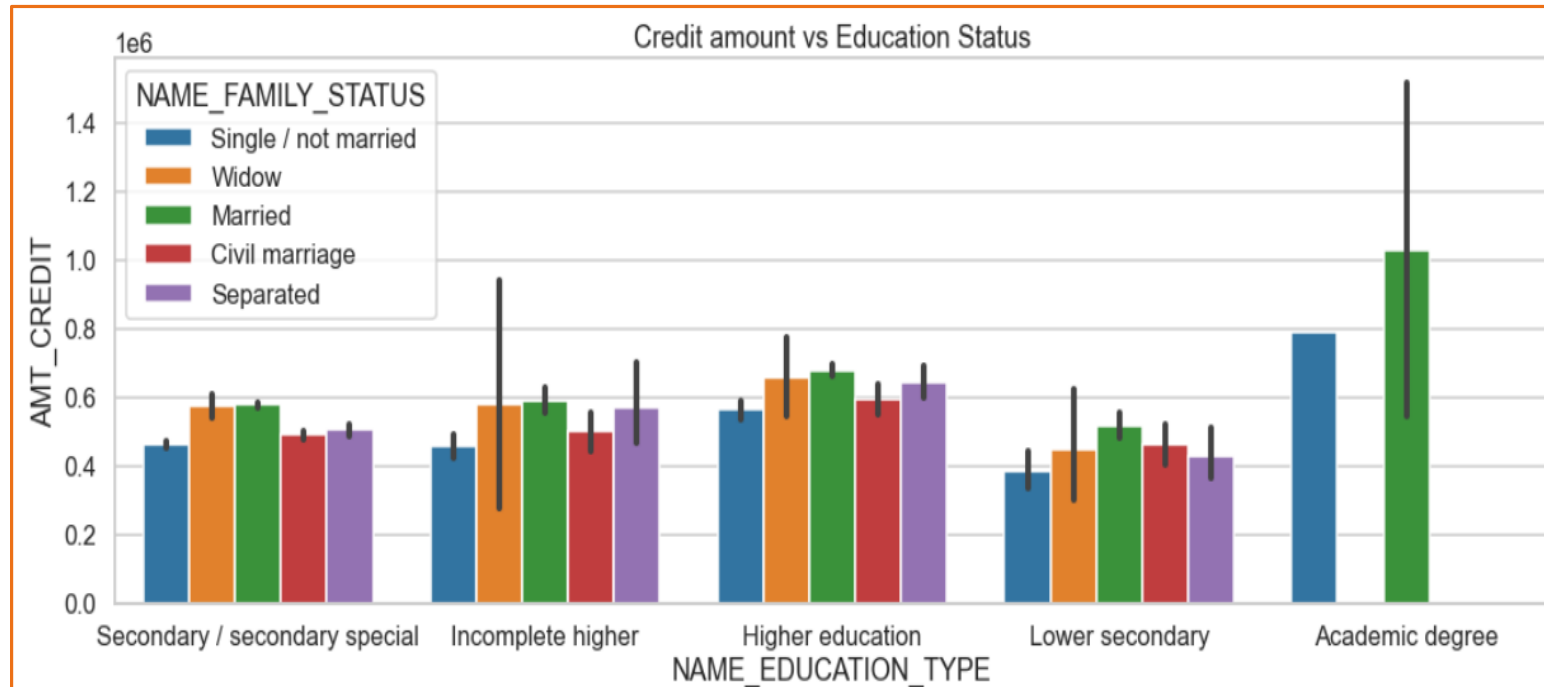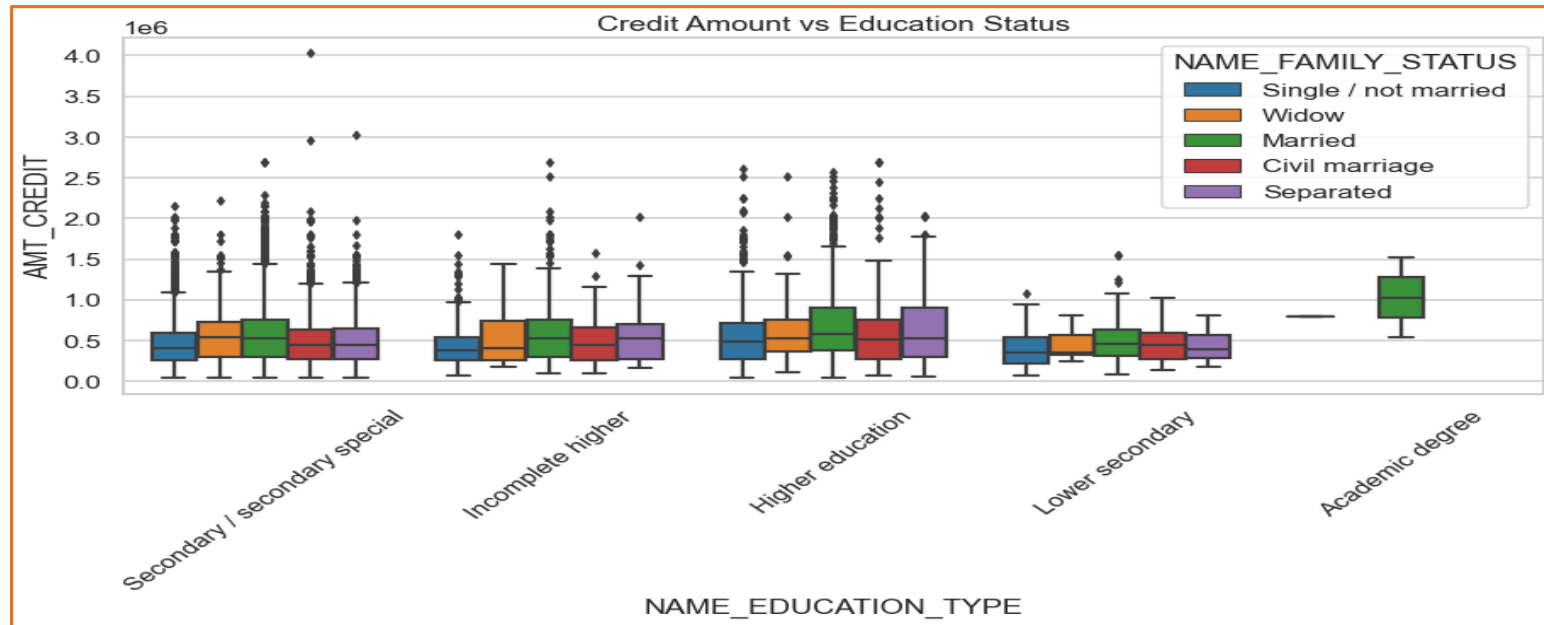
# Bi-variate Analysis for Target 1

## CREDIT AMOUNT VS EDUCATION STATUS

- In all the section of education Status married people have taken more credits.
- For Academic degree section data of only married people is available and don't have any outliers.
- Most outliers are from Secondary special , incomplete higher and higher education.
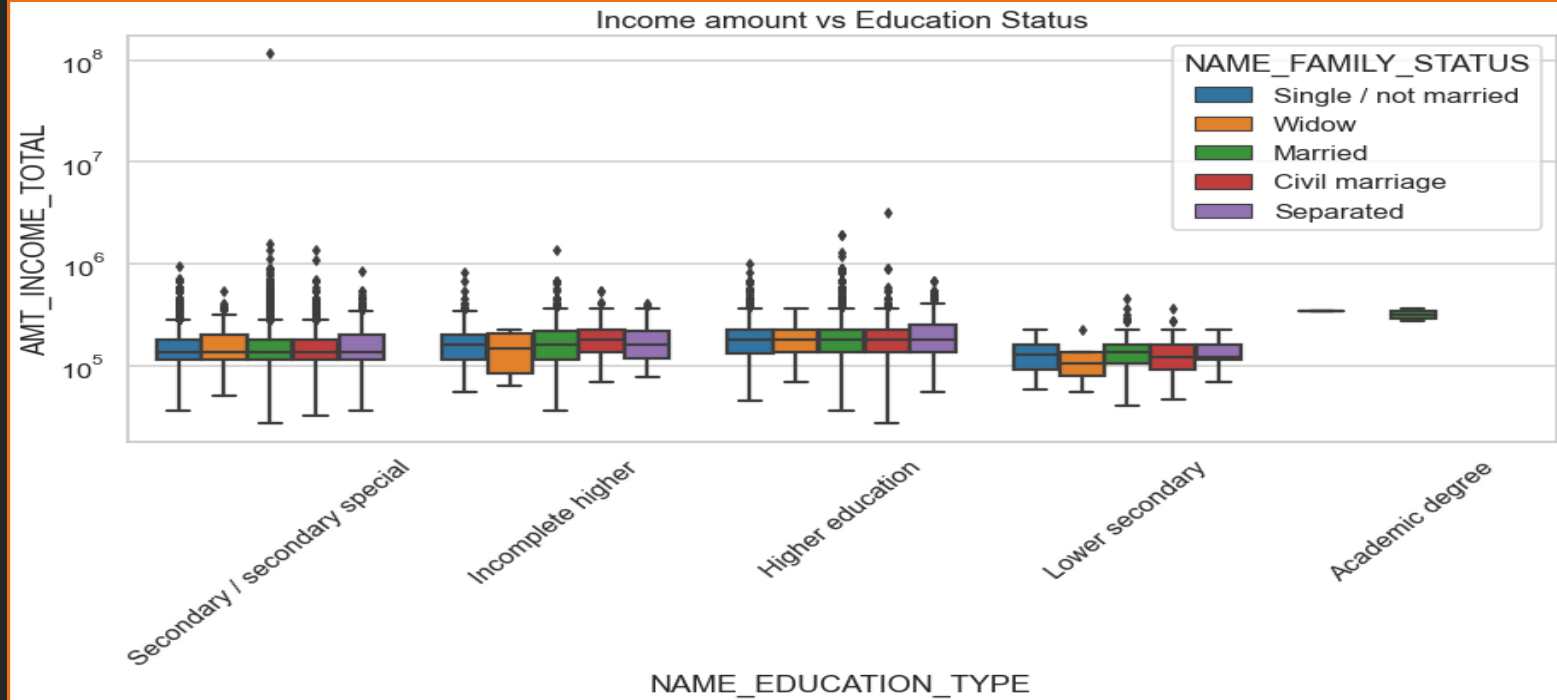
## AMT_CREDIT & NAME_EDUCATION_TYPE

- Married people have high amount credit irrespective of the education types.
- Lower Secondary have less credit amount than the others.

# Bi-variate Analysis for Target 1(contd..)

**INCOME ANOUNT vs EDUCATION STATUS**
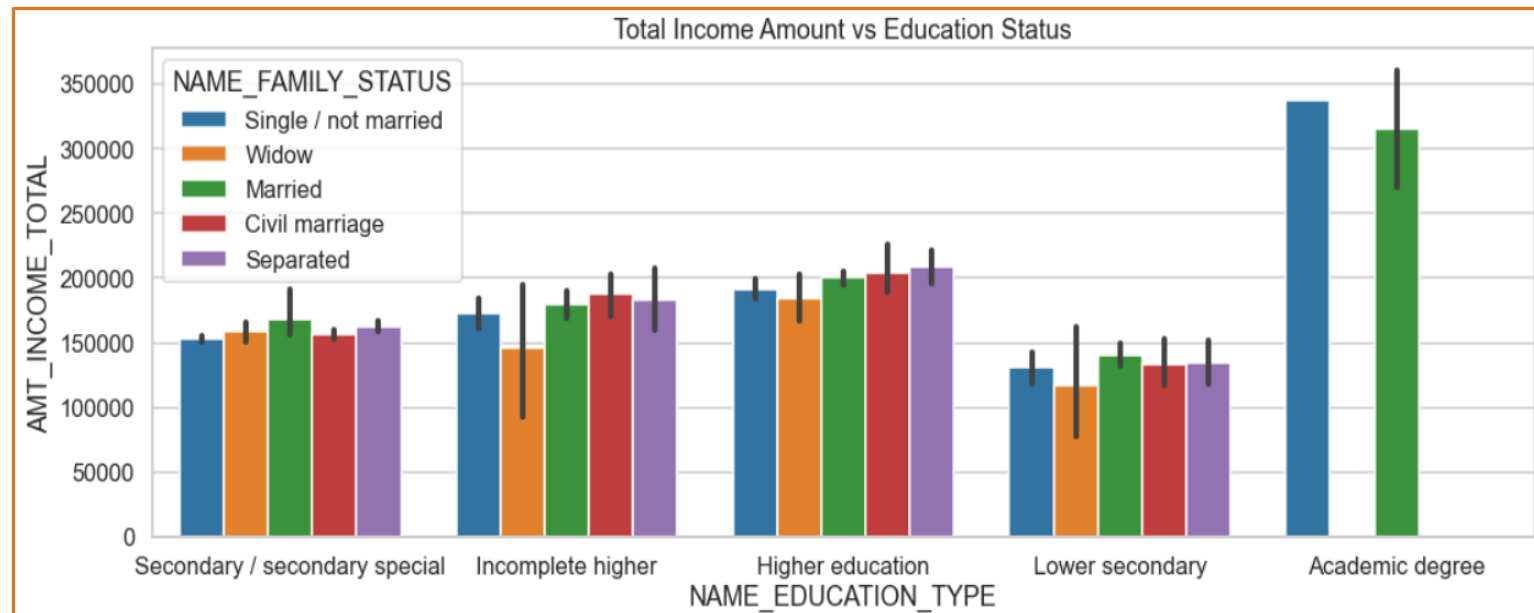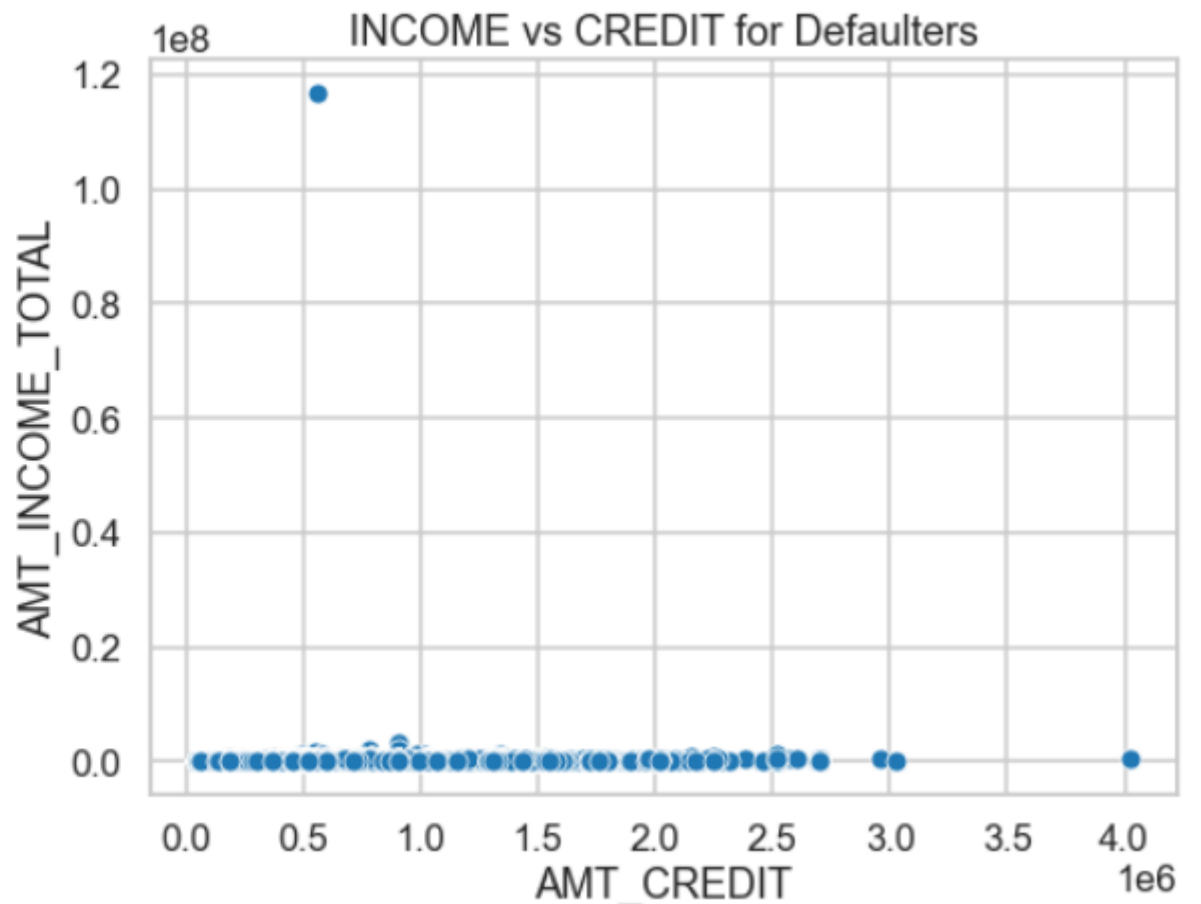
- Married Category has greater number of outliers in all the education categories and having more values in the first quartile.

- Lower secondary have less income amount than others.

**AMT_INCOME_TOTAL & NAME_EDUCATION_TYPE**

- Married and single/not married Category is having Academic Degree have highest amount of income in total

- Lower secondary have less income amount than others.

# Bivariate Analysis  - Defaulters and Non-Defaulters

- The non-defaulters took more credit as compared to their incomes.
- Defaulters have more amount of credit but their income is not high enough so its risky to give credit to the defaulters.

# Bivariate Analysis  - Defaulters and Non-Defaulters

- Defaulters are less if price of good is upto 500k and amount credit is also less than 500k
- Non-Defaulters the density is high when the amount goods is upto 1000k and the amount credit is upto 1000k

# HEAT MAP/Correlation Matrix (TARGET 0)

- Credit amount is inversely proportional to the date of birth, which means Credit amount is higher for low age and vice-versa.
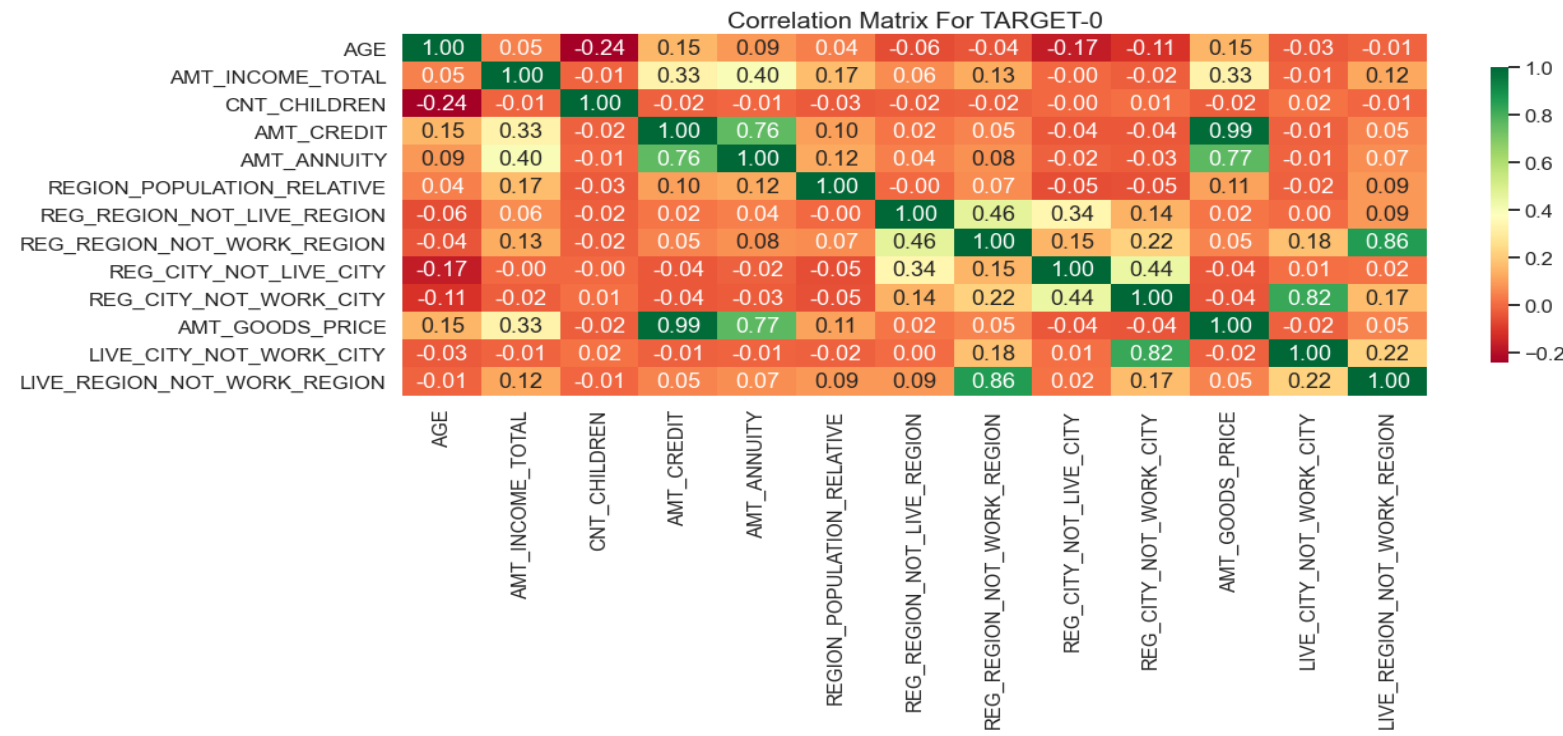- Income amount is inversely proportional to the number of children client have, means more income for less children client have and vice-versa.
- Credit amount is inversely proportional to the number of children client have, means Credit amount is higher for less children count client have and vice-versa.
- clients have less children in densely populated area.
- Credit amount is higher to densely populated area.
- The income is also higher in densely populated area.



Correlation Matrix For TARGET-0

# HEAT MAP/Correlation Matrix For TARGET 1

Most of the observations are same as of target-0 correlation heat map, but few new points observed.

- The client's permanent address does not match contact address are having less children and vice-versa.

- The client's permanent address does not match work address are having less children and vice-versa.



Correlation Matrix For TARGET-1

# Previous Application Data

- Previous application data frame is created using the file previous_application.csv.

- The data was inspected using the shape, info & column functions.

- The columns were checked for missing values. Columns with more than 50% missing values were removed

- Columns containing values like 'XNA', 'XAP' were dropped.

- A new data frame is created by combining the application_data and previous_application data.

- Columns were renamed and unwanted columns were dropped.

# Univariate Analysis

The number of refused applications is quite high.



Distribution of Contract Status

# Univariate Analysis

Diving the data set into four categories based on the contract status. We have the below observations:
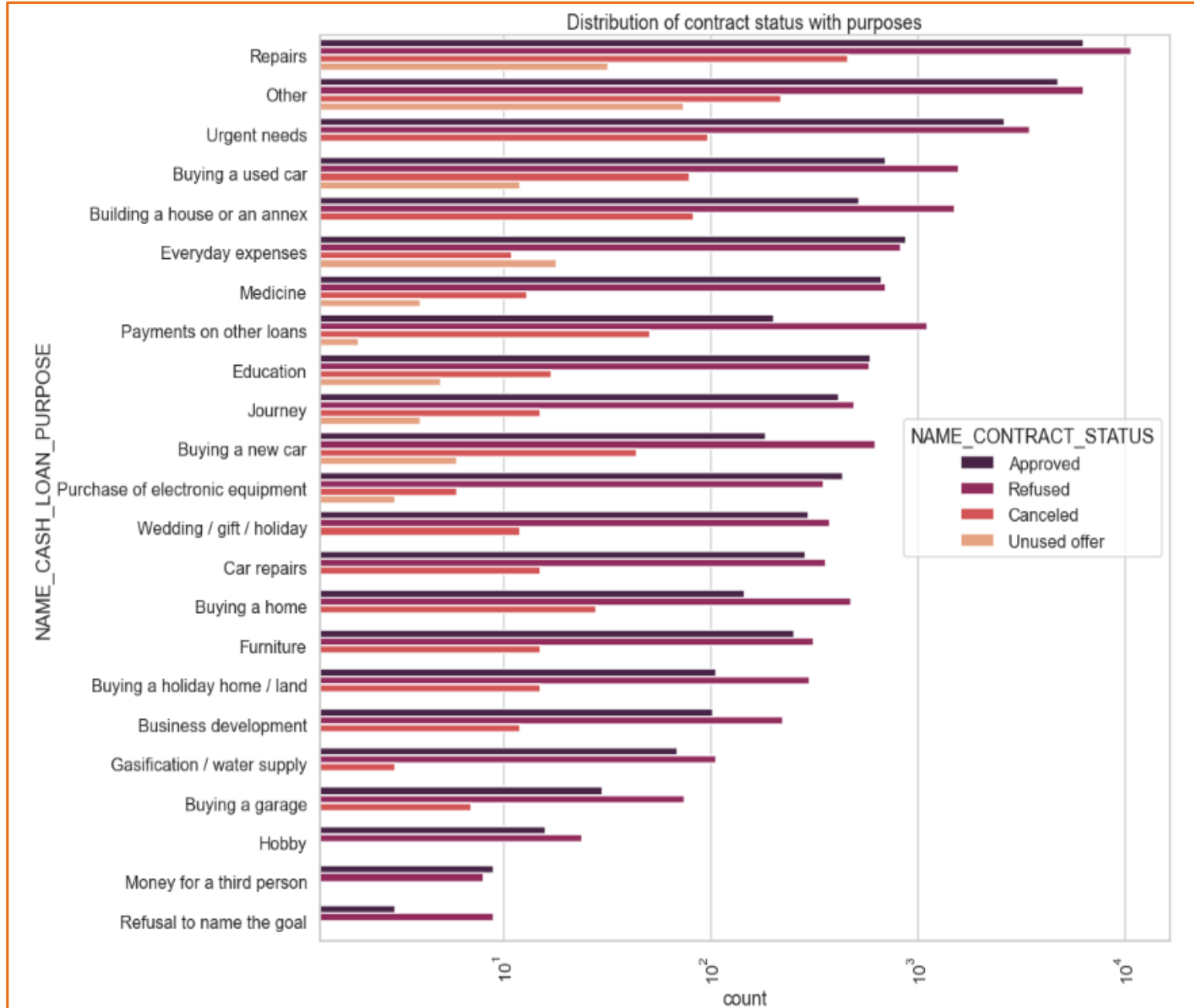
The ratio of defaulter's vs the non-defaulters in all the statuses is almost the same.

# Univariate Analysis – Contract Status & Cash Loan Purpose

- Most rejection of loans came from purpose 'repairs'.

- For Medicine purposes we have equal number of approvals and rejection.

- 'Payments on other loans', 'buying a new car' and 'buying a home' has significantly higher rejection than approvals

- Only 'Money for a third person', 'Purchase of electronic equipment' and 'Education' has higher loan approval than rejections.
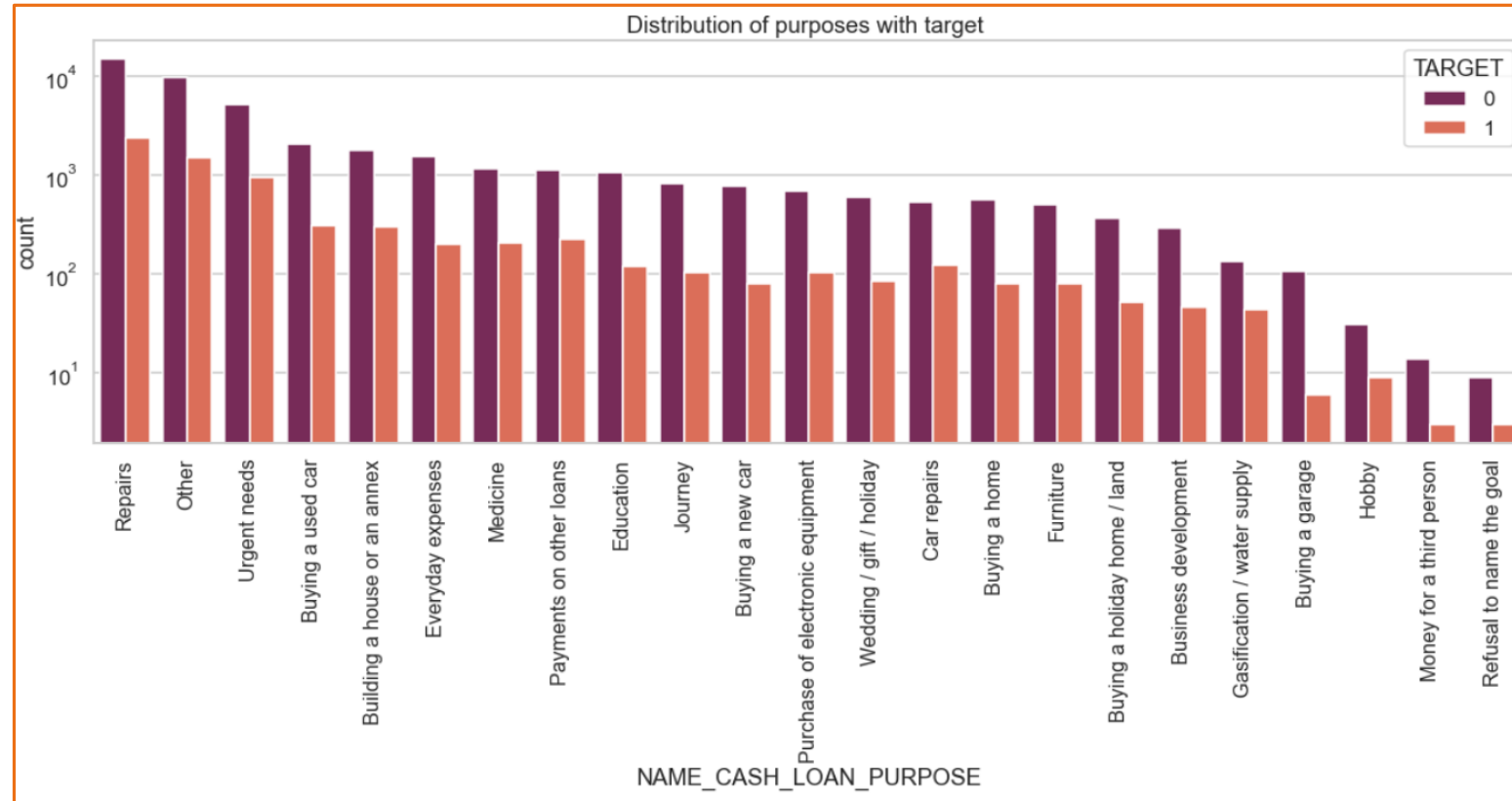
## Distribution of contract status with purposes



Distribution of contract status with purposes

# Univariate Analysis – Cash Loan Purpose & Target

- Loan purposes with 'Repairs' are facing more difficulties in payment on time.

- There are few places where loan payment is significantly higher than facing difficulties. They are 'Refusal to the name goal', 'Money for a third person', 'Buying a garage'. Hence, we can focus on the purposes for which the client has minimal payment difficulties.
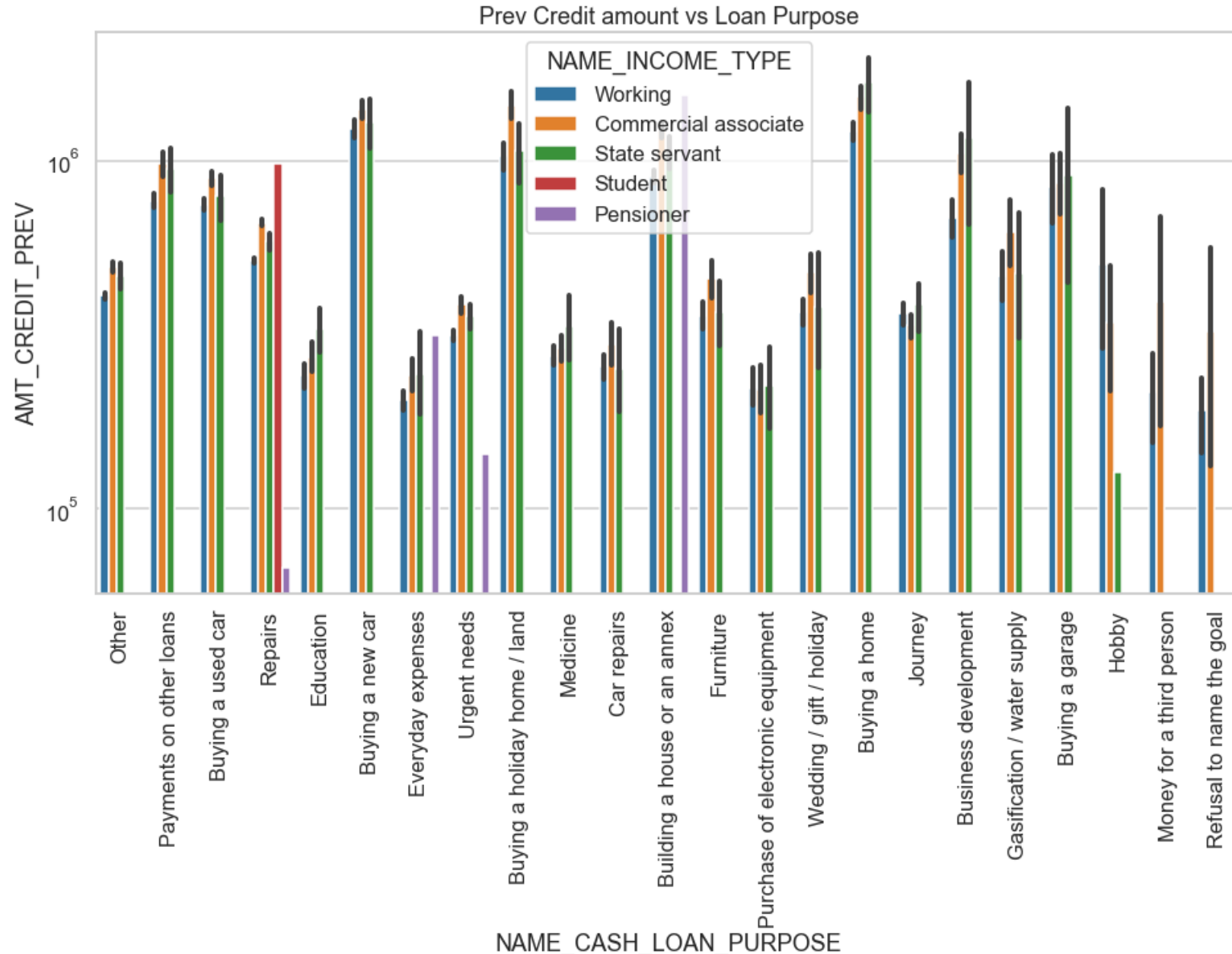
## Distribution of purposes with target

# Bi-variate Analysis – Previous Credit Amount vs Loan purpose

- The credit amount of Loan purposes like 'Buying a home', 'Buying a land', 'Buying a new car' and ' Building a house' is higher.

- Income type of state servants have a significant amount of credit applied

- 'Purchase of electronic equipment', 'Hobby', 'Refusal to the name the goal', 'Everyday expenses' is having less credits applied for.
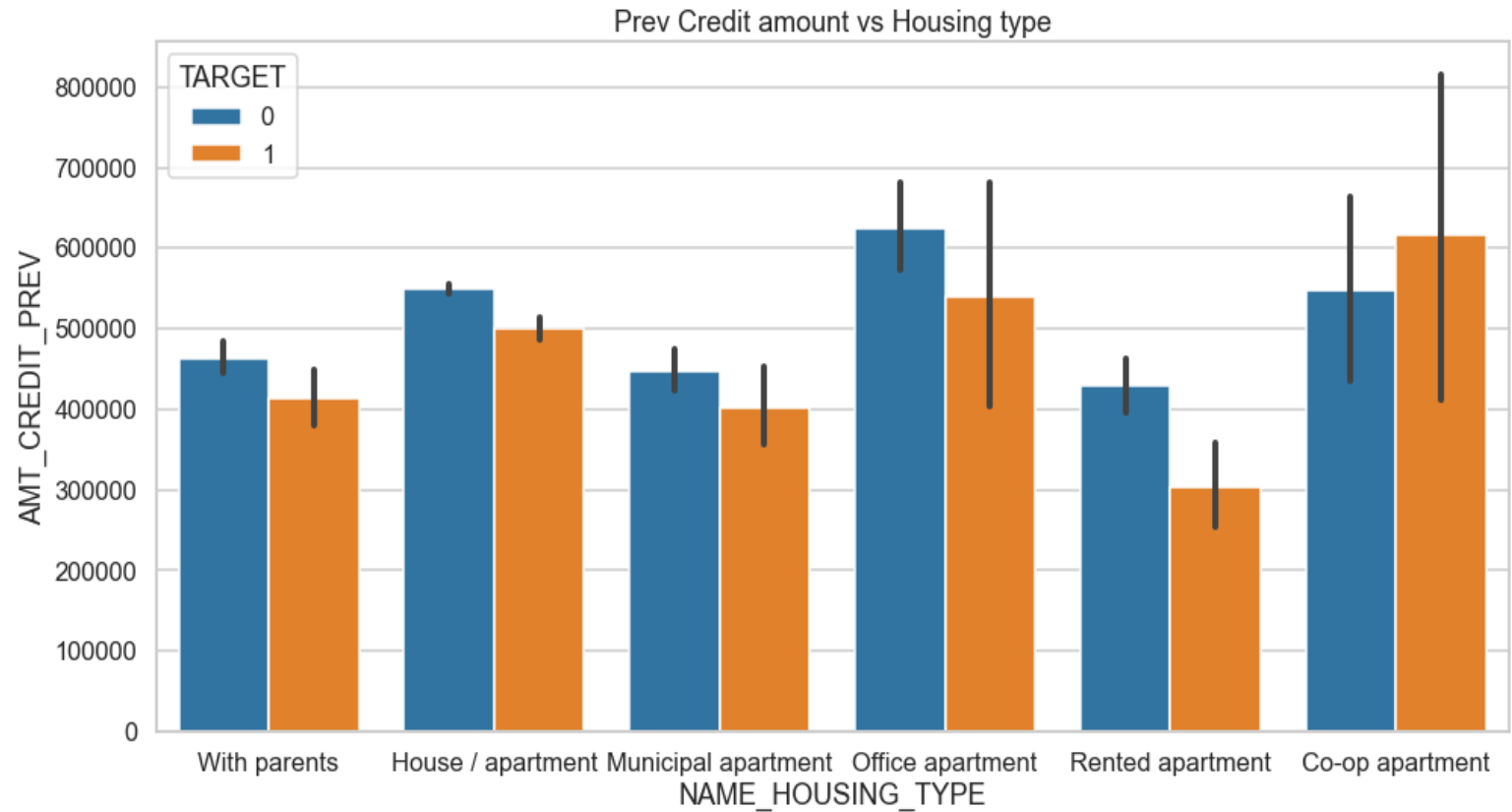


Prev Credit Amount vs Loan purpose

# Bi-variate Analysis

**Previous Credit Amount vs Housing Type**

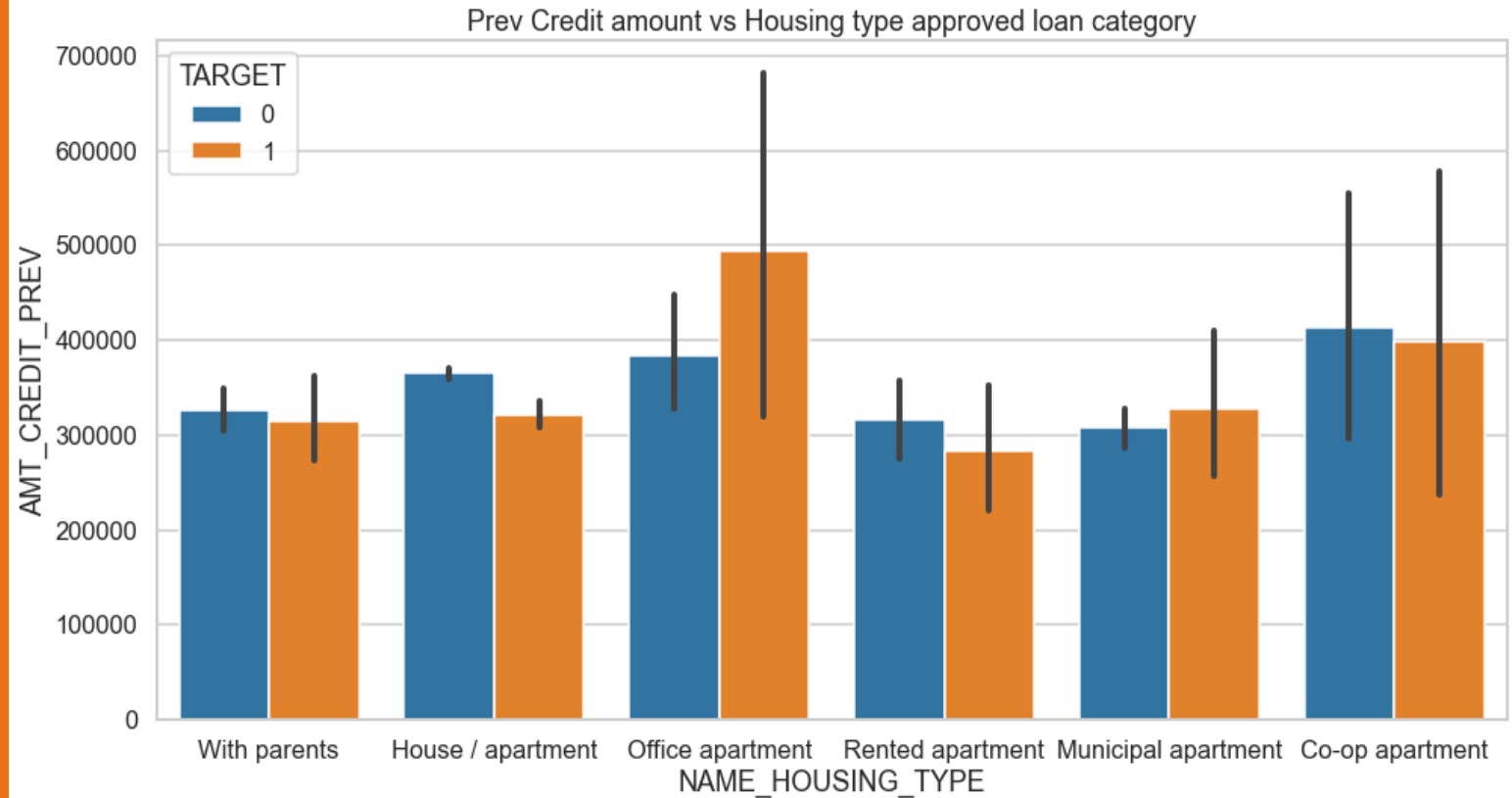- Housing Tppe as "office appartment" have higher credit for target 0.

- Co-op apartment has higher credit for target 1.

- So, we can conclude that the bank should avoid giving loans to the housing type of co-op apartment as they have payment difficulties.

- Bank can focus mostly on housing type with parents or House\appartment or municipal appartment for successful payments.



Prev Credit amount vs Housing type

# Bi-variate Analysis

**Previous Credit amount vs Housing type Approved loan category**

- Clients having office appartment have higher approved rate but must of their payments are unsucessful, so bank should not approve such loan categories.
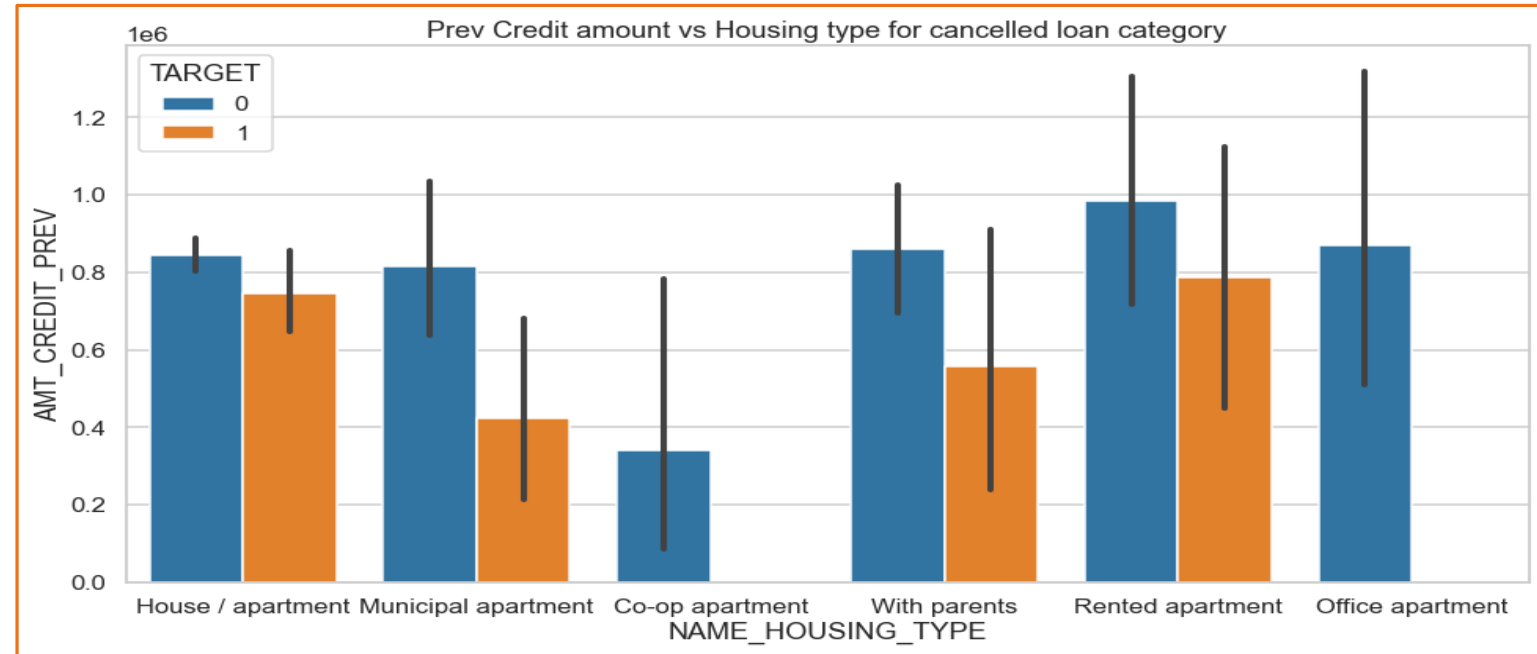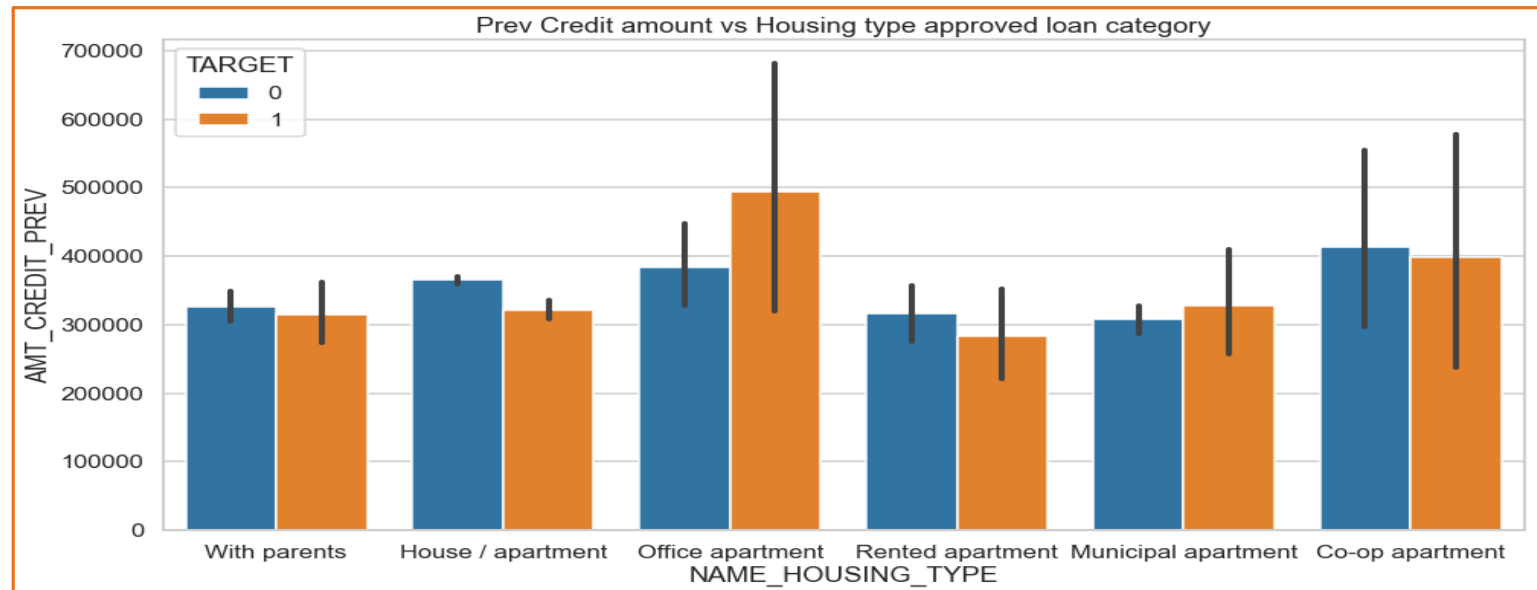


Prev Credit amount vs Housing type approved loan category

# Bi-variate Analysis

**Previous Credit amount vs Housing type for Rejected loan category**

- Bank should refuse the loan application of clients having housing type "co-op apartment" as they have more defaults.

**Previous Credit amount vs Housing type for cancelled loan category**

- Bank should target clients having housing type "with parents" as more no. of applications got cancelled from this section, but they are most successful payment clients



Prev Credit amount vs Housing type approved loan category



Prev Credit amount vs Housing type for cancelled loan category

➢ Banks should focus more on customers with contract type 'Student' ,'pensioner' and 'Businessman' with housing 'type other than 'Co-op apartment' for successful payments.

➢ There are less loan applications from customers having high salary and no repayment issues.

➢ Banks should focus less on customers with income type as 'Working' as they have a greater number of unsuccessful payments.

➢ Customers with loan purpose as "Repair" have higher number of unsuccessful payments.

➢ Get as much as clients with housing type as 'With parents' as they are having least number of unsuccessful payments.

➢ The customers whose applications were previously approved and have housing type as office apartment seem to have a greater number of unsuccessful payments.

➢ Bank should not reject the applicants having housing type 'house/apartment', 'muncipal apartment', 'with parents' as more successful payments achieved from these categories.

# Conclusion