# Problem Statement /  Use Case

Bank and financial institutions are supposed to provide loans to their customers however they need to take an intelligent decision before disbursing any loan. If the applicant becomes a defaulter then it is a financial loss or if Banks do not offer loan services to their customers then It is a loss of potential business.

The given data contains information about loan applications at the time of applying for a loan, it contains 2 scenarios:
1. Client with Payment difficulties.
2. All other Cases.

# Business Objective

The company wants to find the important column or attributes which can help to identify the possibility of defaulters i.e. variables that are strong indicators of defaulters.

Also, this case study will identify patterns that indicate if a customer has difficulty in paying their installment which may help to take corrective actions and decisions higher interest rates, reduced loan amount, or denying the loan.

# Data Set

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 on whether the previous application had been **approved, Canceled, Refused, or Unused offer.**

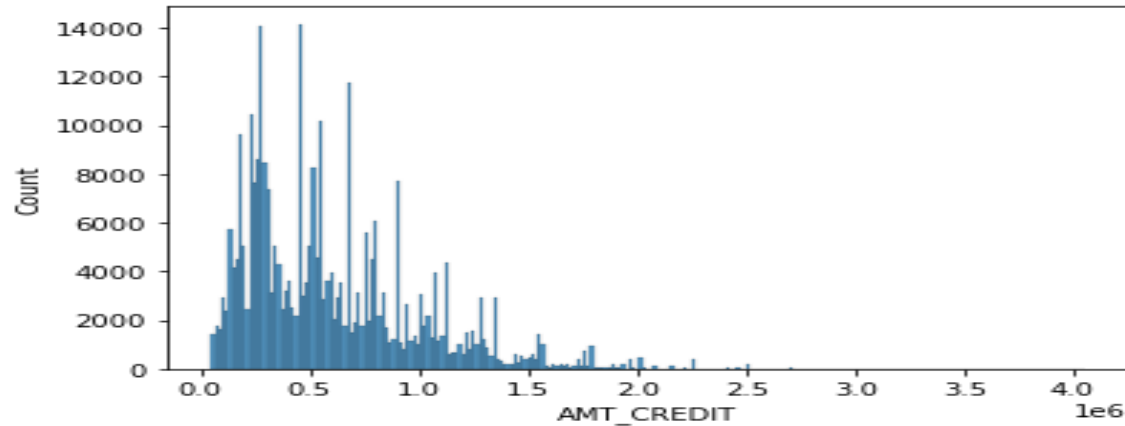3. *'columns_description.csv'* is a data dictionary that describes the meaning of the variables.

# Analysis of Application Data

**Approach Followed:-**

1. Load the application_data.csv
2. Check the structure and the Metadata
3. Find the Missing value % in each column
4. Drop the columns having higher Missing Value %(>45% in this dataset) as higher missing values will impact analysis.
5. Identification/Imputation of Missing values of the remaining columns.
6. Identification/Division of important columns for further Analysis into Continuous and Categorical Columns
7. Outlier Identification/Analysis through Boxplot.
8. Outlier Treatment.
9. Univariate Analysis for Continuous Data
10. Bivariate Analysis for Continuous Vs Categorical Data
11. Finding the Data Imbalance.
12. Segmented Analysis for TARGET
13. Multivariate Analysis for Continuous Data
14. Multivariate Analysis for Categorical Data
15. Assumption: XNA value in Organization Type is considered as Null

# Univariate Analysis for Continuous Columns in Application data

Distribution `of AMT_CREDIT`



Maximum Clients get amount credited is < 50K

Distribution of AMT_GOODS_PRICE



Largest no. of Clients gets loan against 5L
AMT_GOODS_PRICE



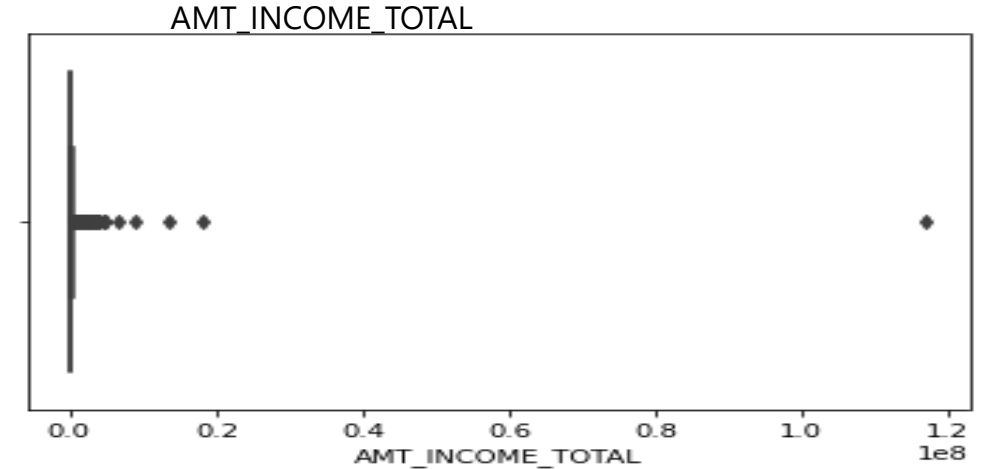Distribution of REGION_POPULATION_RELATIVE



Distribution of Days of birth
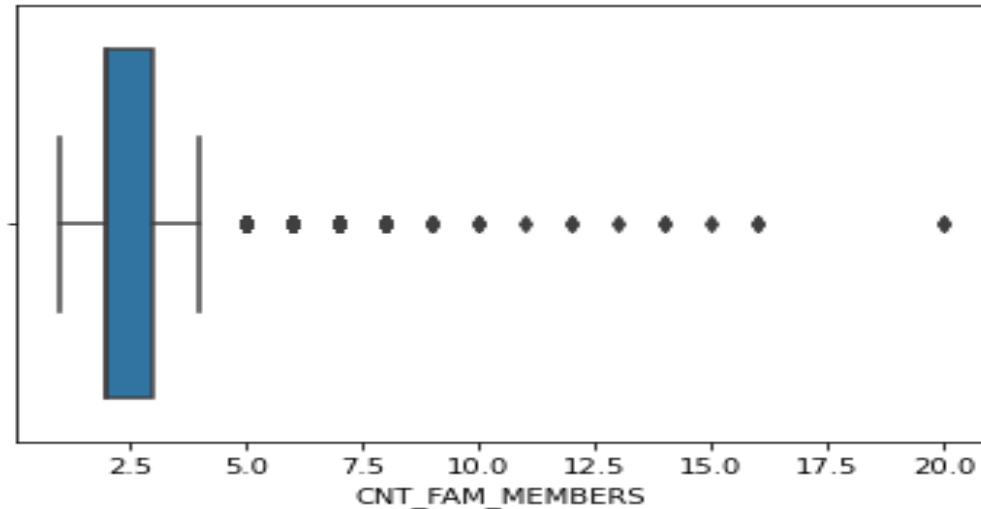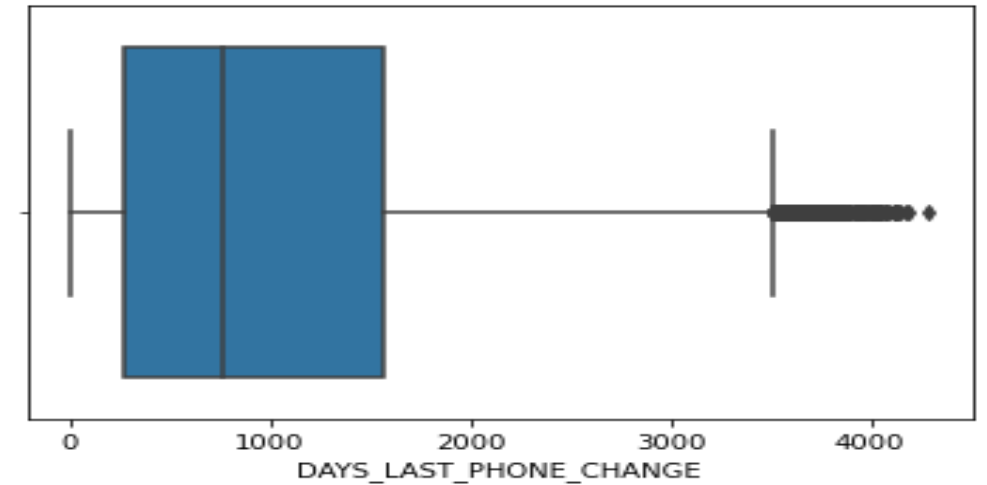
# Outlier Analysis (Boxplot) in Application data



CNT_Children

Boxplots clearly show the values above 2.5 are outliers

AMT_INCOME_TOTAL

Applicants with Income above 900K are outliers.

Applicants with 5 or more family members are clearly outliers

Applicants with `DAYS_LAST_PHONE_CHANGE` above 3514.0 are outliers

# Bivariate Analysis for Continuous Vs Categorical Data



Client staying in Region Ration 1 have higher EXT Source2 data.

Loaned Amount is higher for Businessman

Amount Annuity is higher for highly skilled occupation Type

Amount Credited to Academic Degree is higher
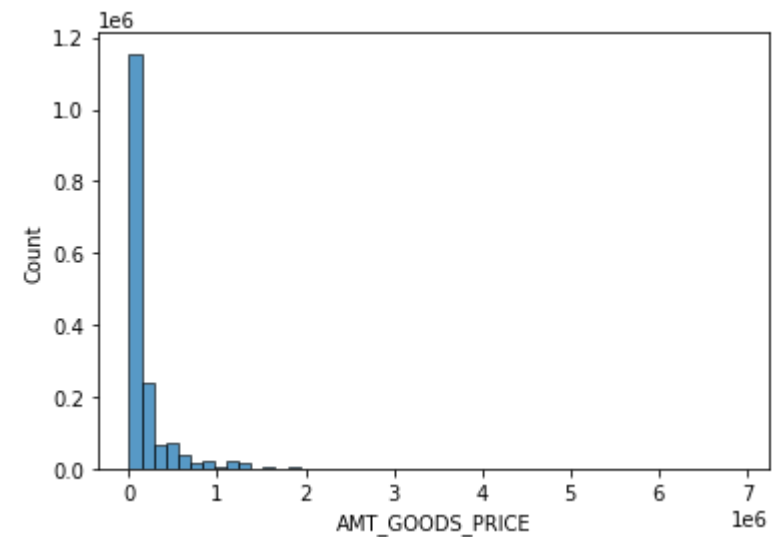
# Bivariate Analysis for Categorical Vs TARGET

(Defaulters/Having difficulties in payment (TARGET==1),Non Defaulters/All other cases (TARGET ==0)



Defaulters are higher in Cash Loans

Defaulters are higher for Secondary/Secondary Special

# Multivariate Analysis for Continuous Columns



AMT_CREDIT has a high correlation with AMT_GOOD_PRICE

AMT_CREDIT has a high correlation with AMT_ANNUITY

AMT_ANNUITY has a high correlation with AMT_GOOD_PRICE

CNT_CHILDREN has a high correlation with CNT_FAM_MEMBERS

# Multivariate Analysis for Continuous Columns –Pair Plot

# Multivariate Analysis for Categorical Columns

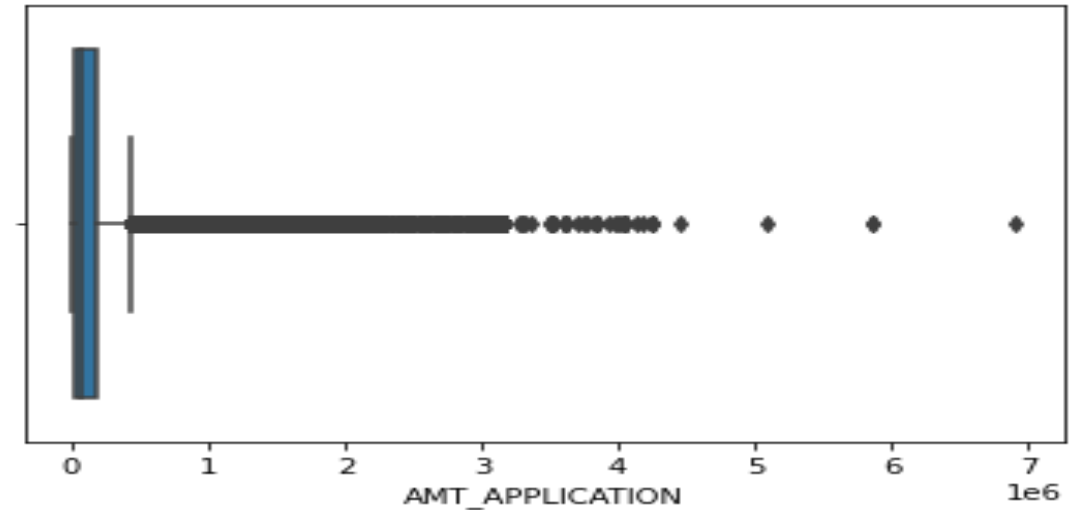# Analysis of Previous Application Data
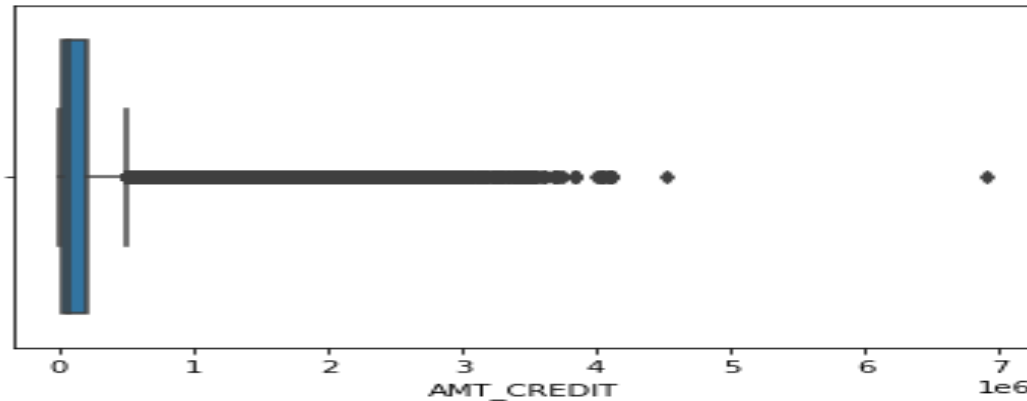
# Univariate Analysis for Continuous Columns in Previous Data
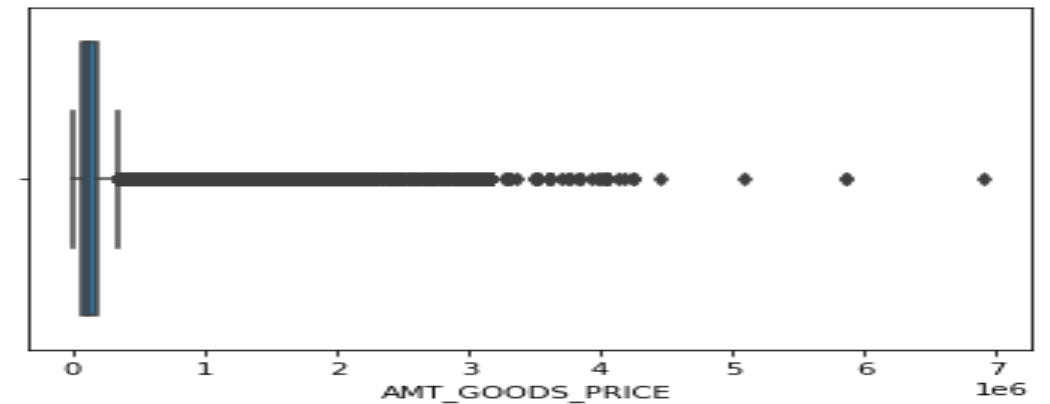
# Outlier Analysis (Boxplot) in Previous Data



Outliers for AMT_ANNUITY
IQR = 9276.930000000002 : Floor = -6368.298750000004
Capping = 30739.421250000007

Outliers for AMT_APPLICATION
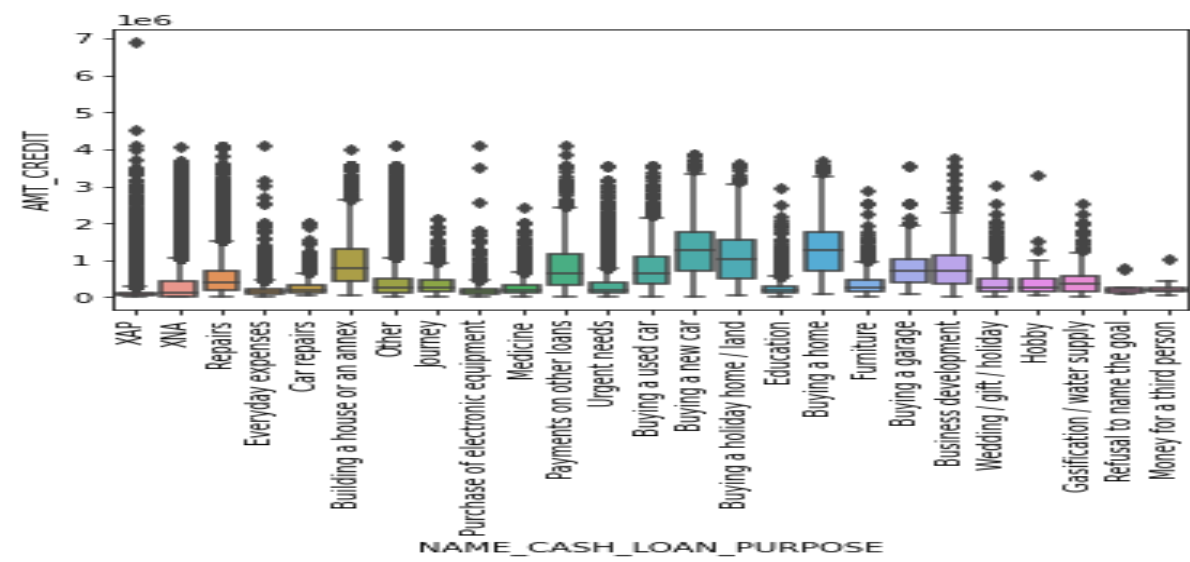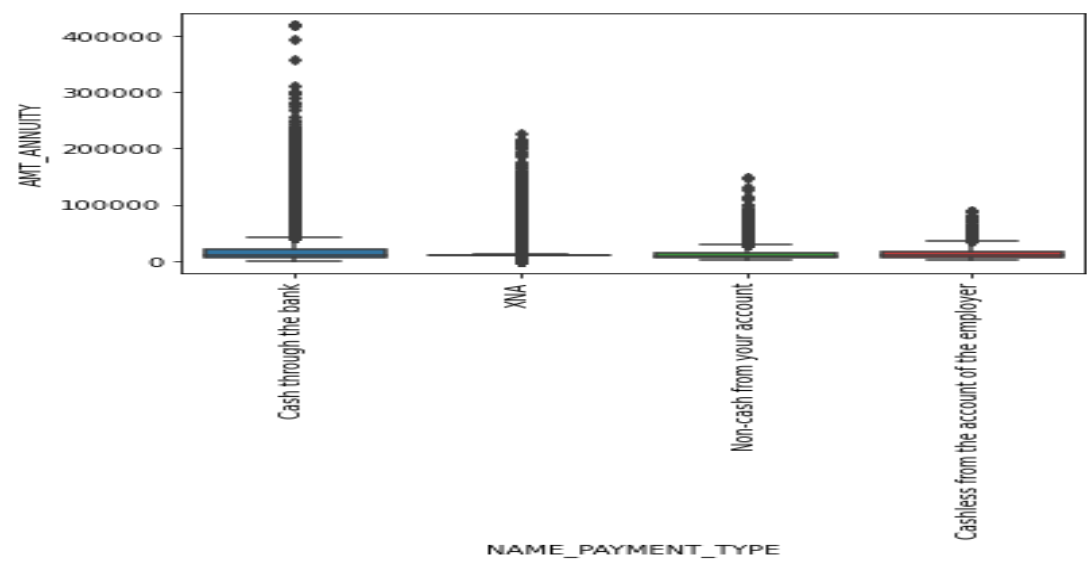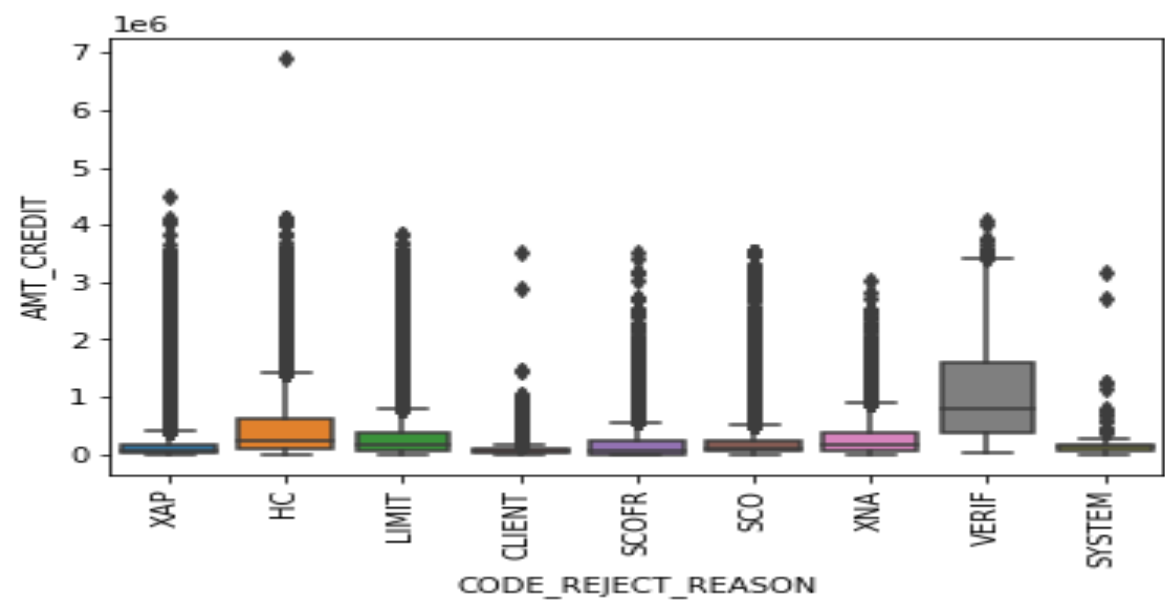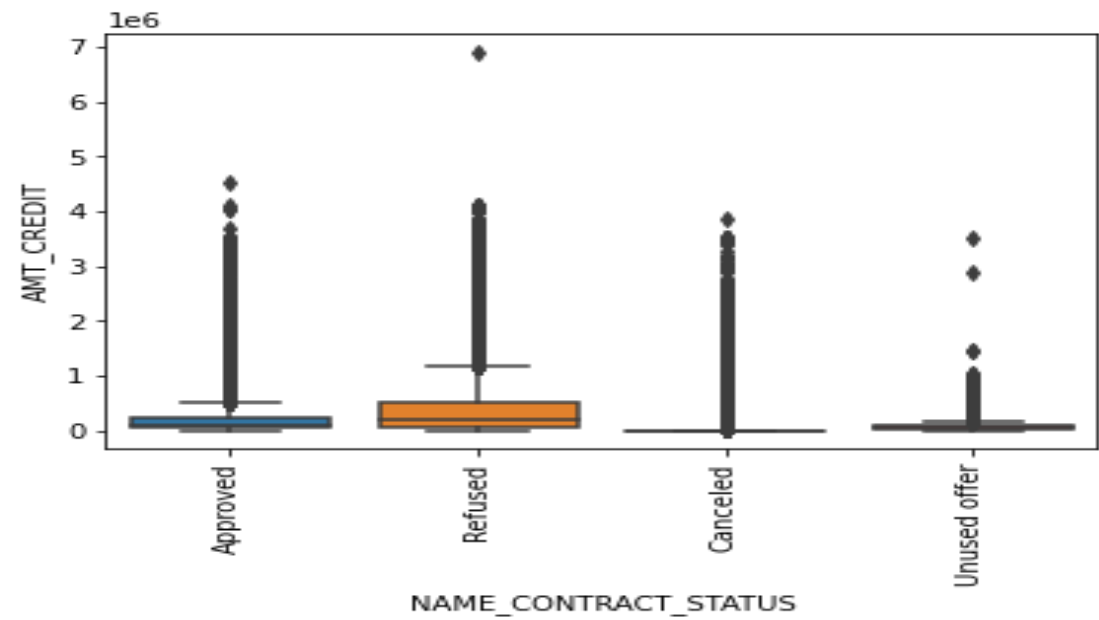IQR = 161640.0 : Floor = -223740.0 : Capping = 422820.0

Outliers for AMT_CREDIT
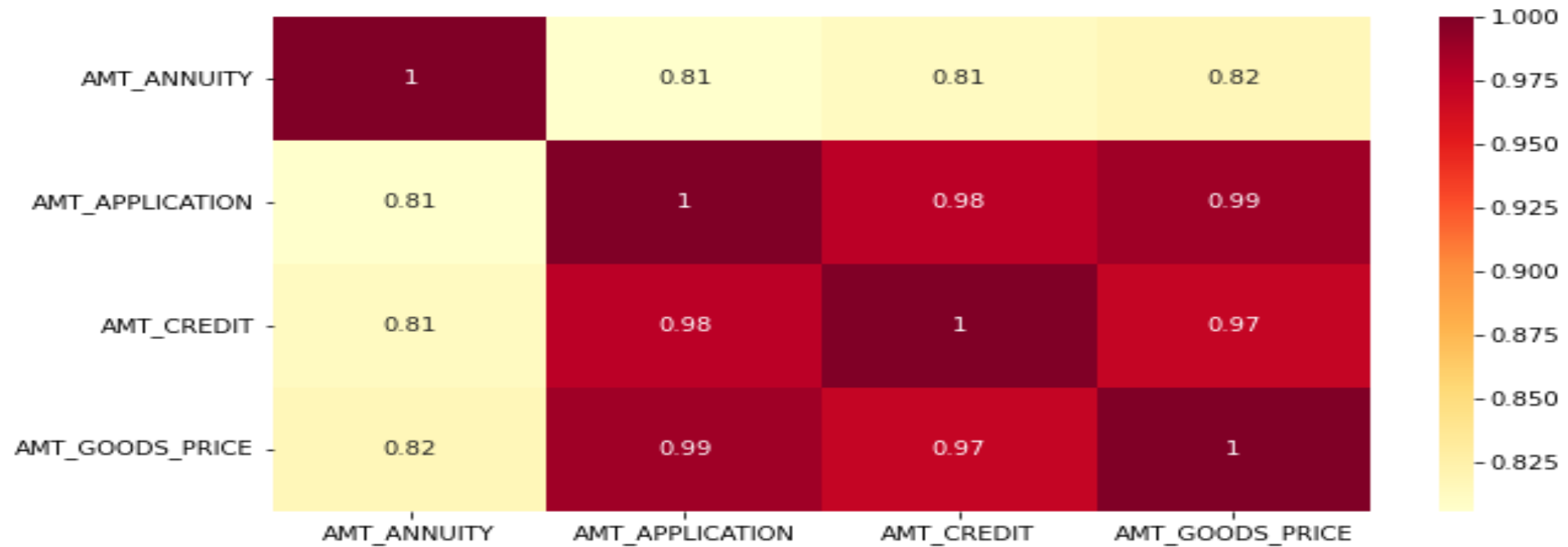IQR = 192258.0 : Floor = -264226.5 : Capping = 504805.5

Outliers for AMT_GOODS_PRICE
IQR = 112905.0 : Floor = -101857.5 : Capping = 349762.5

# Bivariate Analysis in Previous Data

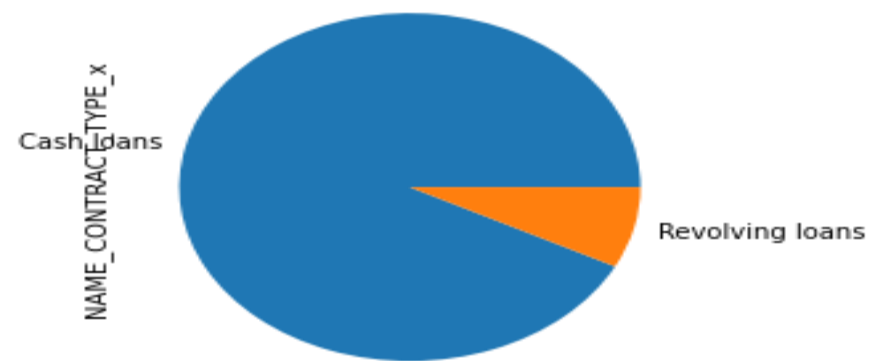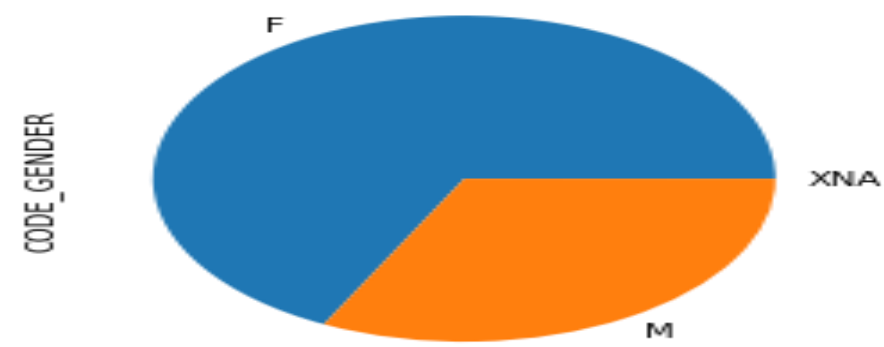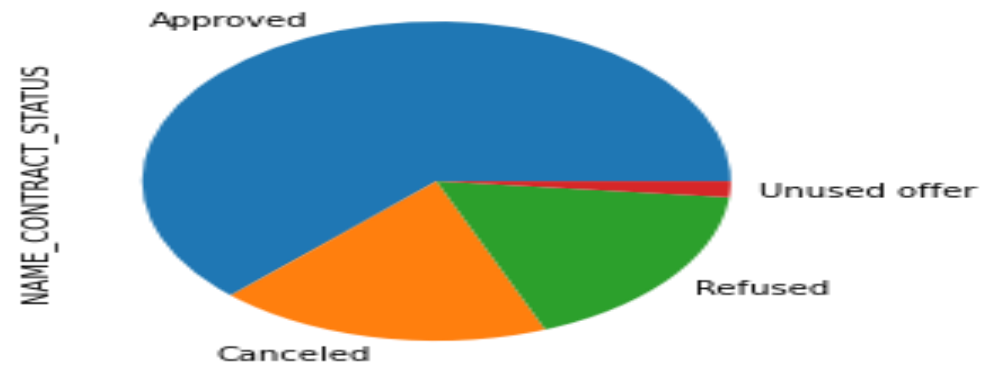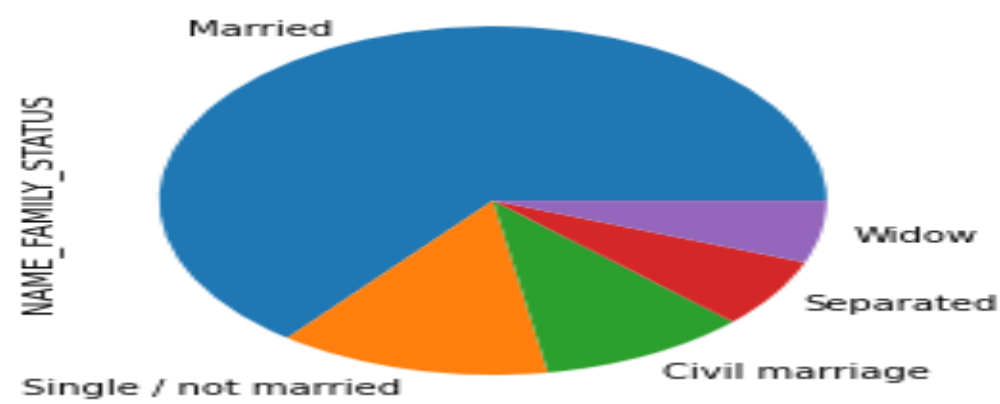# Multivariate Analysis of Previous Application Data



- There is a high correlation between AMT_APPLICATION and AMT_GOODS_PRICE
- There is a high correlation between AMT_APPLICATION and AMT_CREDIT
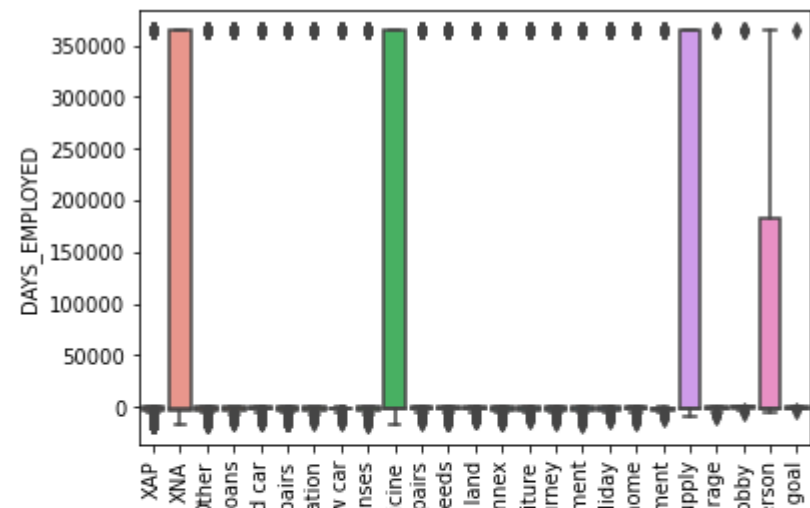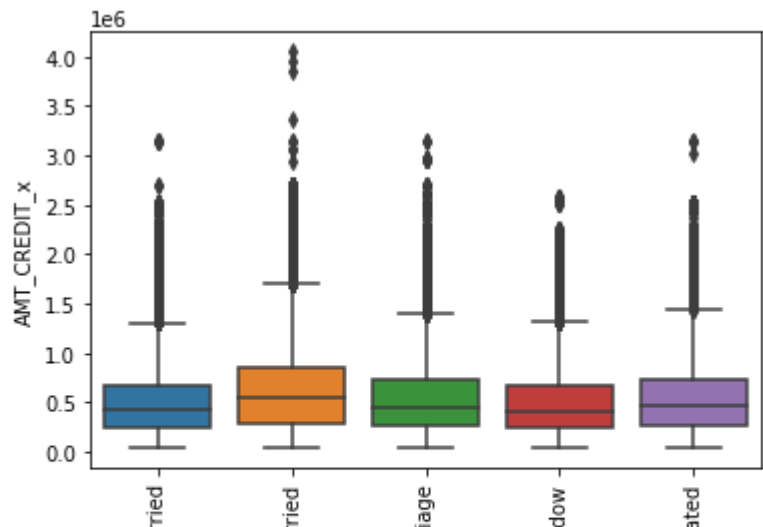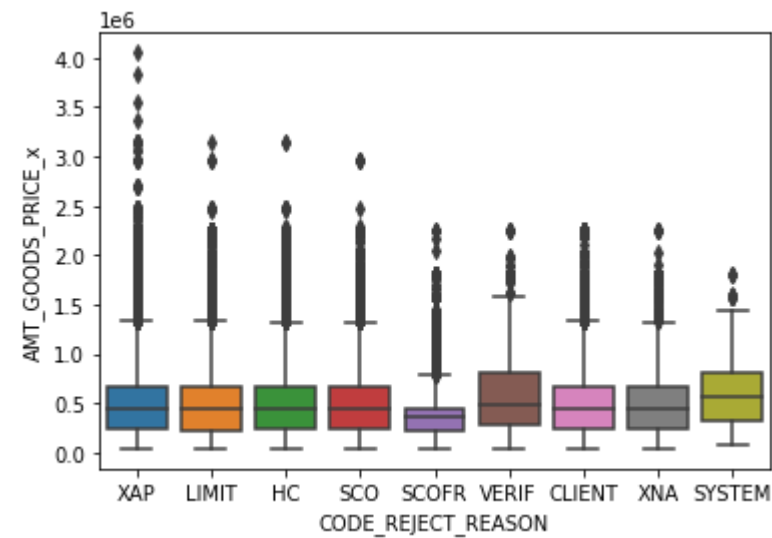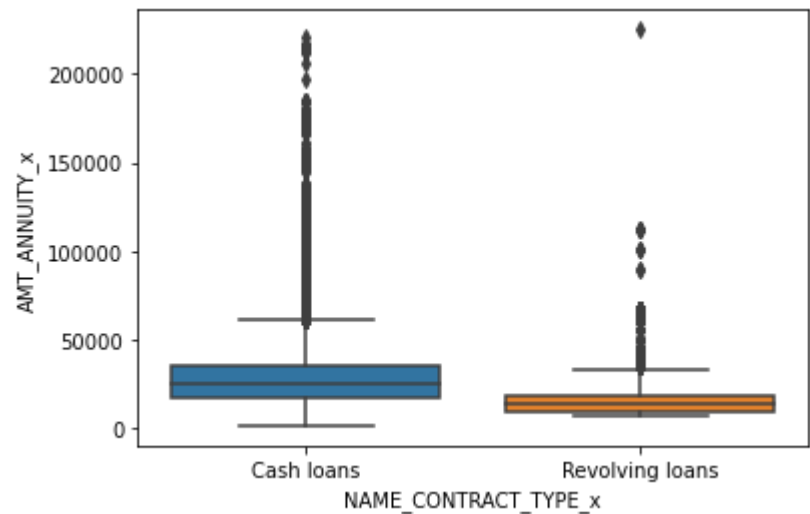- There is a high correlation between AMT_GOODS_PRICE and AMT_CREDIT

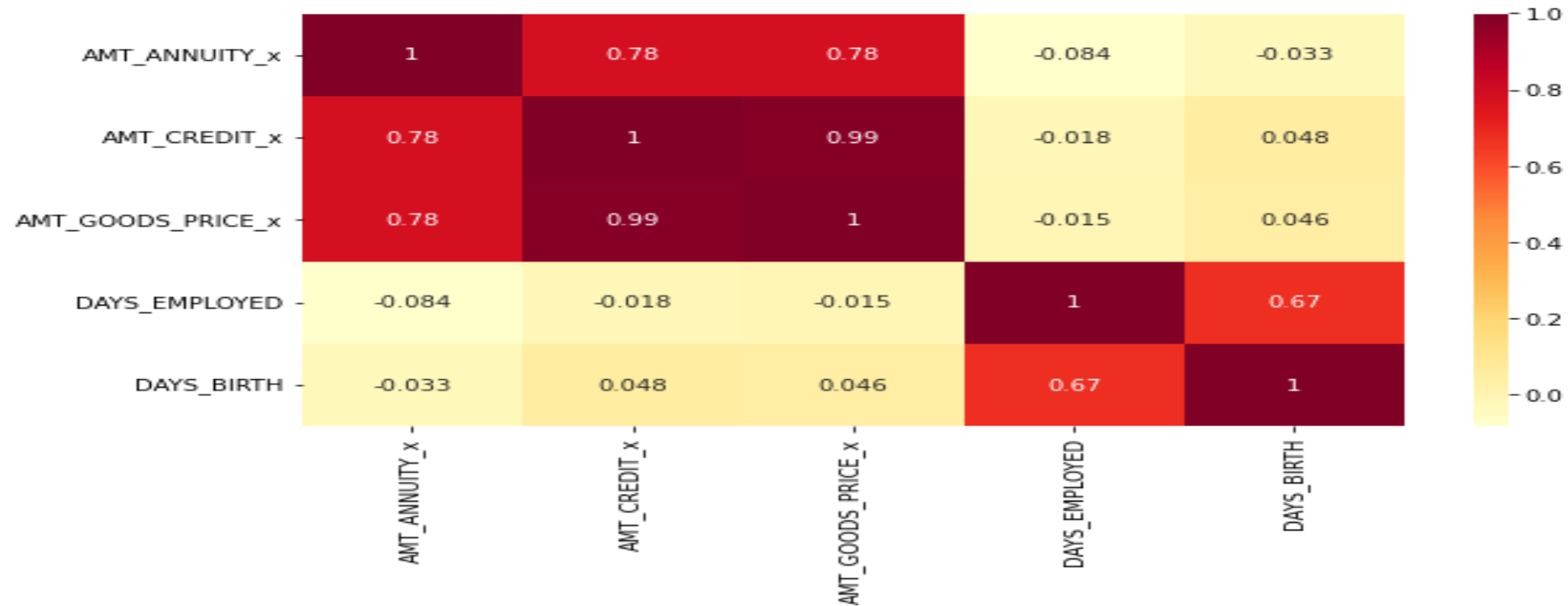# Analysis of Merged Dataset( Application Data + Previous Data)

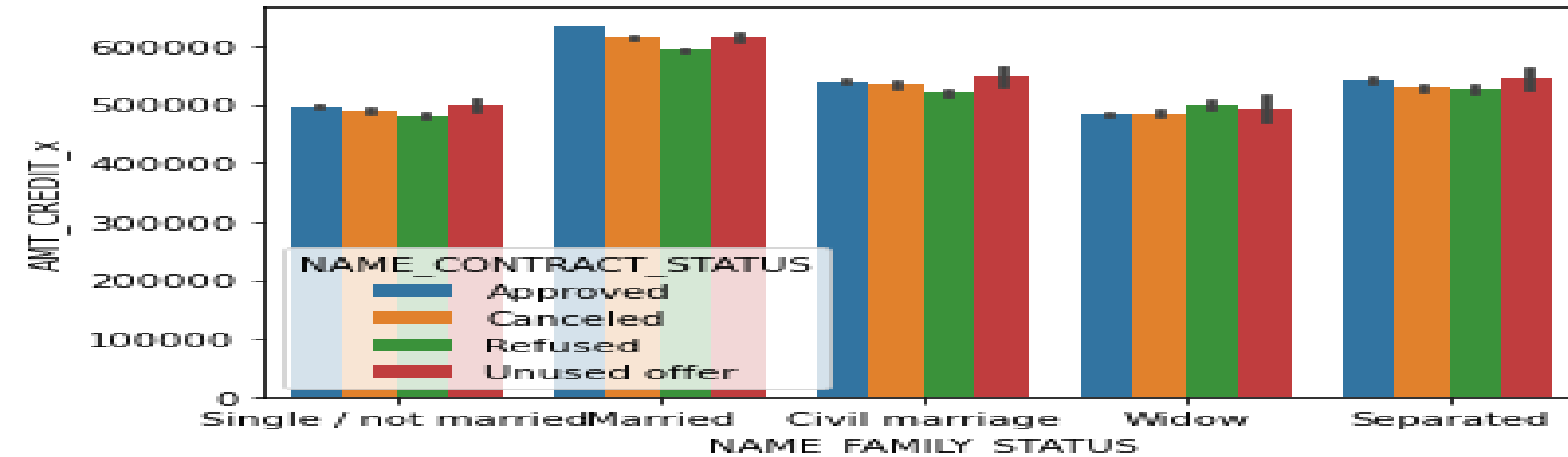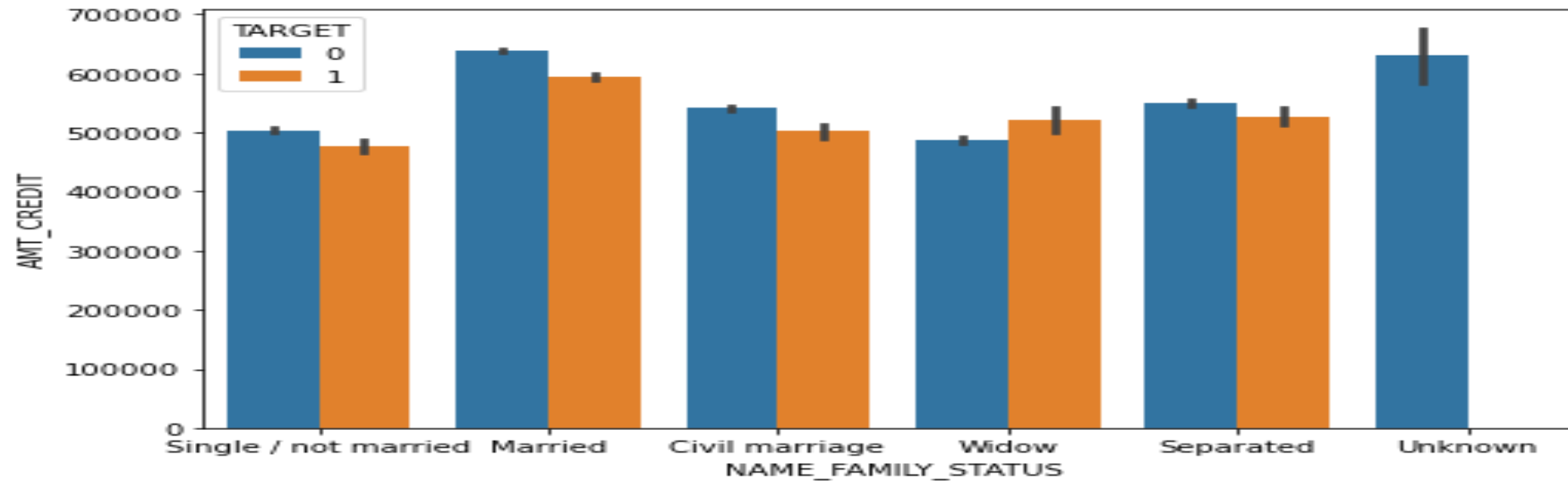# Univariate Analysis of Merged Dataset

Bivariate Analysis of Merged Dataset

# Multivariate Analysis of Merged Dataset


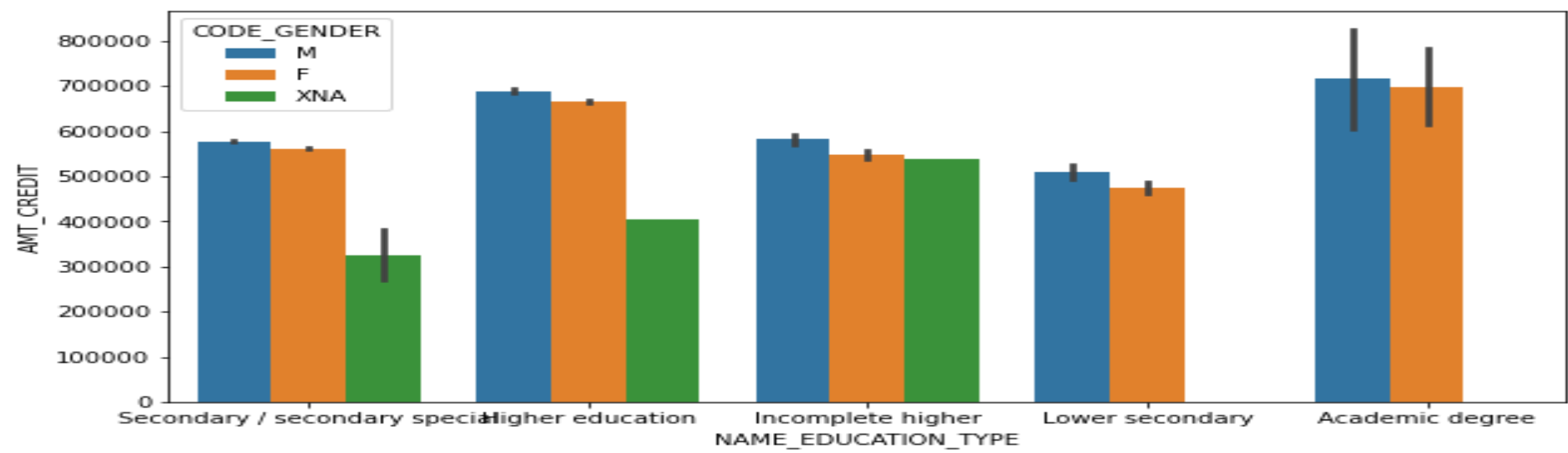
- There is a high correlation between AMT_CREDIT_x and AMT_GOODS_PRICE_x
- There is a good correlation between AMT_ANNUITY_x and AMT_CREDIT_x
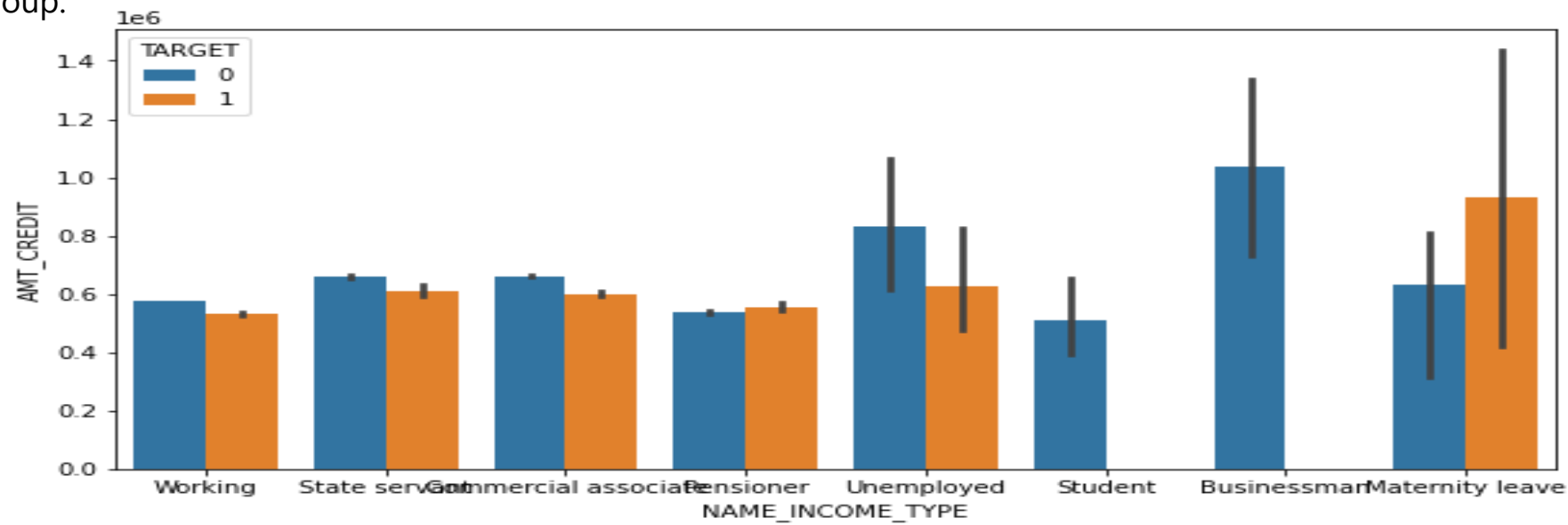- There is a high correlation between AMT_GOODS_PRICE _x and AMT_ANNUITY_x

## Conclusion

- We can approach the **client category as married people** for granting loans.
- It is clearly indicated from the graph plotted that married people have the highest number of the approval.

We can approach the **males with academic degrees followed by males with higher education** for granting loans.



We can approach the **Students & businesses man** for granting loans as we can see that there are no defaulters in this group.

# Additional analysis based on other features