# CREDIT EDA CASE STUDY

BY – 1. ANIKET CHAWARE

2. SHUBHAM GOUR

# INTRODUCTION

This case study aims to give us an idea of applying EDA in a real business scenario. In this case study, we develop a basic understanding of risk analytics in banking and financial services and understand how data is used to minimise the risk of losing money while lending to customers.

# PROBLEM STATEMENT

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. 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.
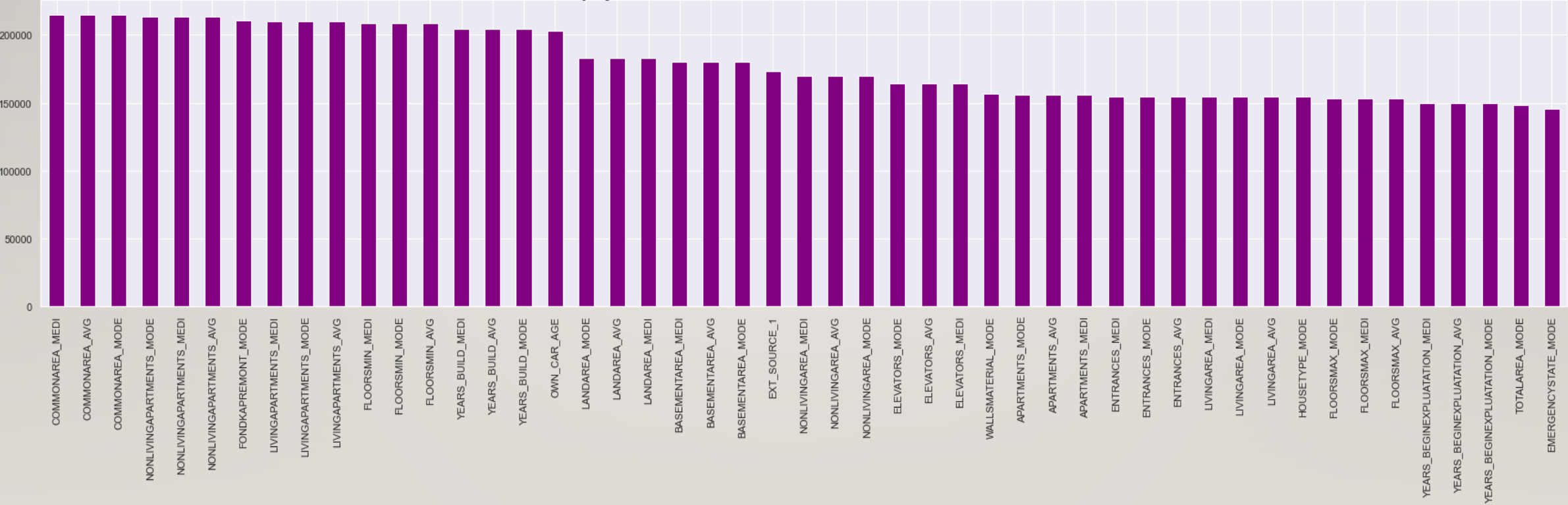
- Two types of risks are associated with the bank's decision:
  - If the applicant is likely to repay the loan, then not approving the loan results in a loss of business to the company
  - If the applicant is not likely to repay the loan, i.e. he/she is likely to default, then approving the loan may lead to a financial loss for the company.

# DATA CLEANING

# LIST OF EMPTY COLUMNS VALUES ARE MORE THAN 35%



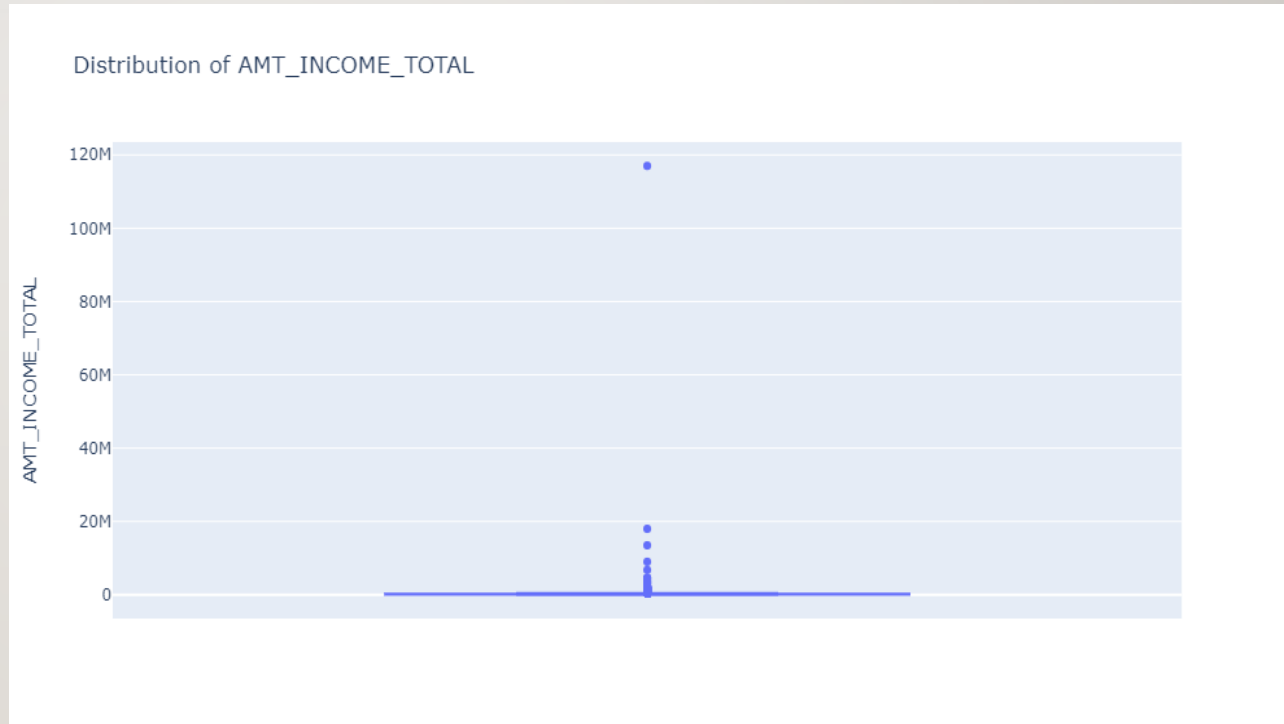List of Empty Columns counts values are more than 35%

# FINDING OUTLIERS
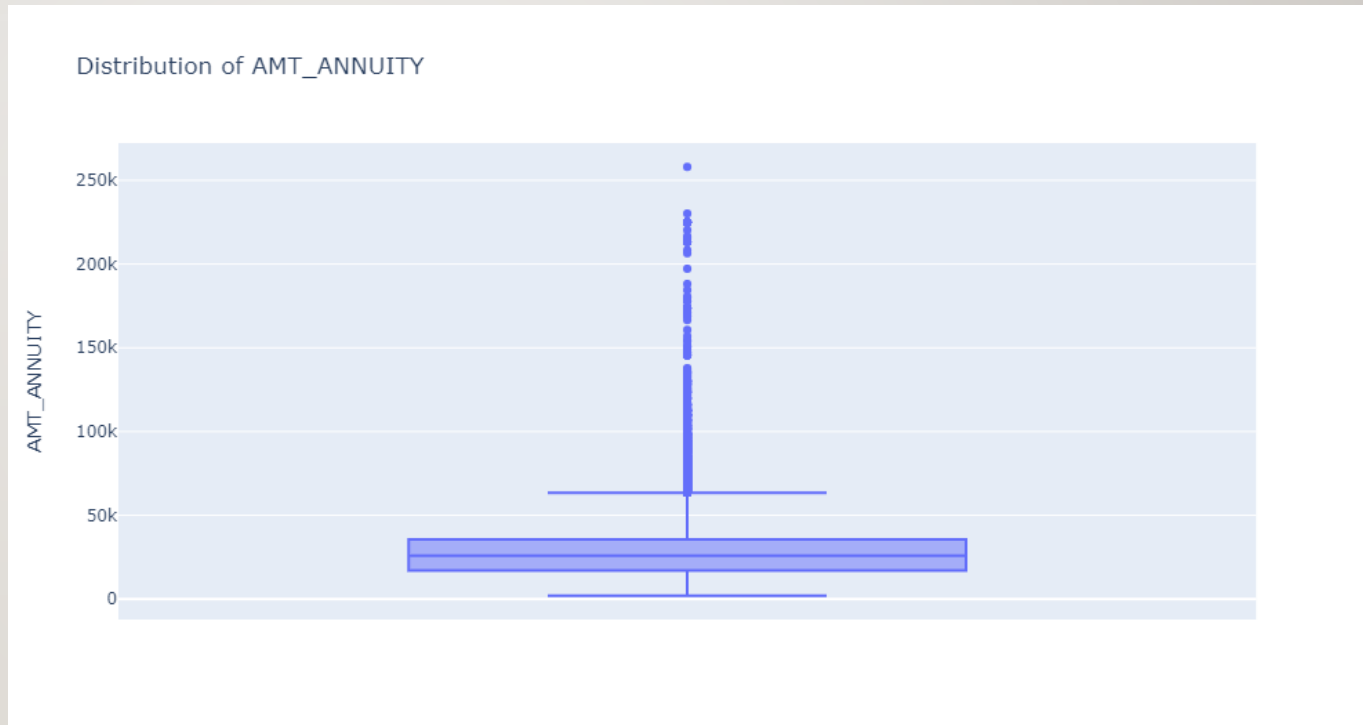
# DISTRIBUTION OF AMT_INCOME_TOTAL

Observtion :

Here, in the column 'AMT_INCOME_TOTAL' which tells us the income of the client. We observe a value of 117M which is surely an outlier.



Distribution of AMT_INCOME_TOTAL

# DISTRIBUTION OF AMT_ANNUITY

Observation:

Here, in the column 'AMT_ANNUITY' which tells the loan annuity. We observe a value which is greater that 258000 which is surely an outlier.
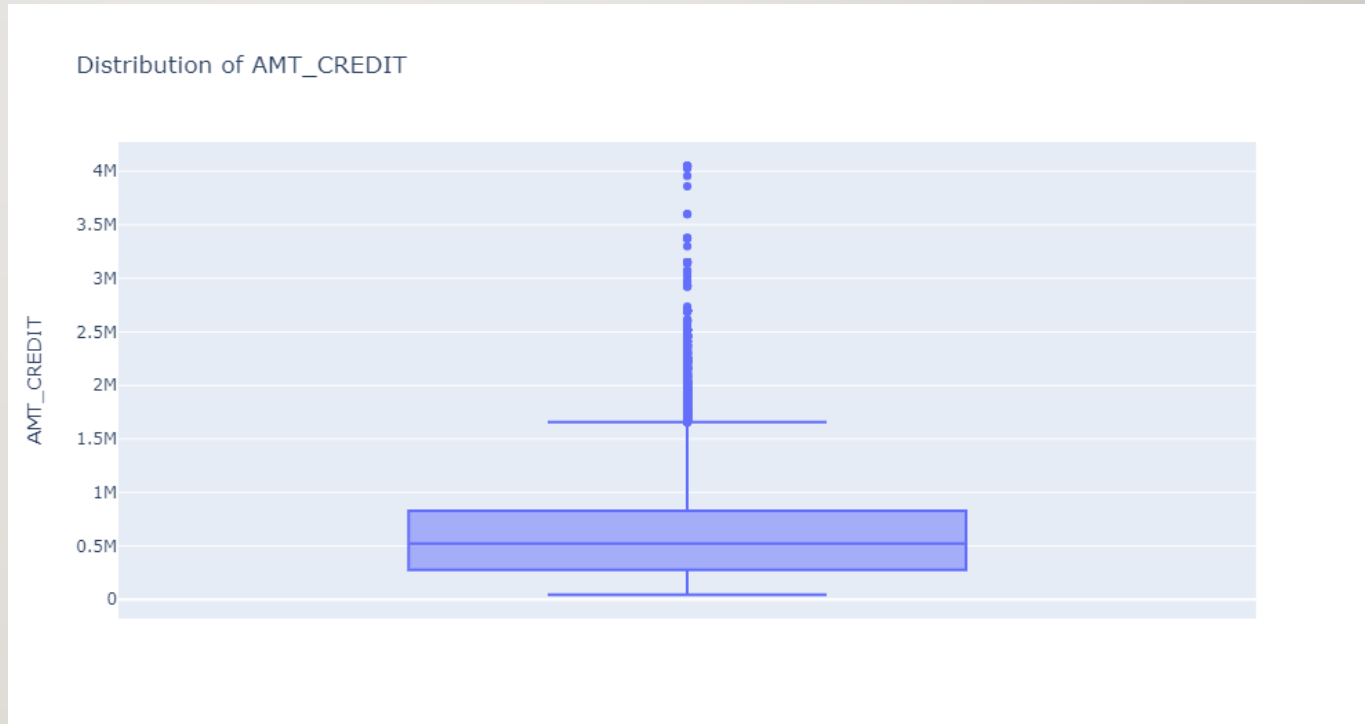


Distribution of AMT_ANNUITY

# DISTRIBUTION OF AMT_CREDIT

Observation:

Here, in the column 'AMT_CREDIT' which tells the Credit Amount of the client. We observe a value which is greater that 4.05M which is surely an outlier.
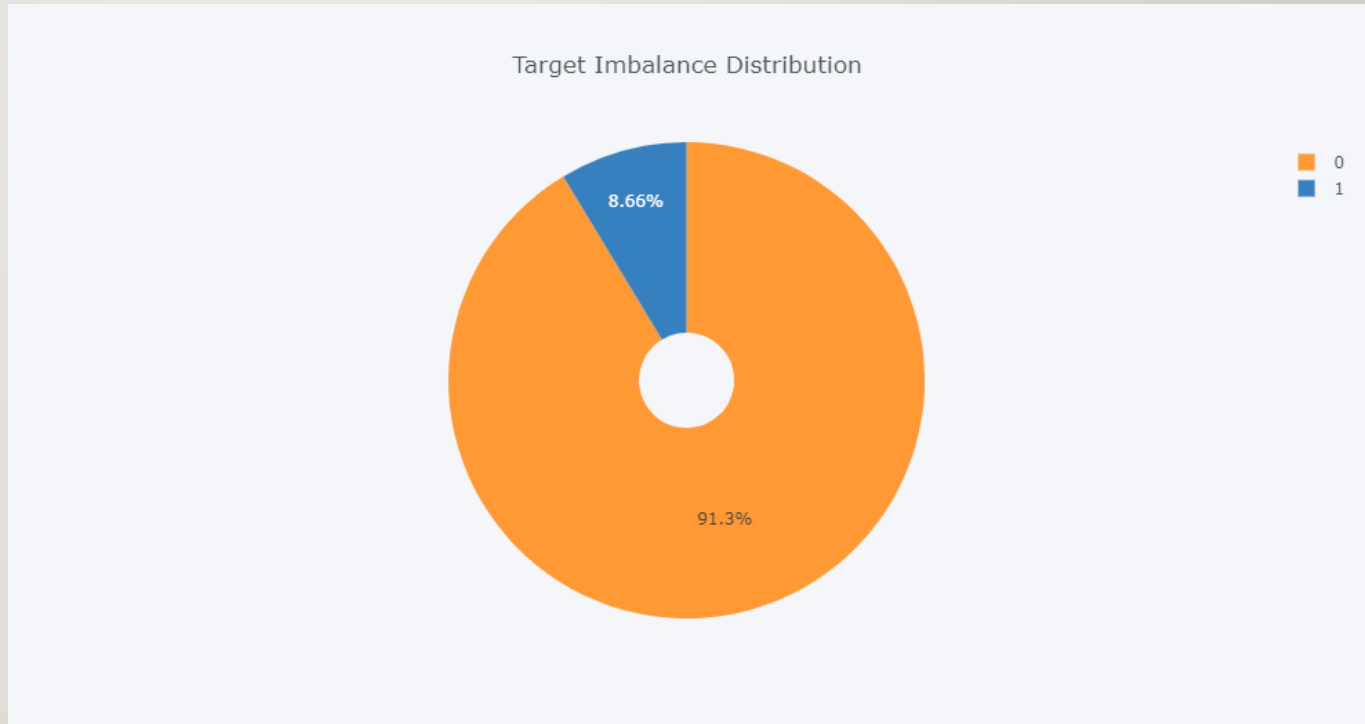


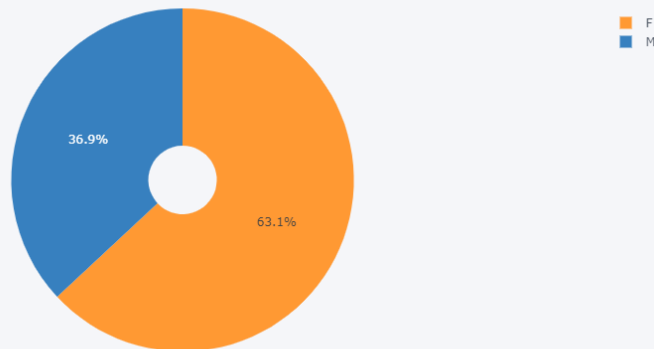Distribution of AMT_CREDIT

# TARGET IMBALANCE DISTRIBUTION

Observation:

Figure shows the Distribution of Target Imbalance and the Imbalance ratio is 10.55

# GENDER DISTRIBUTION
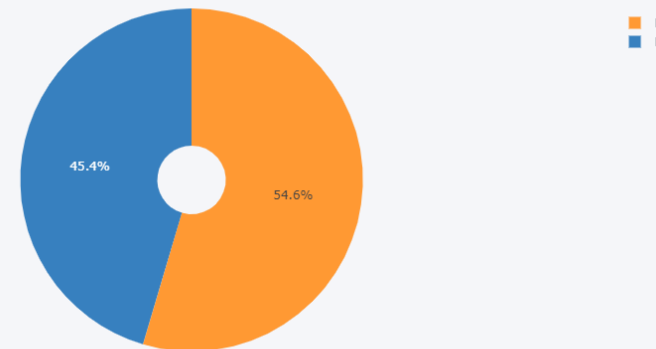


Gender Distibution of Loan Payment With NO Difficulties

F
M

36.9%

63.1%

Gender Distribution of Loan Payment Difficulties

F
M

45.4%

54.6%

Observation :

Comparing the Payment Difficulties and Non Payment Difficulties on the basis of Gender, we observe that Females are the majority in both the cases although there is an increase in the percentage in Male Payment Difficulties from Non-Payment Difficulties.
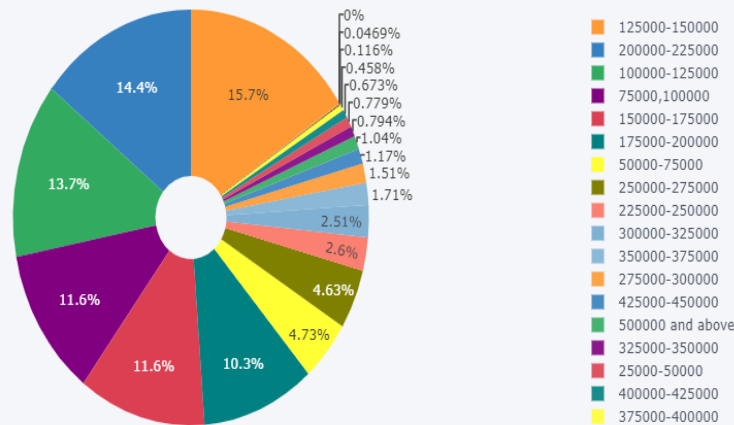
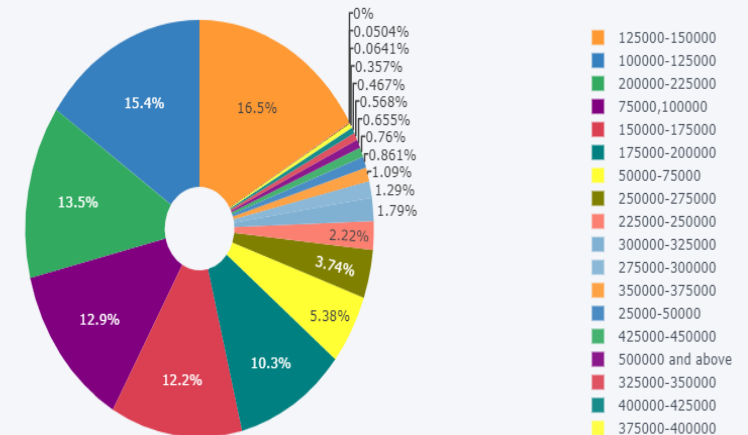# INCOME SOURCES



Obseravtion :
We observe a decrease in the percentage of Payment Difficulties who are pentioners and an increase in the percentage of Payment Difficulties who are working when compared the percentages of both Payment Difficulties and non-Payment Difficulties.

# INCOME RANGE



Observation :
We observe an increase in the percentage of Loan Payment Difficulties whose income is low when compared with the percentages of Payment Difficulties and Loan-Non Payment Difficulties

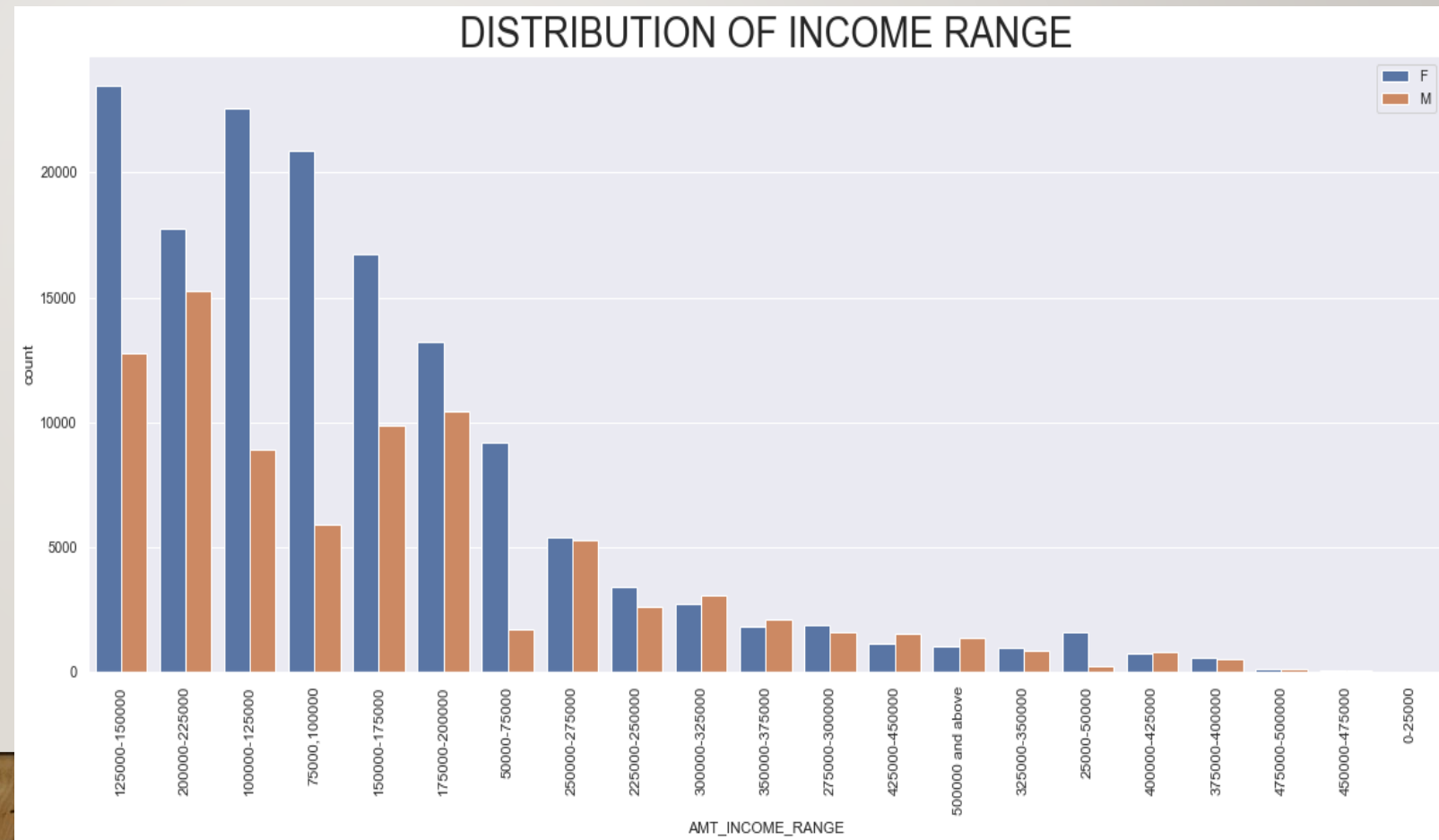# Categorical Analysis of Numerical Variables on the basis of 'Target' Variable

# TARGET 0

# DISTRIBUTION OF INCOME RANGE

Observation:

1. Female counts are higher than male.
2. Income range from 100000 to 200000 is having more number of the credits.
3. This graph show that females are more than male in having credits for that range.
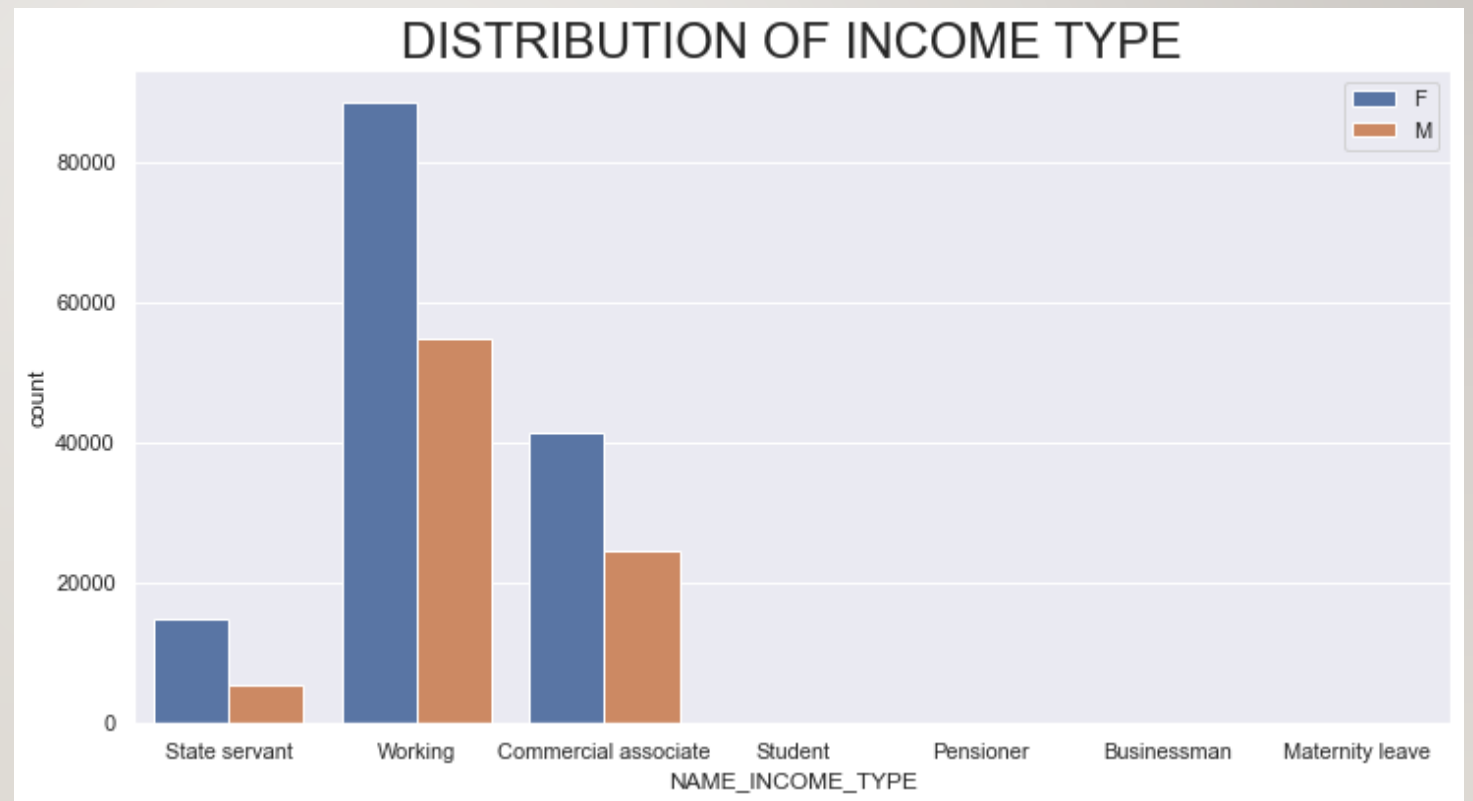4. Very less count for income range 400000 and above.

# DISTRIBUTION OF INCOME TYPE

Observation:.

1. 'Working', 'Commercial Associate', and 'State Servant' has the high the number of credits than others.
2. Females are having more number of credits than male.
3. Less number of credits for income type 'Student' ,'Pensioner', 'Businessman' and 'Maternity leave'.
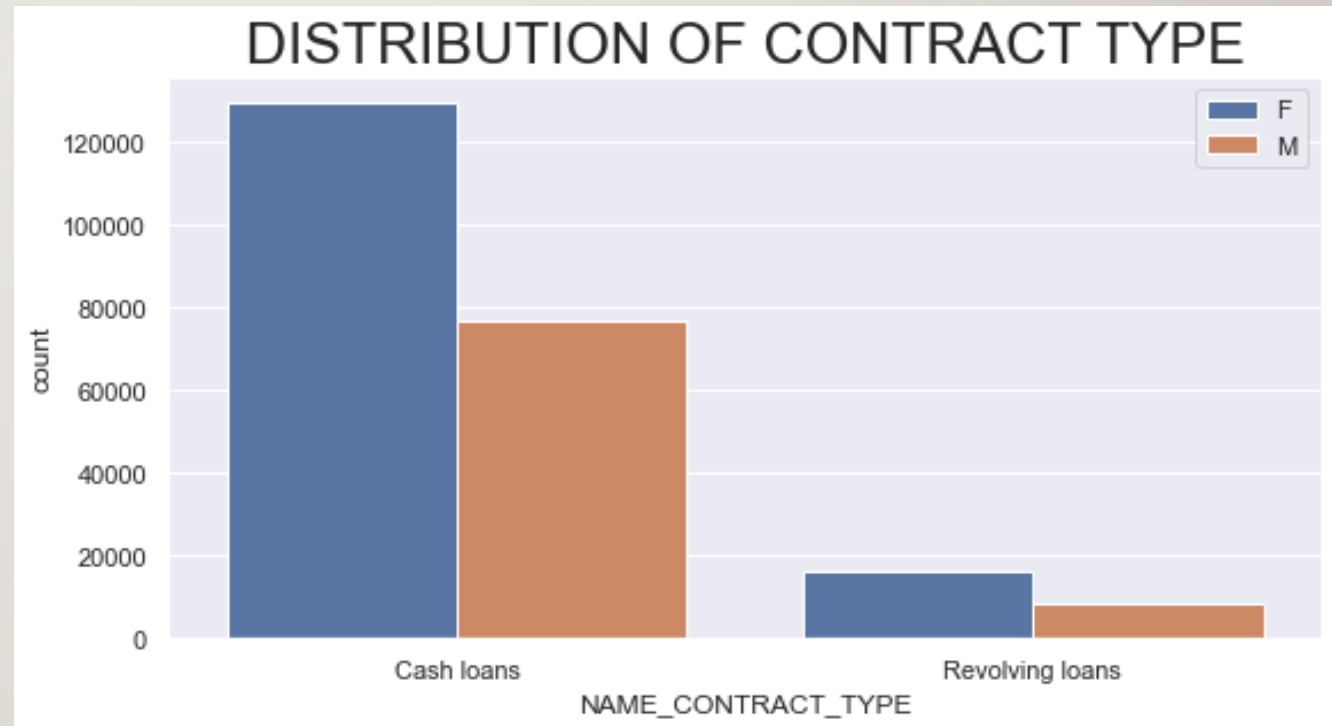
# DISTRIBUTION OF CONTRACT TYPE

Observation:

1.'Cash loans' is having higher number of credits than 'Revolving loans' contract type.
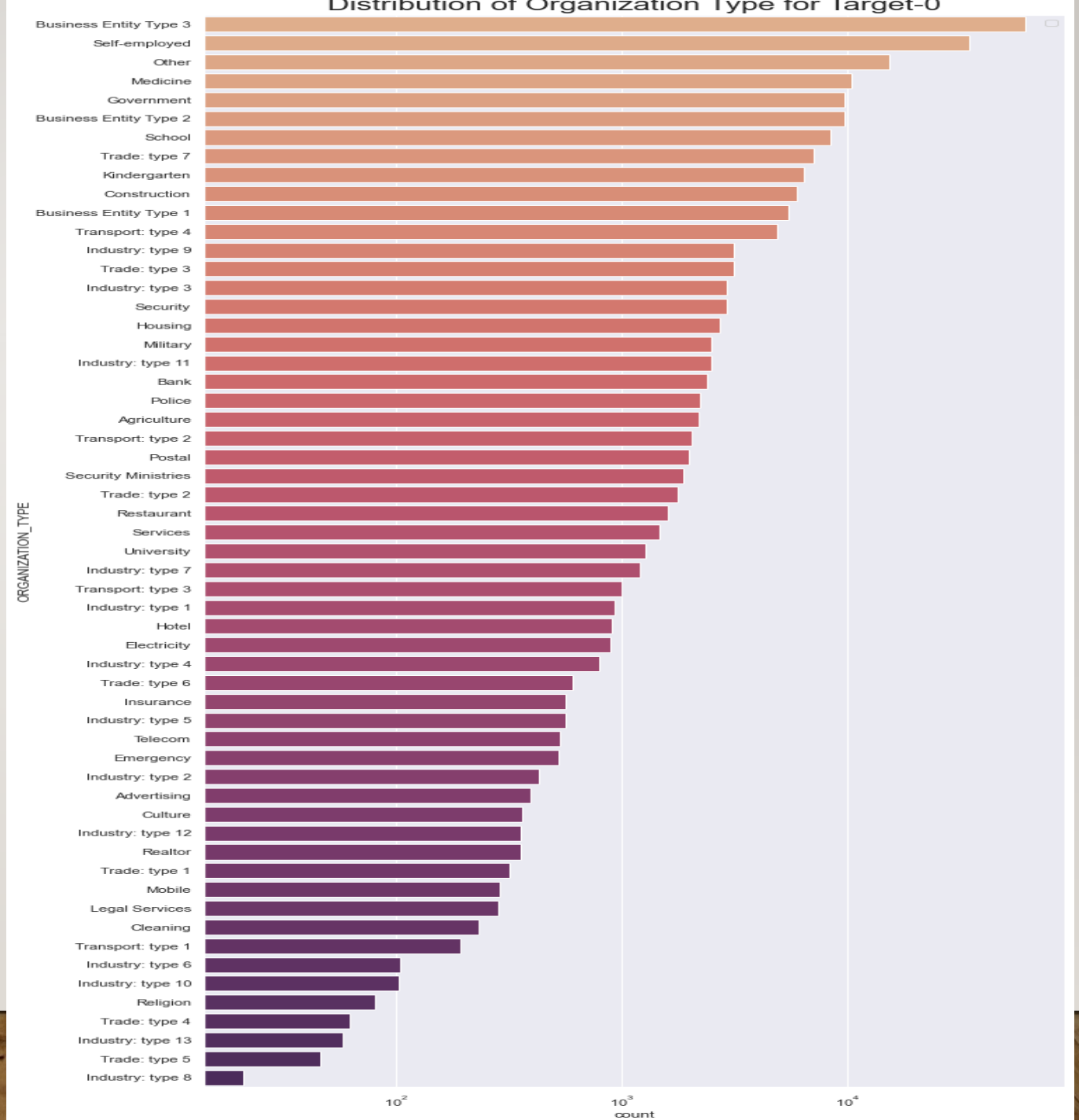
2.Females are leading in this.

# DISTRIBUTION OF ORGANIZATION TYPE

Observation:

1.Clients which have applied for credits are from most of the organization type 'Business entity Type 3' , 'Self employed', 'Other' , 'Medicine' , 'Government' and 'Business entity Type 2'.

2.Very Less Clients from 'Industry Type 8'



Distribution of Organization Type for Target-0

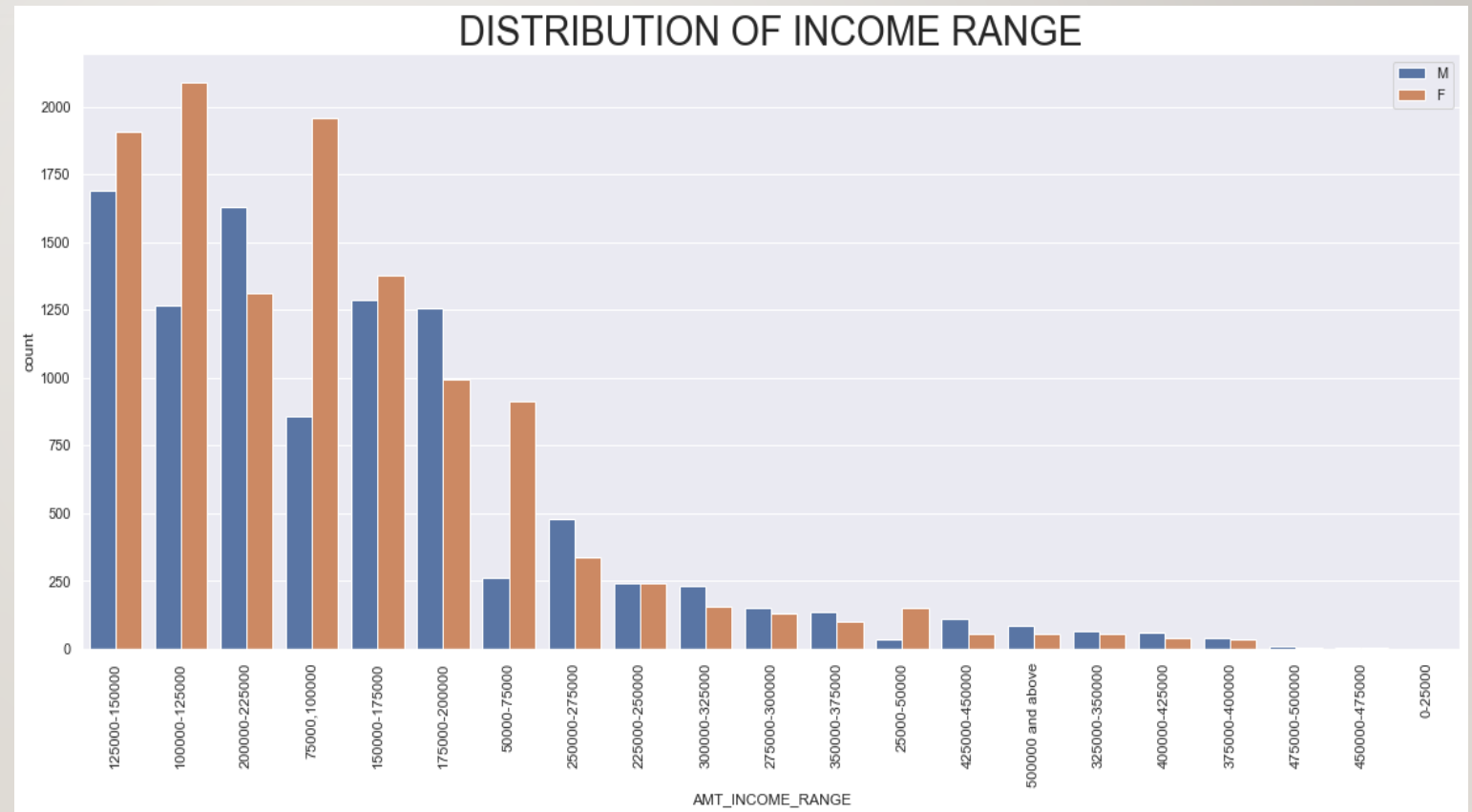# Categorical Analysis of Numerical Variables on the basis of 'Target' Variable

TARGET 1

# DISTRIBUTION OF INCOME RANGE

Observation:

1. Female counts are higher than male.
2. Income range from 100000 to 200000 is having more number of credits.
3. This graph show that females are more than male in having credits for that range.
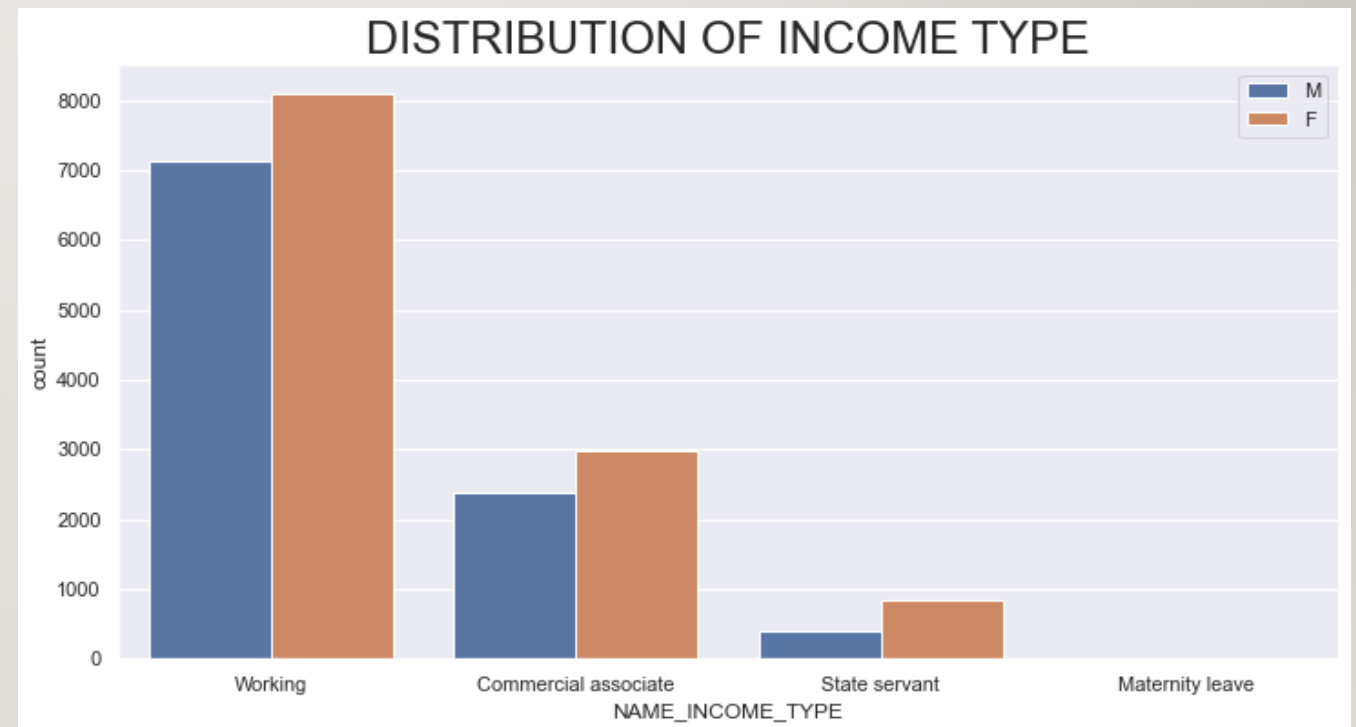4. Very less count for income range 400000 and less.

# DISTRIBUTION OF INCOME TYPE

Observation:

1. 'Working', 'Commercial associate', and 'State Servant' has higher number of credits than other i.e. 'Maternity leave.
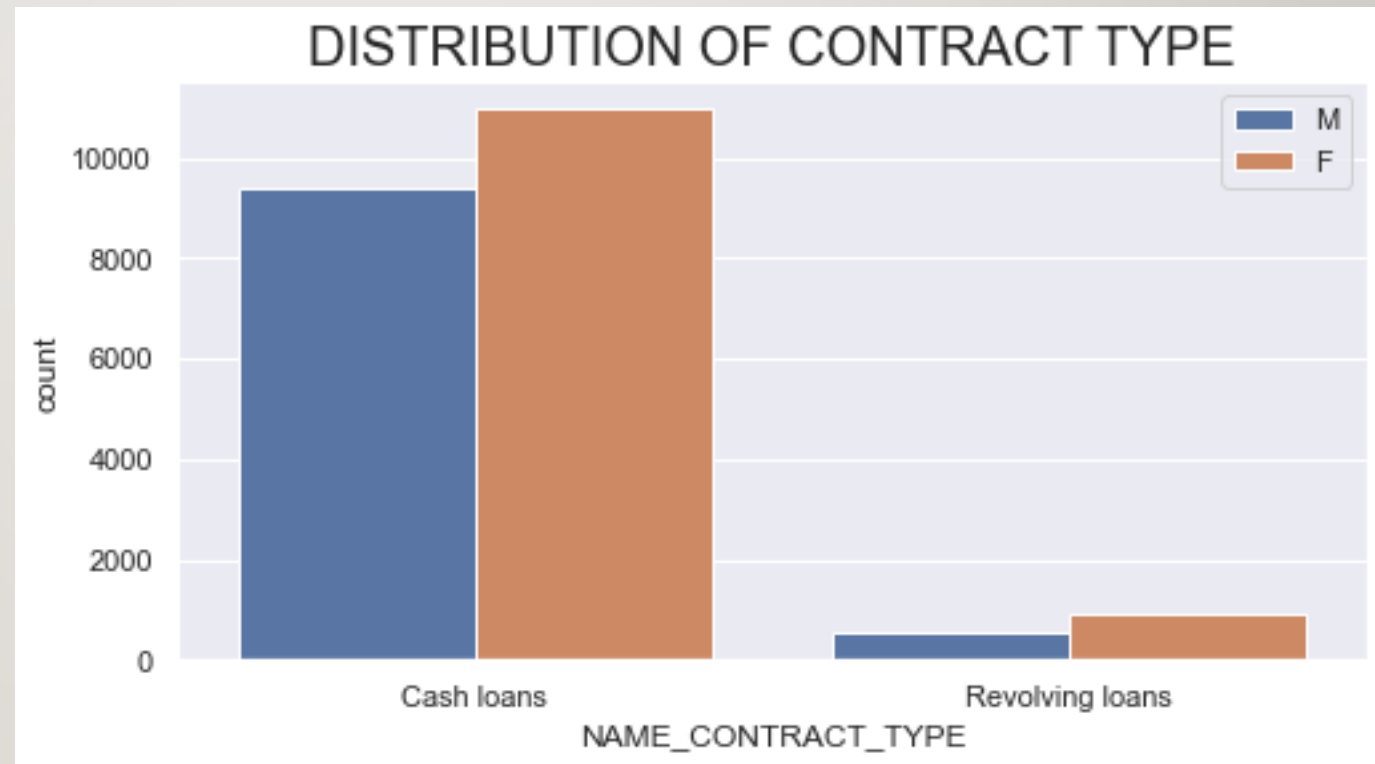
2. Females are having more number of credits than male.



DISTRIBUTION OF INCOME TYPE

# DISTRIBUTION OF CONTRACT TYPE

Observation:

1.'Cash loans' is having higher number of credits than 'Revolving loans' contract type.
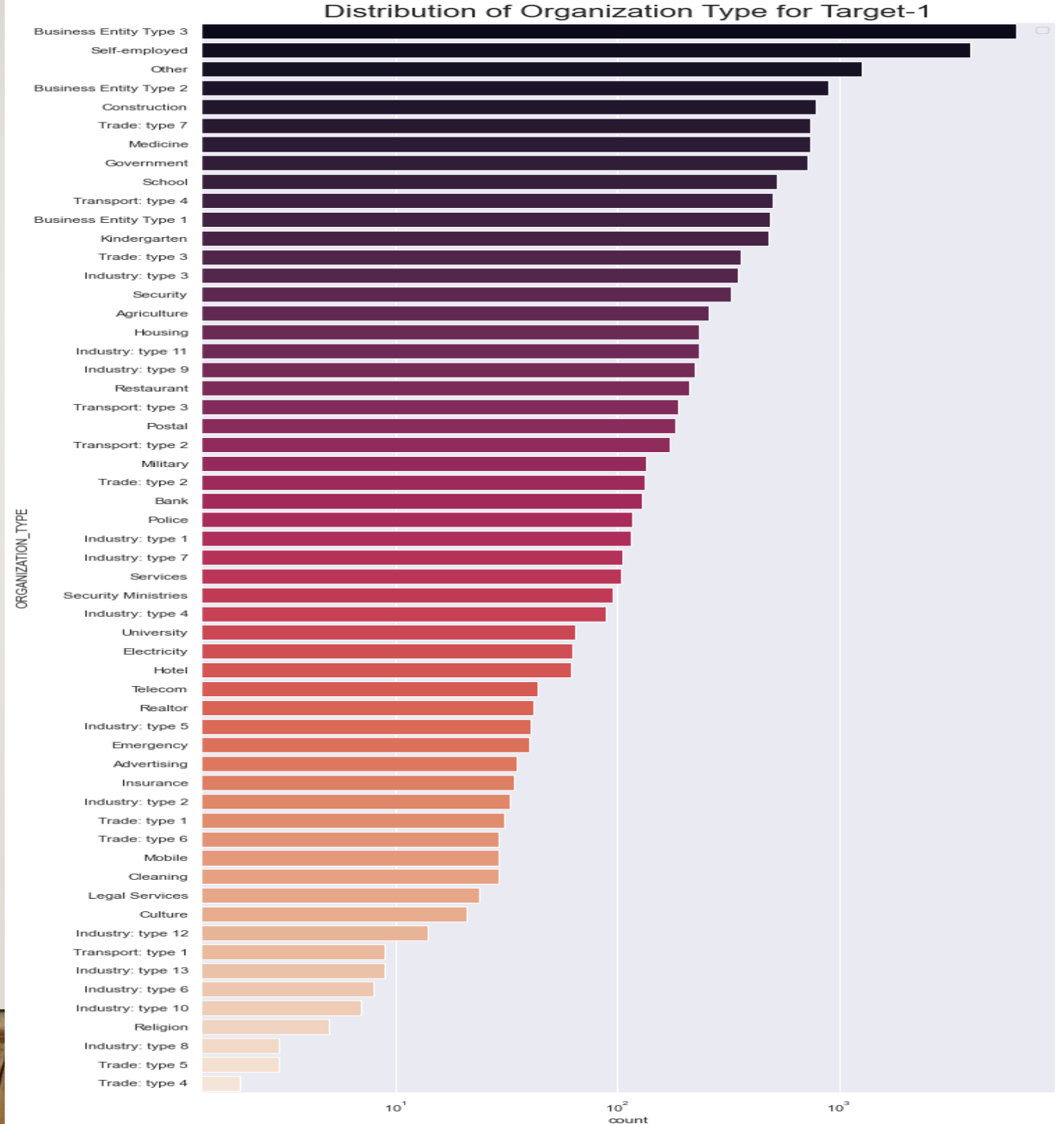
2.Females are leading for applying credits.

# DISTRIBUTION OF ORGANIZATION TYPE

Observation:

1.Clients which have applied for credits are from most of the organization type 'Business entity Type 3' , 'Self employed', 'Other' and 'Business entity Type 2'.

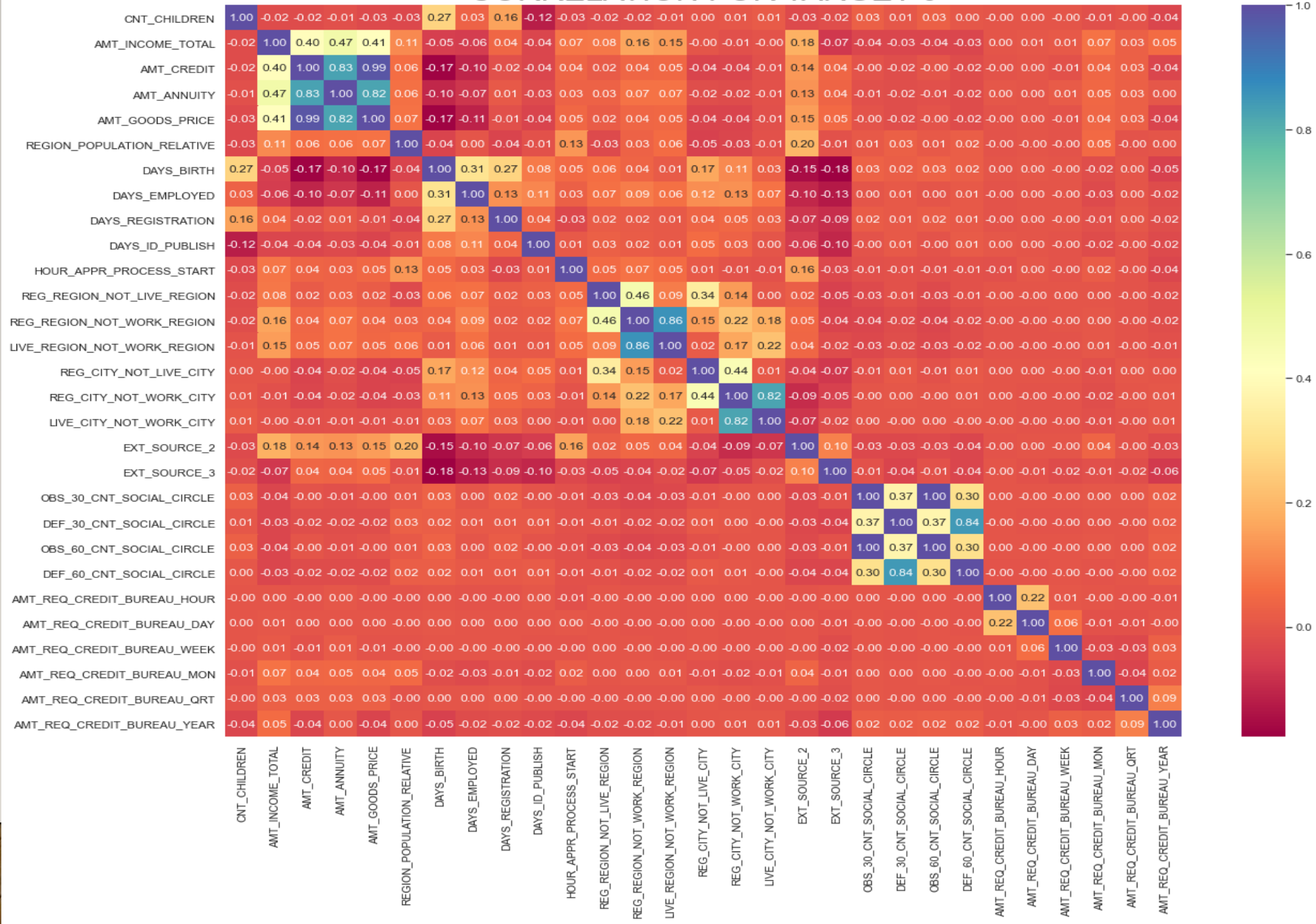2.Very Less Clients from 'Trade Type 4' , 'Trade Type 5' and 'Trade type 8'



Distribution of Organization Type for Target-1

# CORRELATION FOR TARGET 0

Observation:

1. Credit amount is inversely proportional to the date of birth, which means Credit amount is higher for low age and vice-versa.

2. 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.

3. Income amount is inversely proportional to the number of children client have, means more income for less children client have and vice-versa.

4. Less children client have in densely populated area.

5. The income is higher in densely populated area.

6. Credit amount is higher in densely populated area.

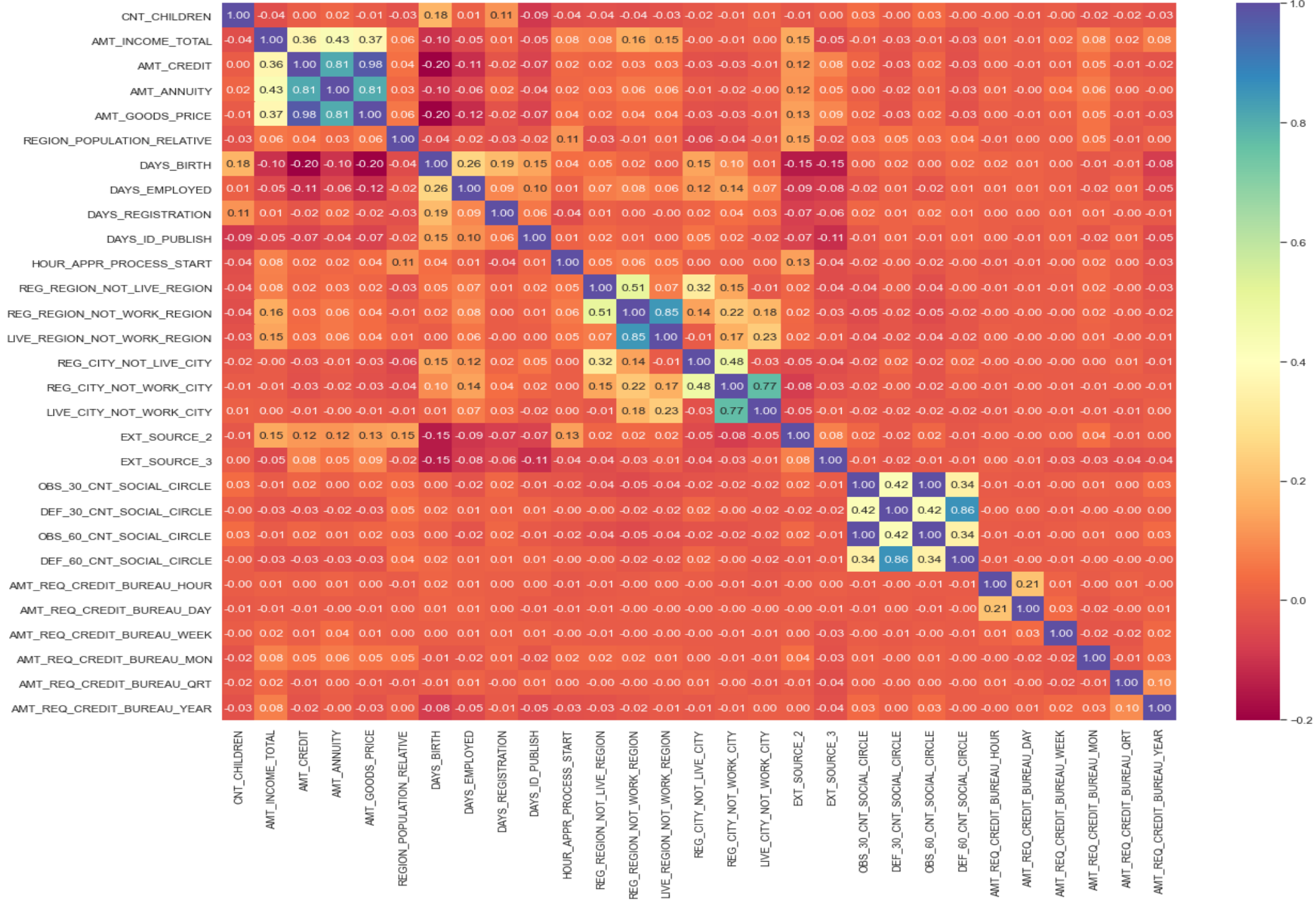CORRELATION FOR TARGET 0

# CORRELATION FOR TARGET 1

Observation:

Heat map for Target 1 is also having little bit same observation just like Target 0. But for few points are different. They are listed below.

1. Client's permanent address does not match to contact address which are having less children and vice-versa.

2. Client's permanent address does not match to work address which are having less children and vice-versa.
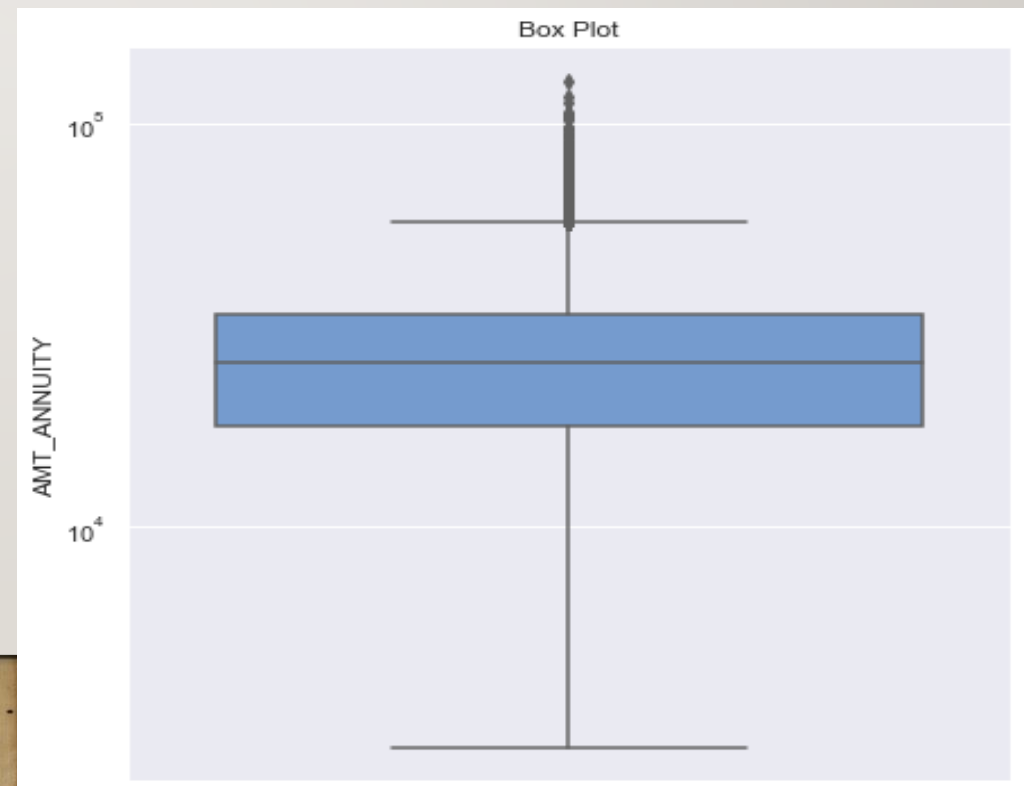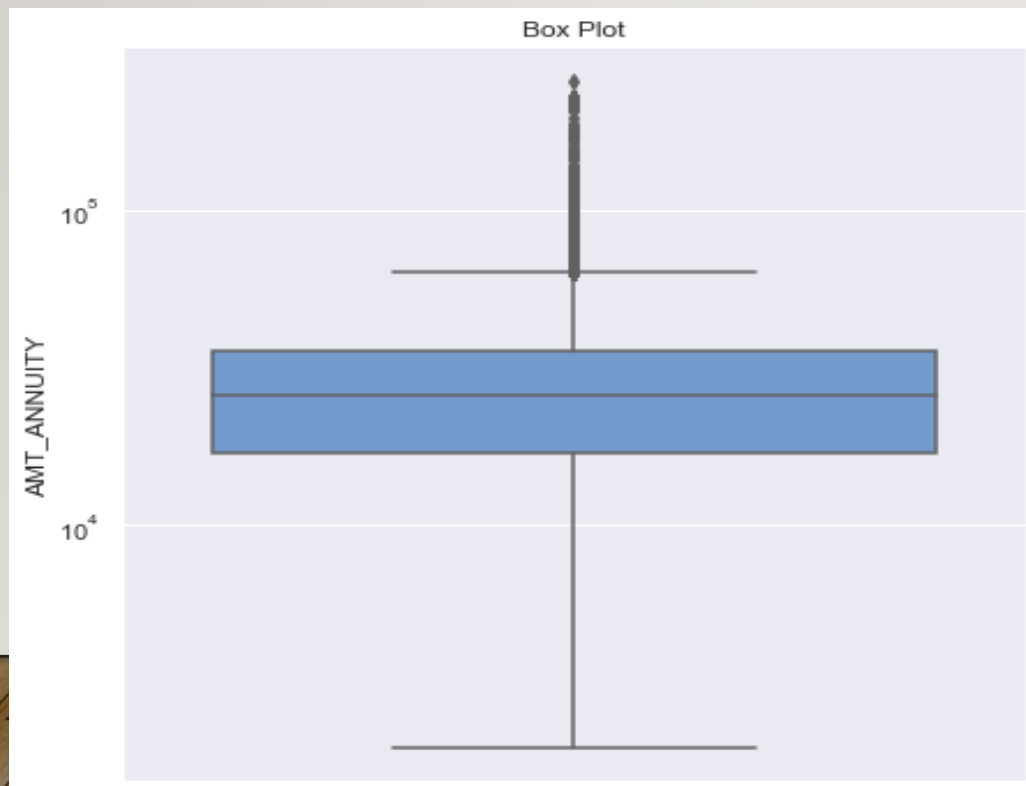
CORRELATION FOR TARGET 1

# Univariate Analysis of Numerical Variables
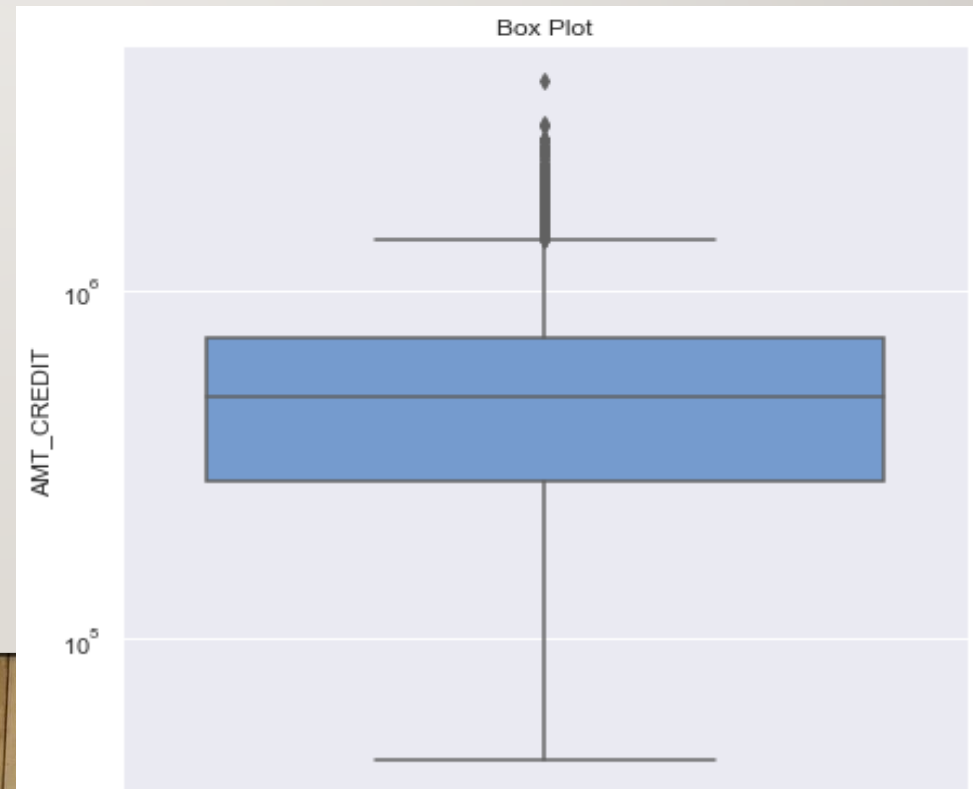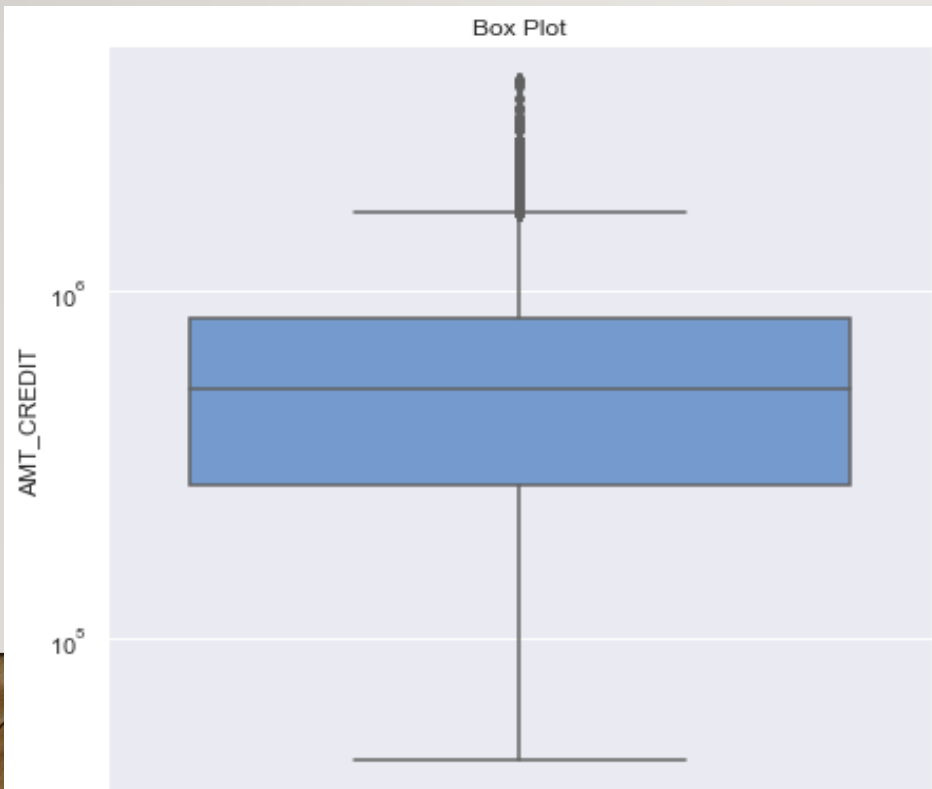
# LOAN ANNUITY

Observations:

1. Some outliers are noticed in income amount.

2. The first quartile is bigger than third quartile for annuity amount which means most of the annuity clients are from first quartile.
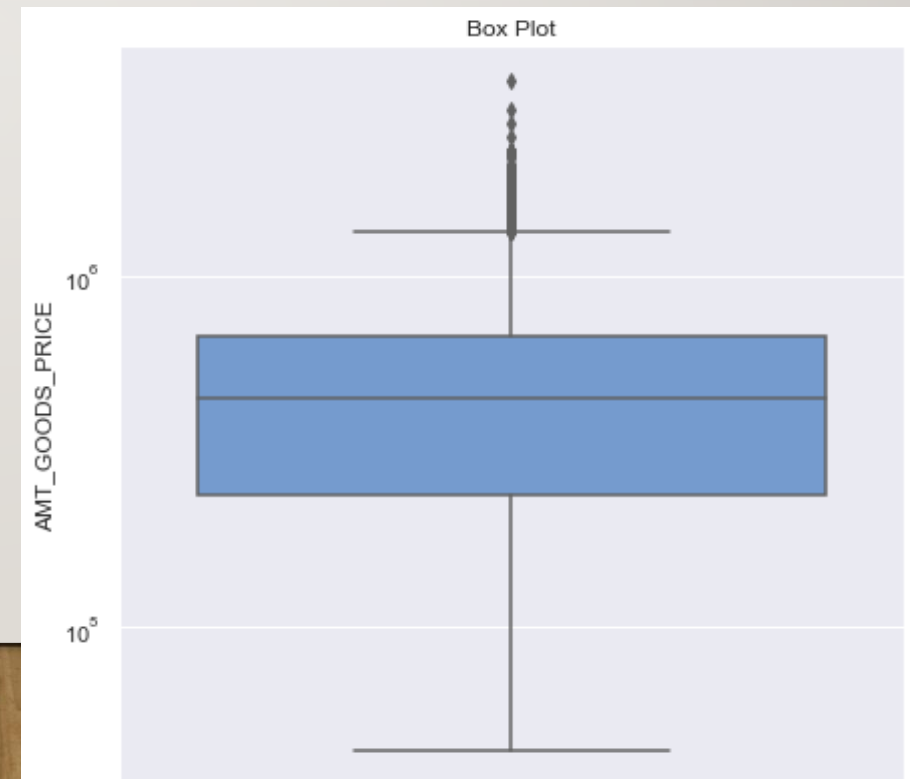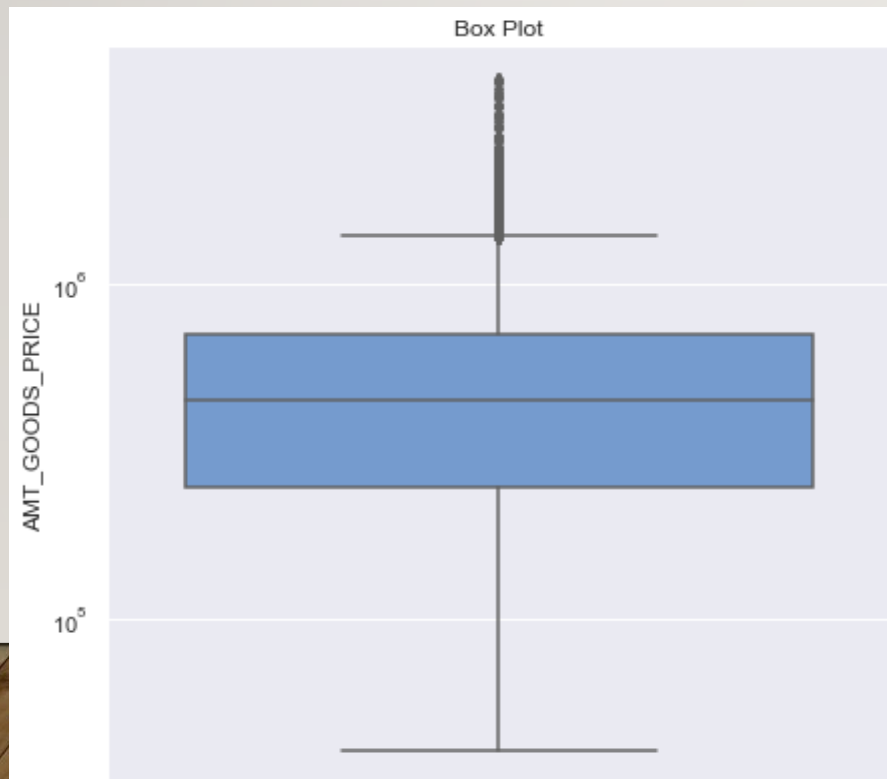
# CREDIT AMOUNT

Observations:

1. Some outliers are noticed in income amount.

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

# GOODS PRICE

Observations:

1. Some outliers are noticed in income amount.

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

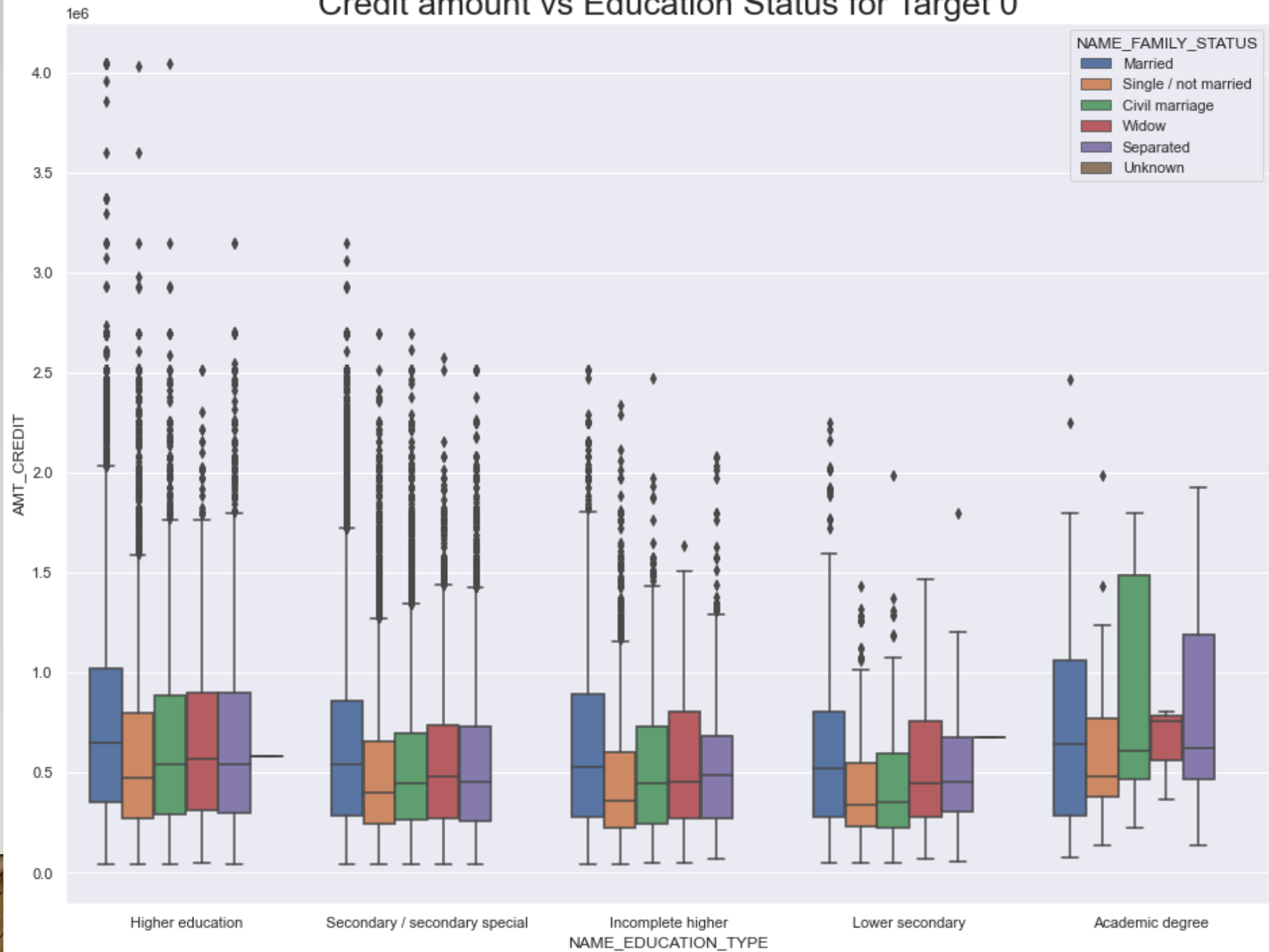# Bivariate Analysis for Numerical Variables

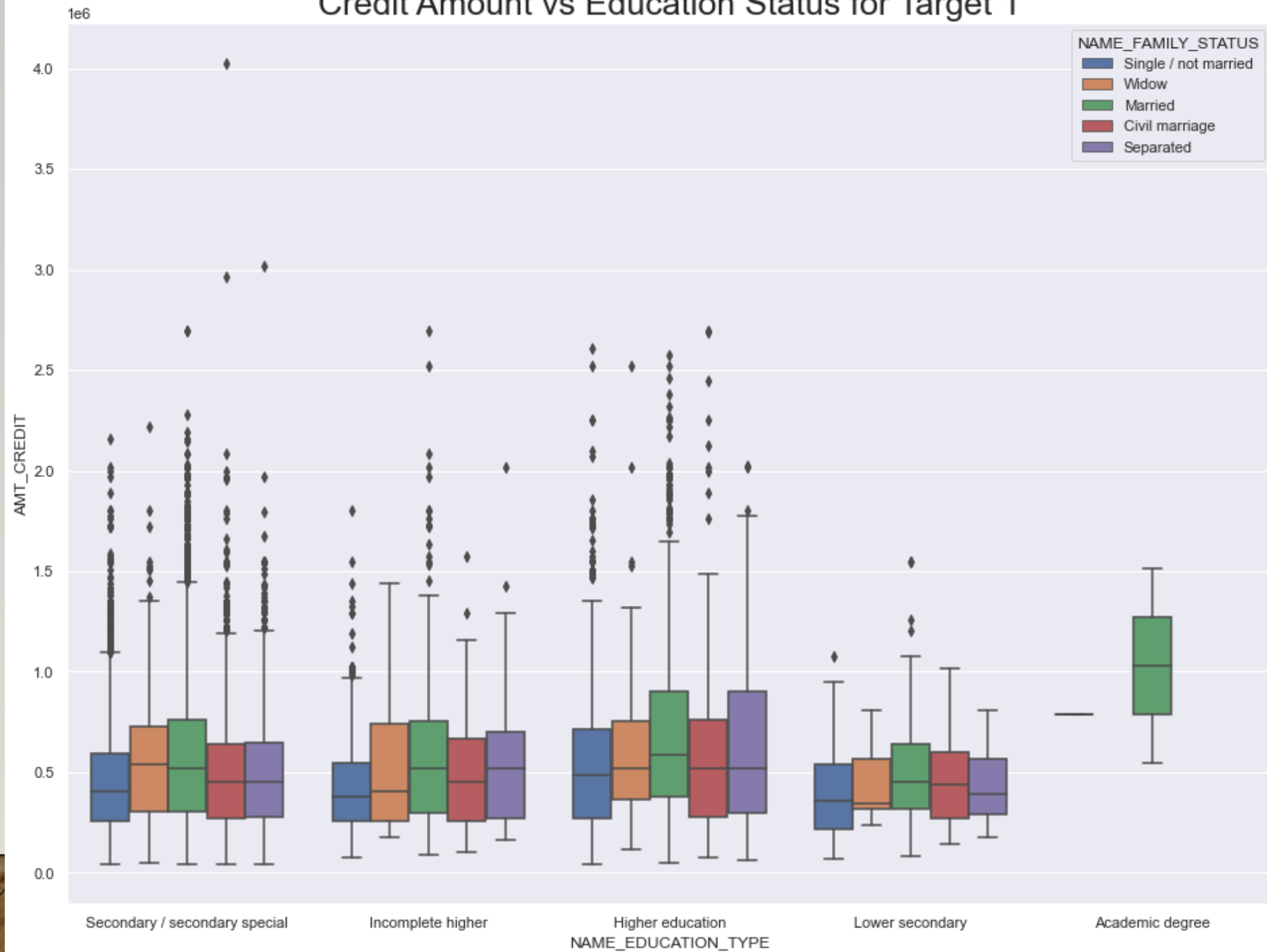# CREDIT AMOUNT VS EDUCATION STATUS

Observation :

- Here we can see that the range of customers without payment of Academic degree is higher than the customer of with payment. And the rest of the Education type is almost same for both the cases.

Credit amount vs Education Status for Target 0
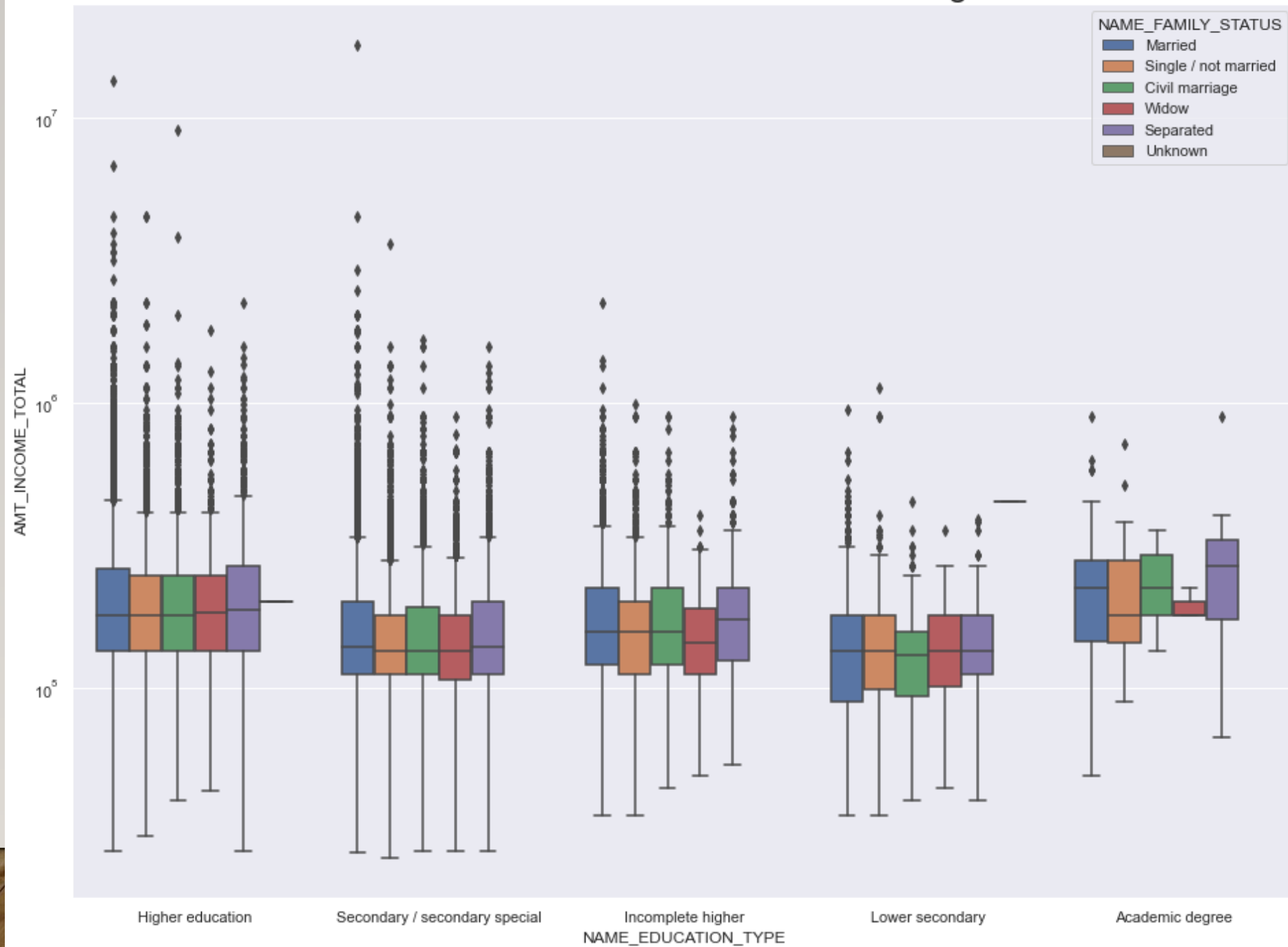
Credit Amount vs Education Status for Target 1

# INCOME AMOUNT VS EDUCATION STATUS
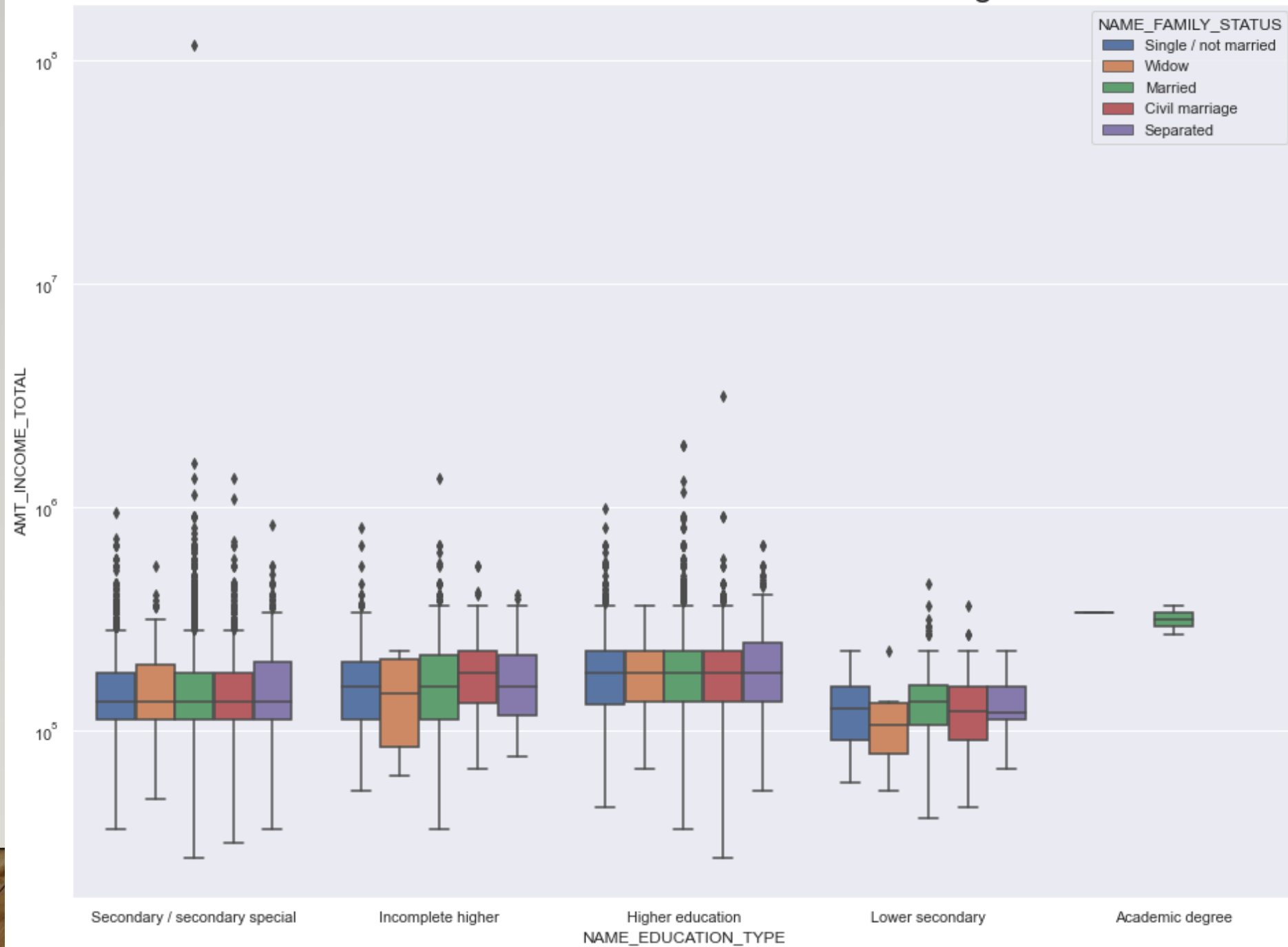
Observation:

Here we can see that the customers without payment is having more outliers as compare to the customer with payment.
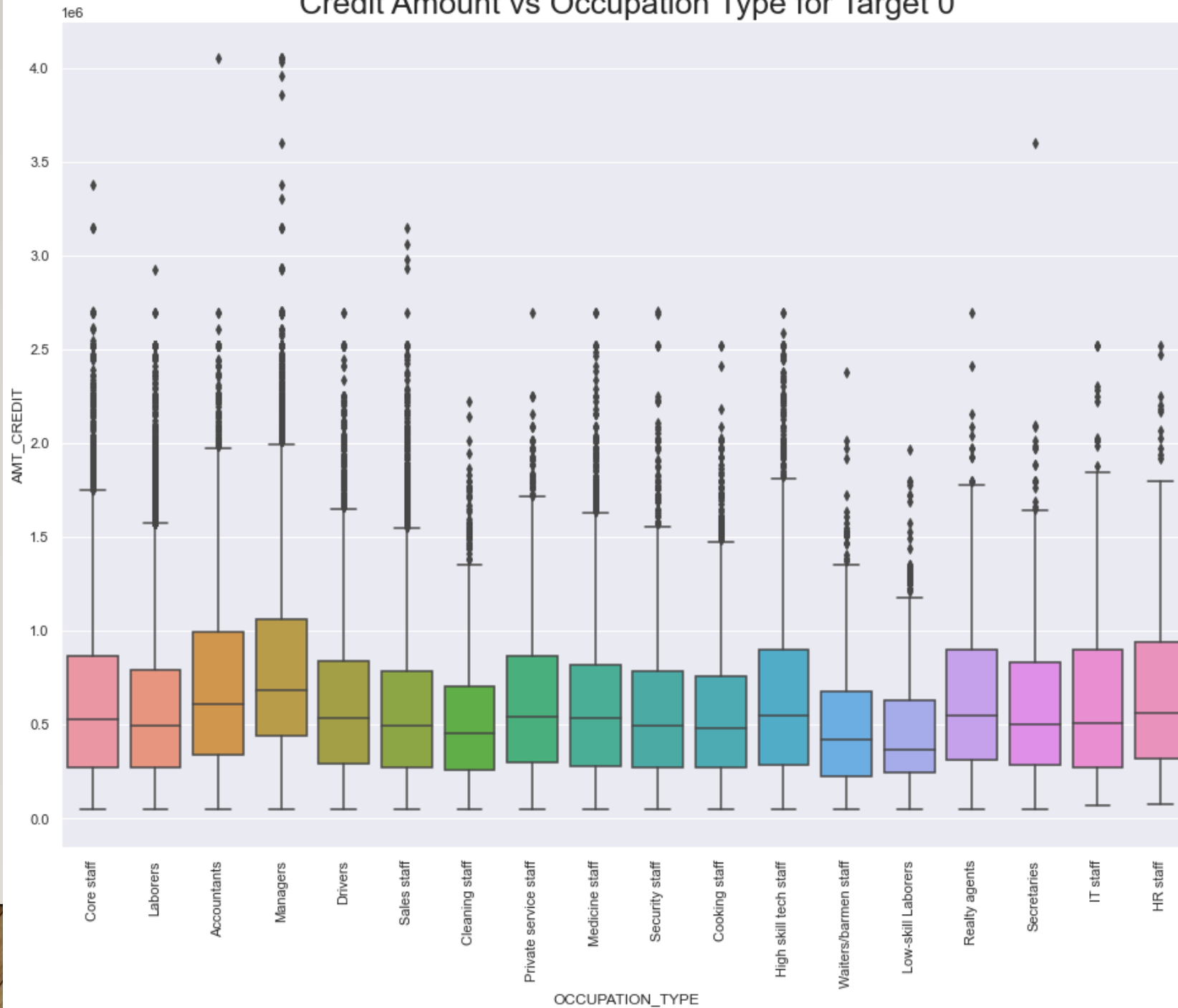
Income amount vs Education Status for Target 1
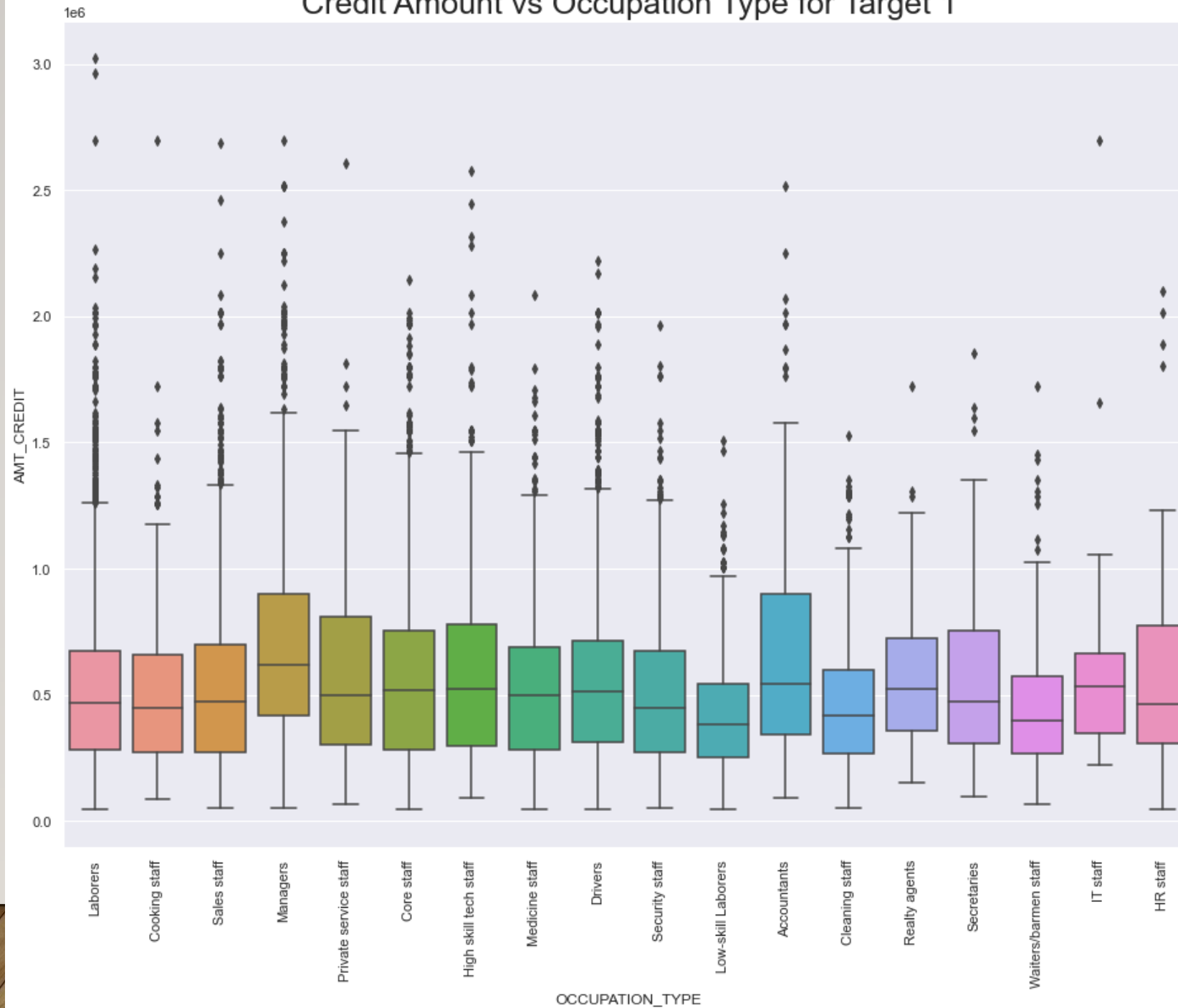
# CREDIT AMOUNT VS OCCUPATION TYPE

Observation :

Here we can see that the range of the customers without payment more as compare to the customers with payment.

Credit Amount vs Occupation Type for Target 1

# Bivariate Analysis of Numerical vs Numerical Variables
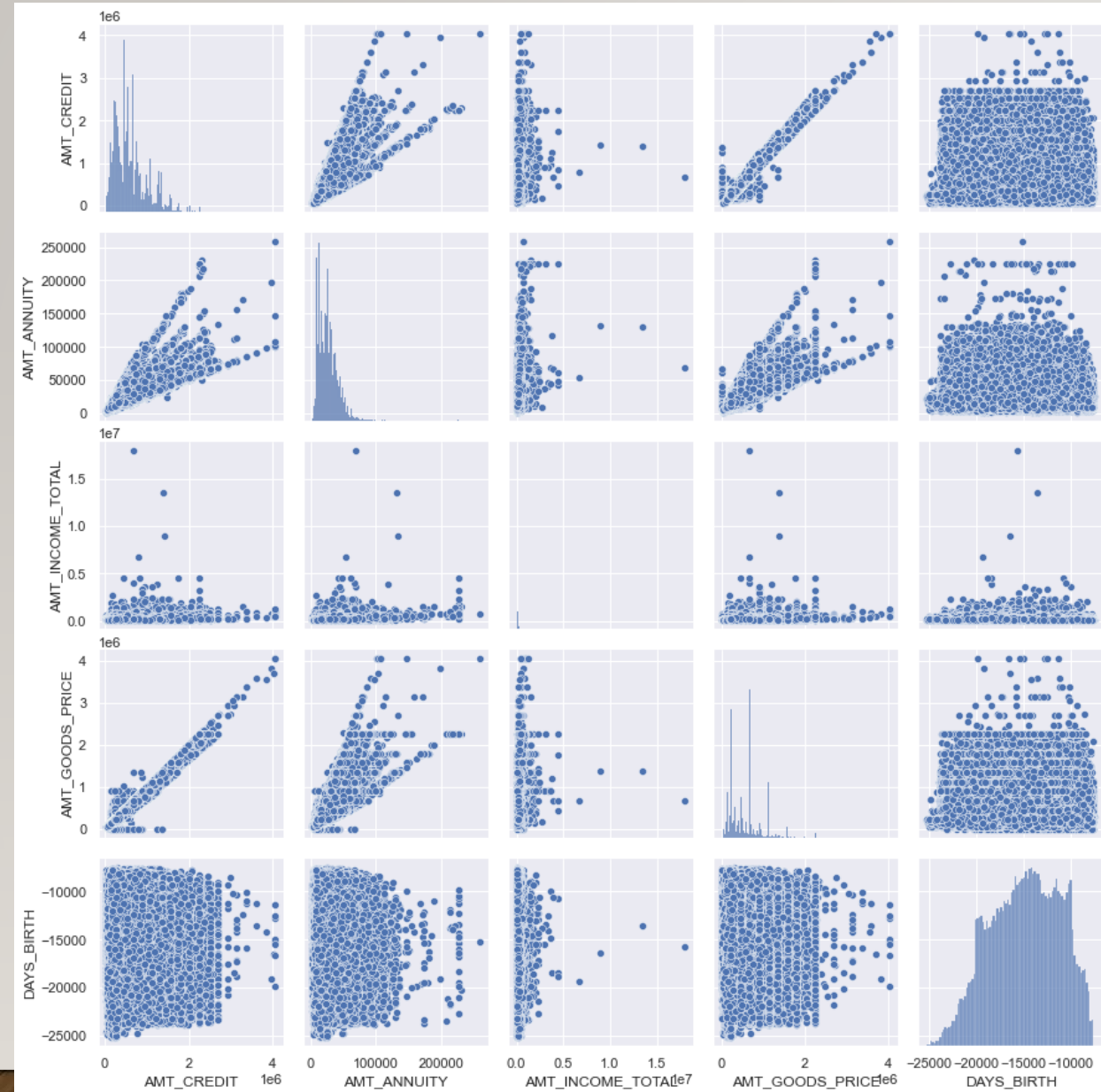
# PAIR PLOT FOR TARGET 0 AND TARGET 1
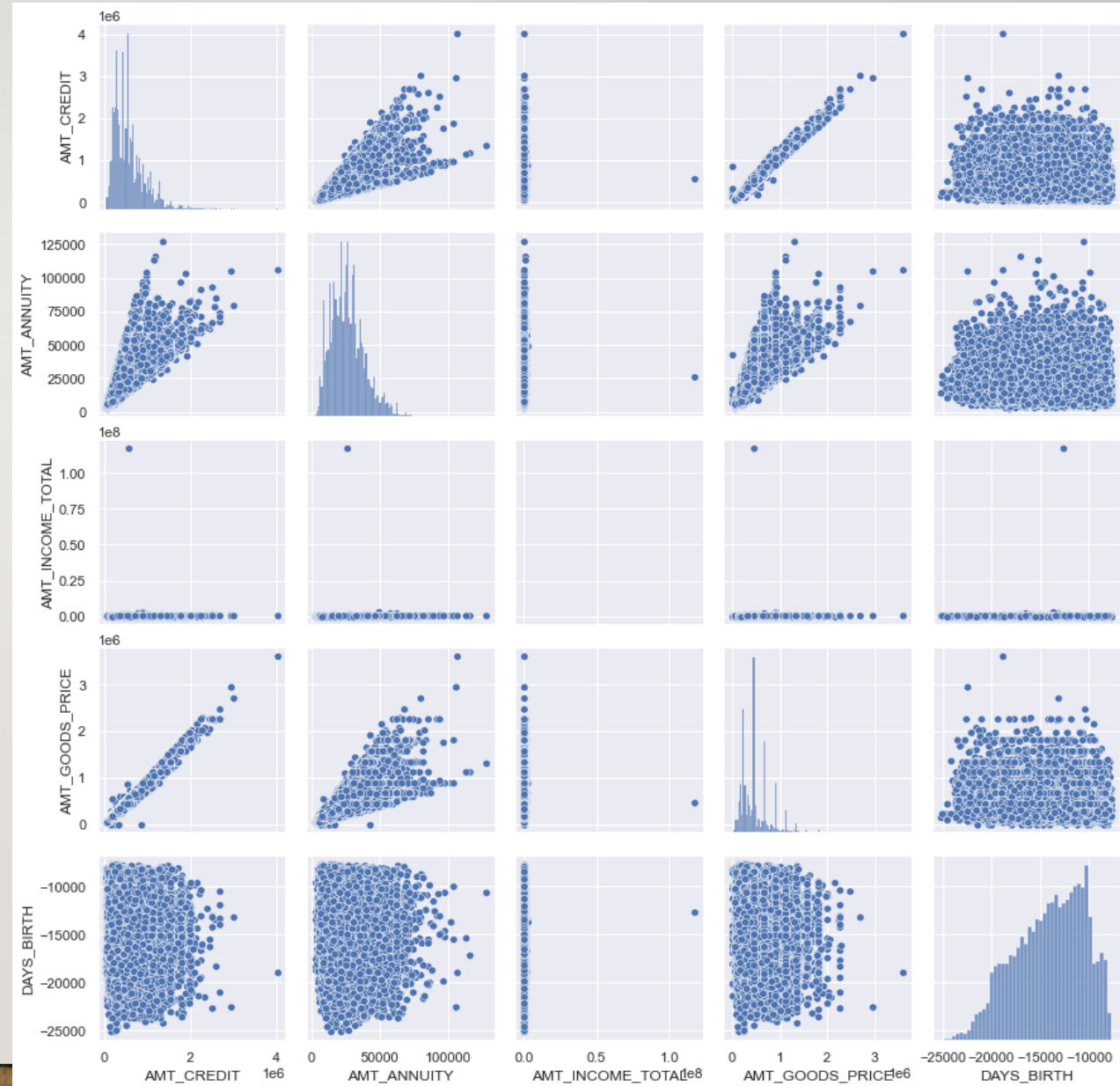
Observation:

We observe that there is a high correlation between credit amount and goods price. There appears to be some deviancies in the correlation of Loan-Payment Difficulties and Loan Payment with No Difficulties such as credit amount v/s income.
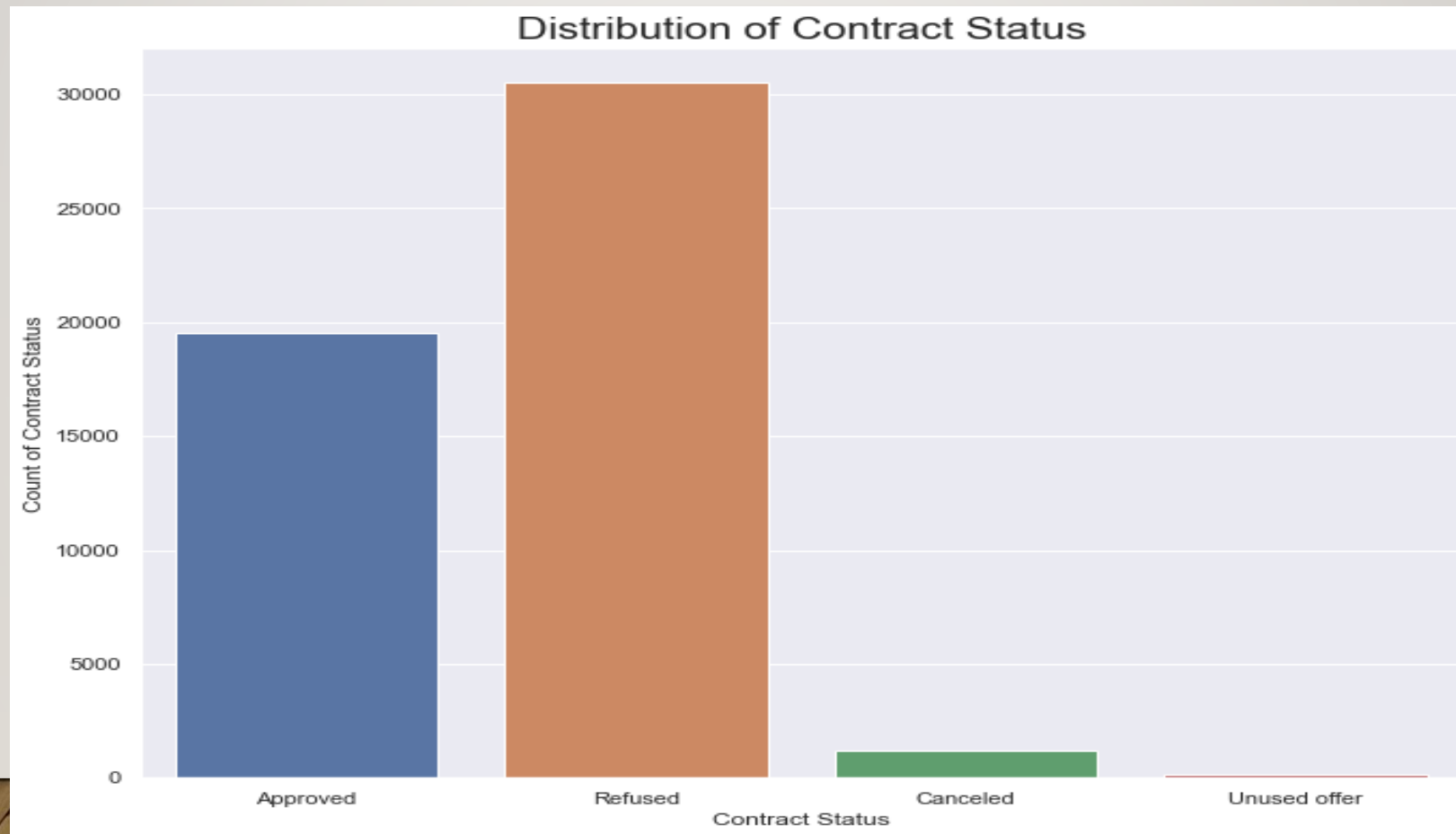
Target 0      Target 1

# Some Univariate Analysis on Previous Application Data
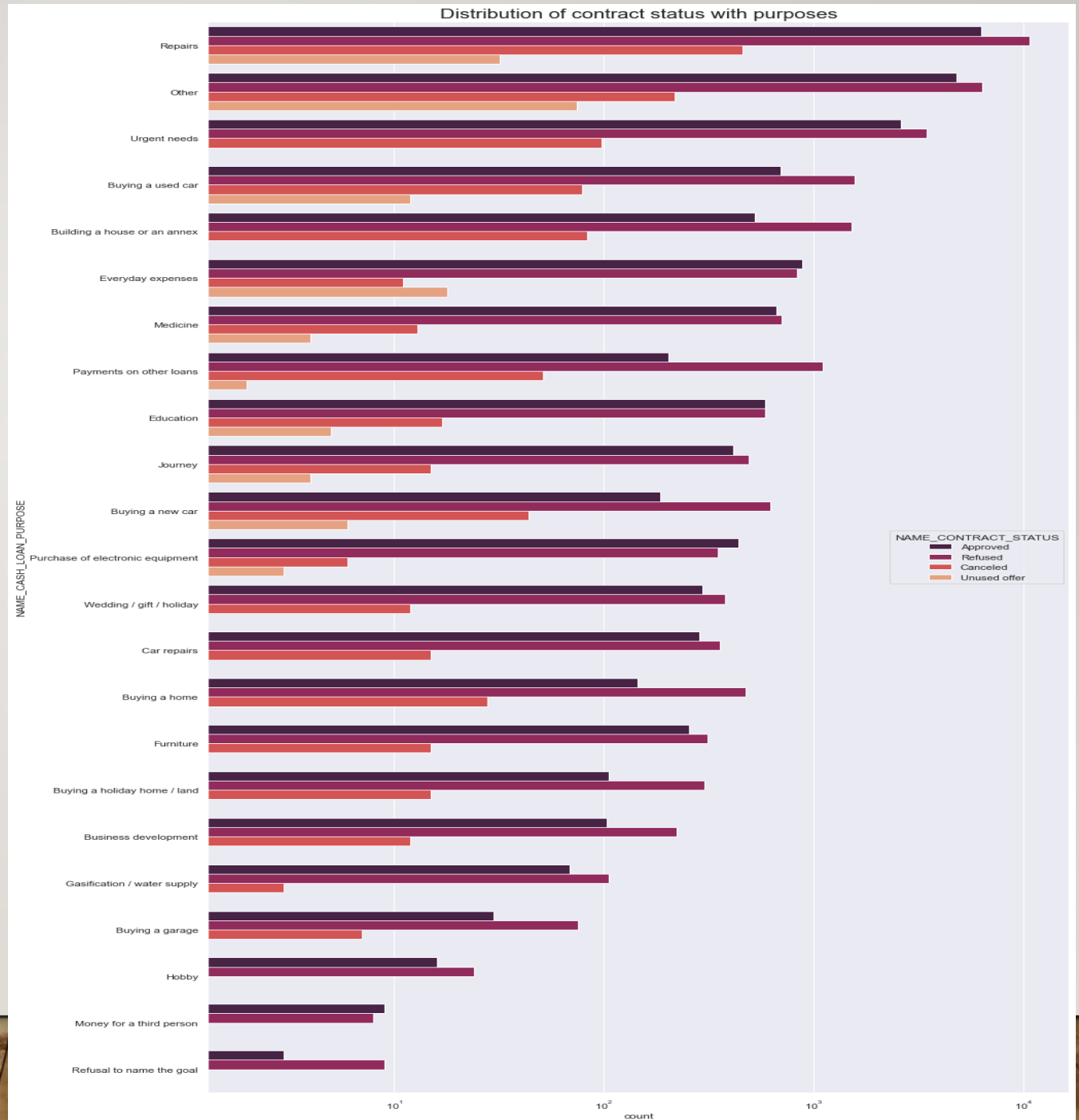
# DISTRIBUTION OF CONTRACT STATUS

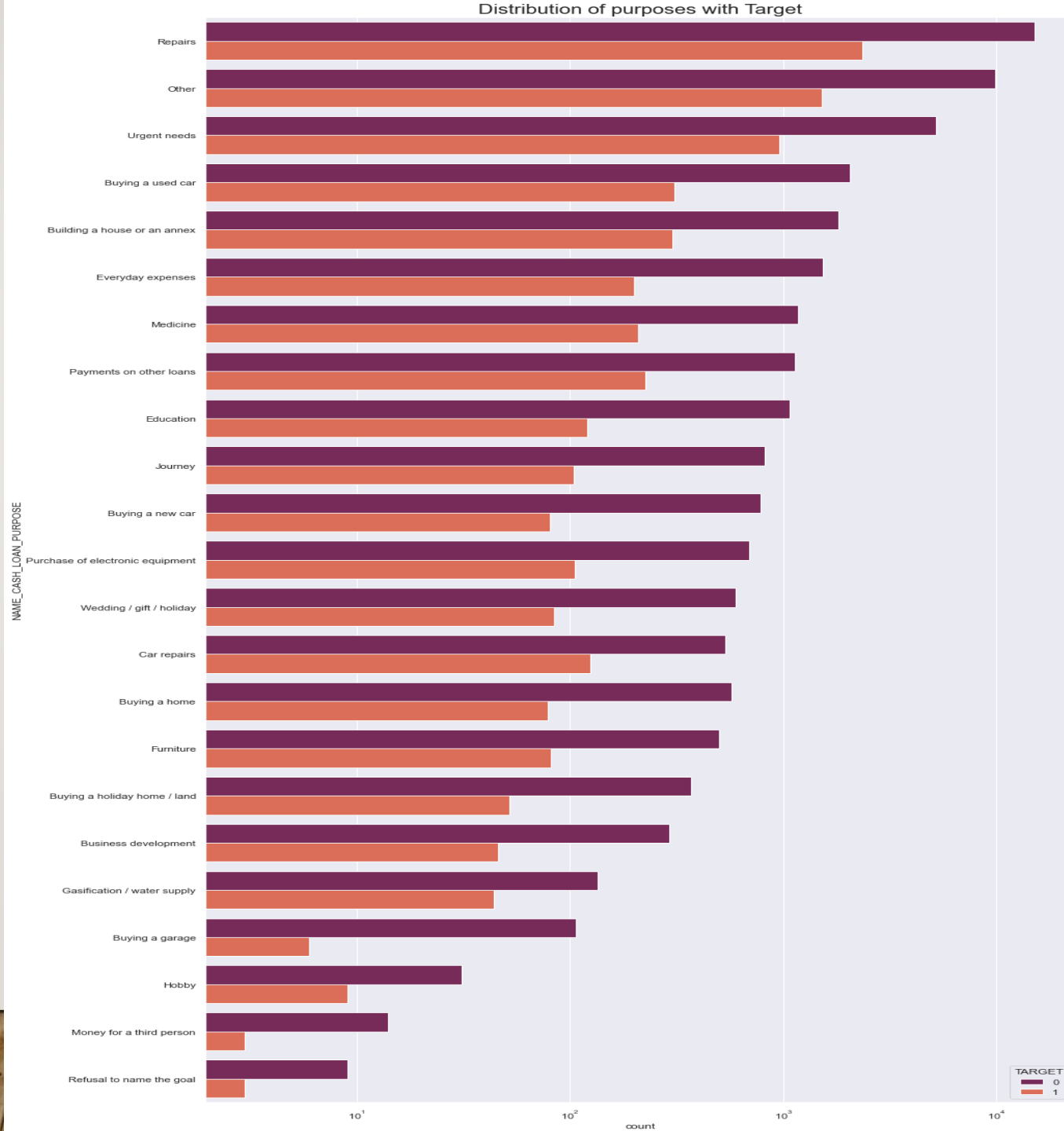# DISTRIBUTION OF CONTRACT STATUS WITH PURPOSES

Observation:

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

2. For education purposes we have equal number of approves and rejection

3. Paying other loans and buying a new car is having significant higher rejection than approves.



Distribution of contract status with purposes

# DISTRIBUTION OF PURPOSES WITH TARGET

Observation:

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

2. There are few places where loan payment is significant higher than facing difficulties. They are 'Buying a garage', 'Business development', 'Buying land', 'Buying a new car' and 'Education' Hence we can focus on these purposes for which the client is having for minimal payment difficulties.
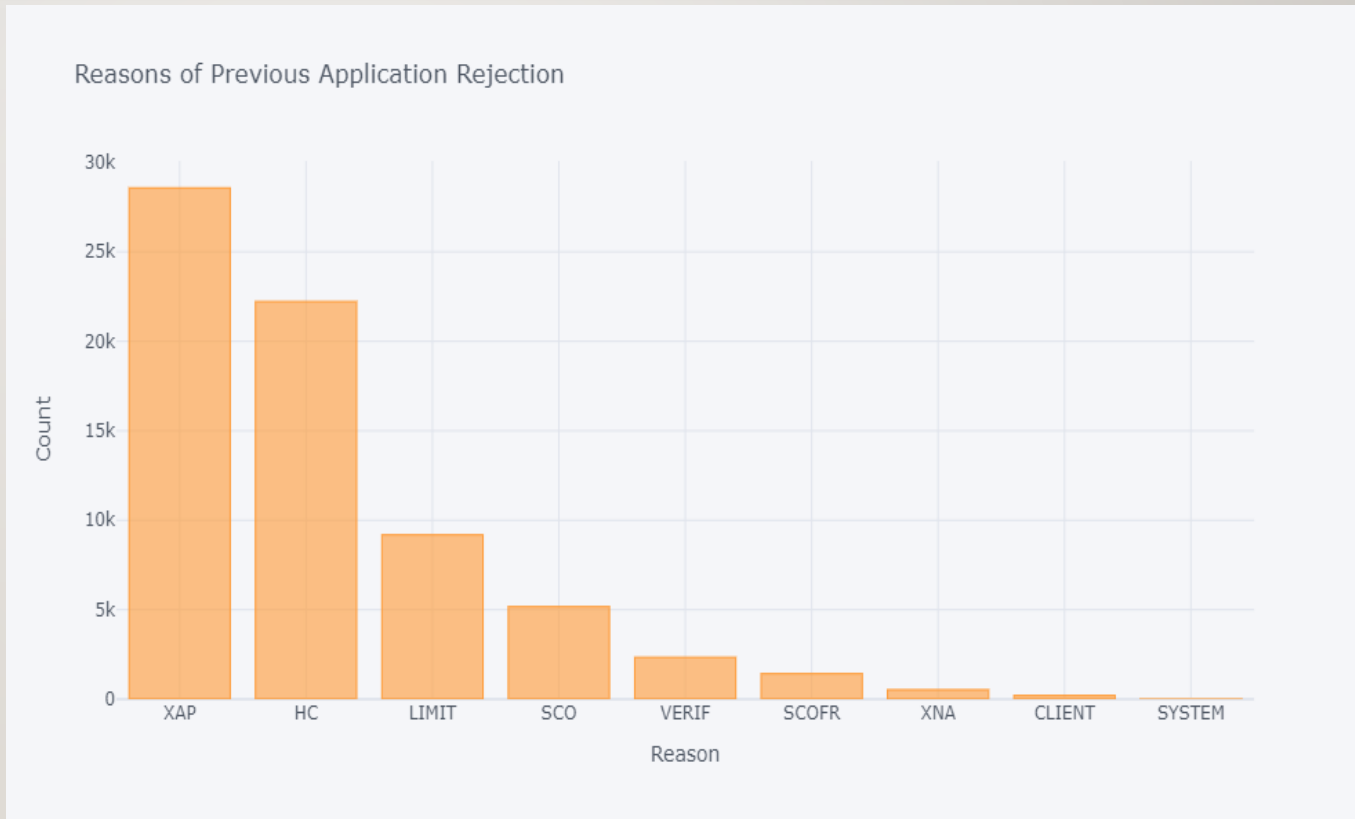


Distribution of purposes with Target

# REASONS OF PREVIOUS APPLICATION REJECTION

Observation:

We observe that XAP is the reason majority of applications got rejected.
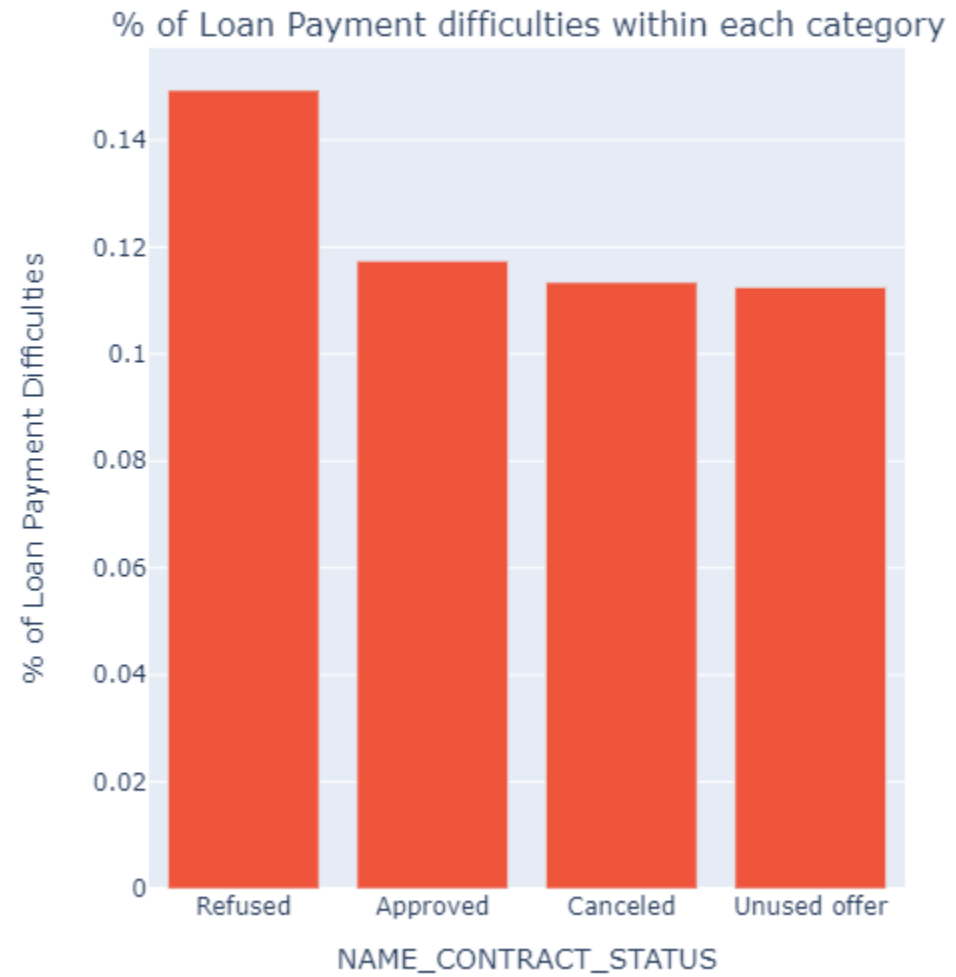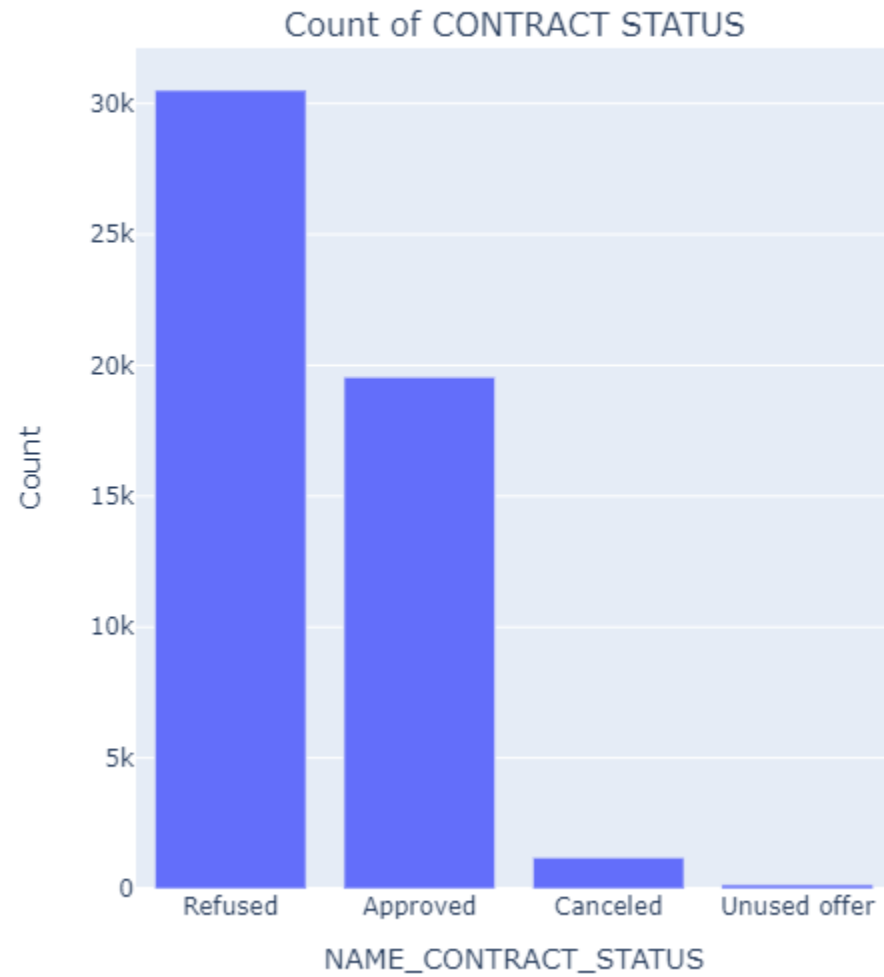
# Bivariate Analysis

# CONTRACT STATUS

Observation:

From the first graph it can be seen that most of the contracts from previous application have been Approved.

It can be clearly seen from the second graph that the -

1. 'Refused' contracts from previous application are the ones who have maximum % of Loan-Payment Difficulties from current application.

2. 'Approved' contracts from previous application are the ones who have minimum % of Loan-Payment Difficulties from current application.
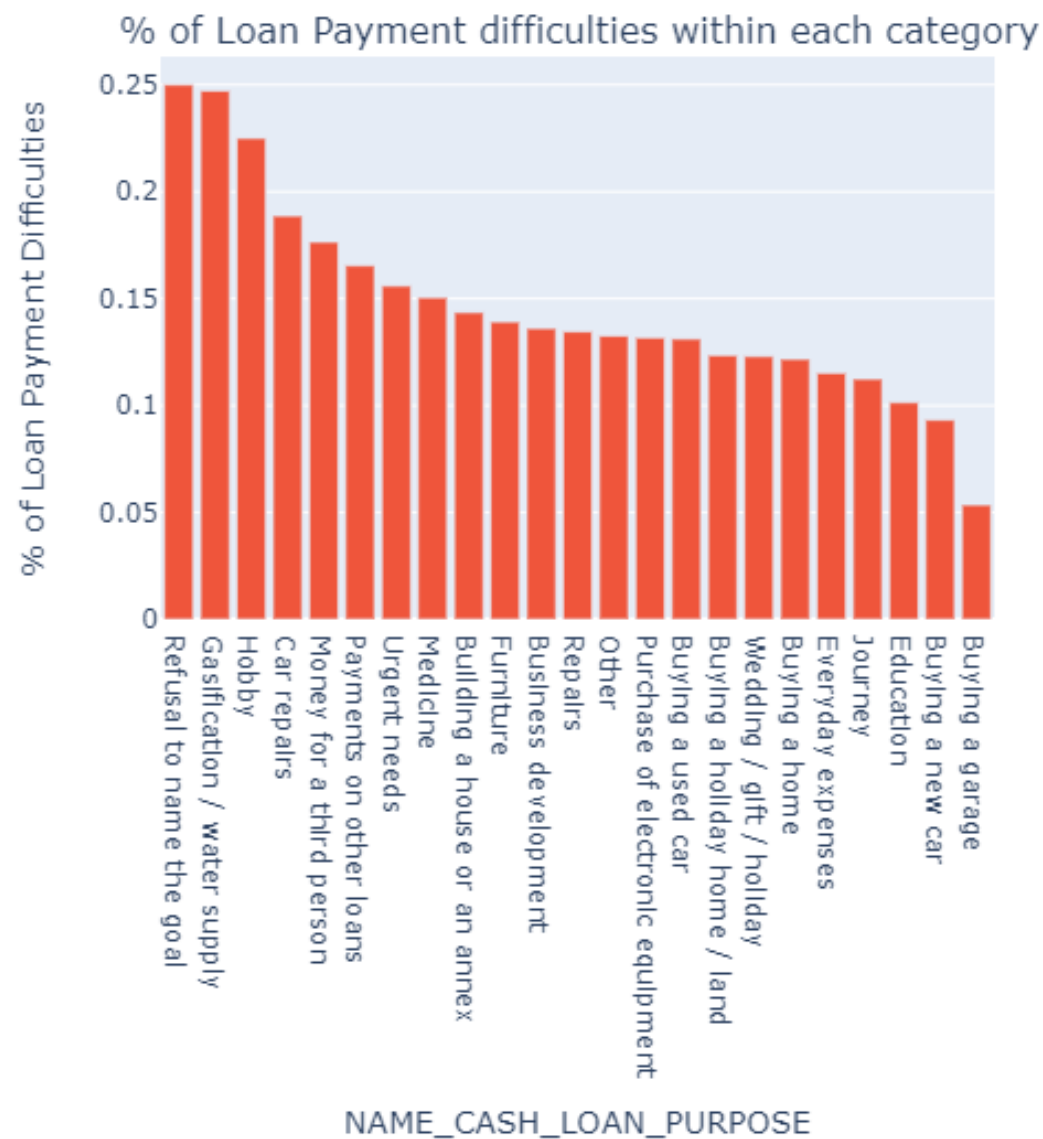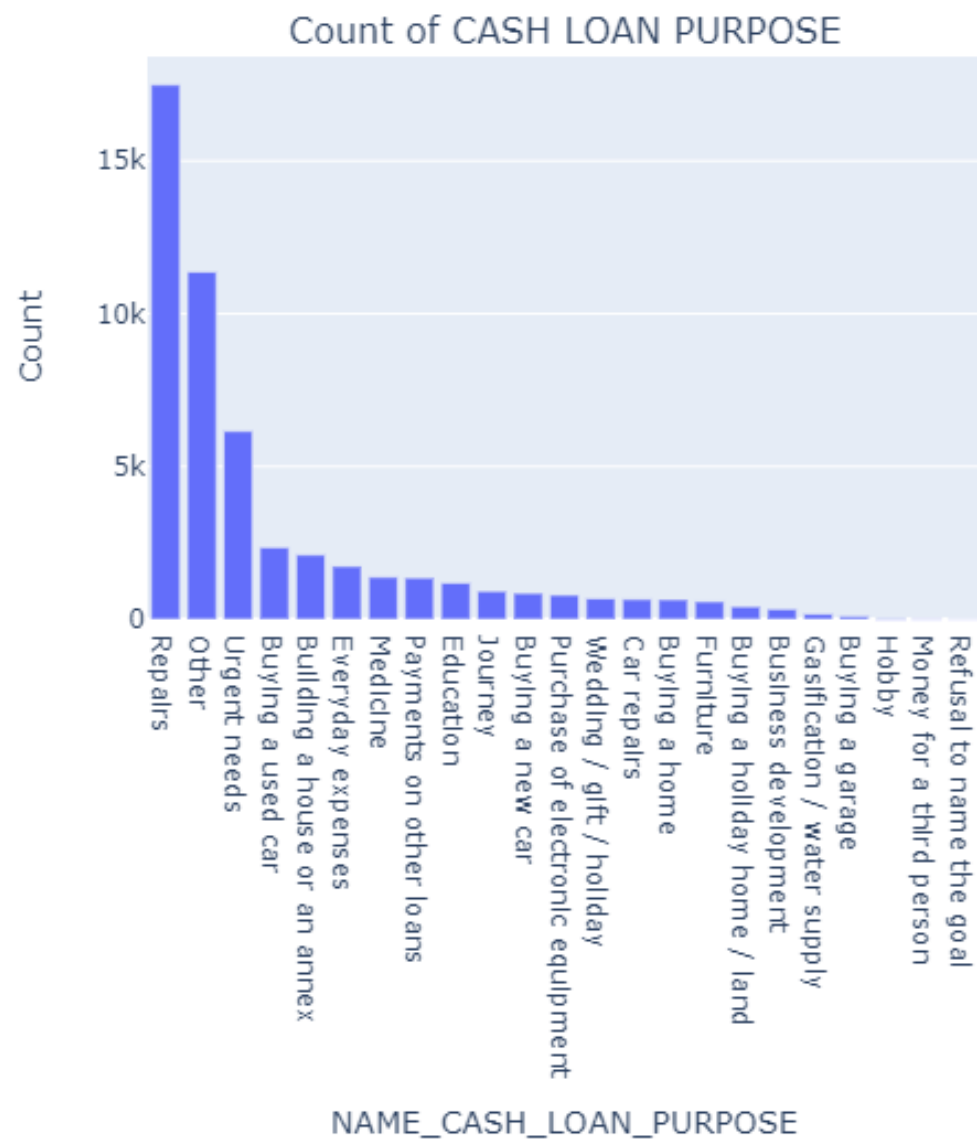
# CASH LOAN PURPOSE

Observation:

From the first graph it can be seen that purpose of cash loan from previous data was maximum for 'Repairs'

It can be clearly seen from the second graph that the

1. 'Refusal to name the goal' for cash loan from previous application are the ones who have maximum % of Loan-Payment Difficulties from current application.

# CONCLUSION

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

- Banks should focus less on income type 'Working' as they are having most number of unsuccessful payments.

- Also with loan purpose 'Repair' is having higher number of unsuccessful payments on time.

# THANK YOU..!!