# EDA ASSIGNMENT

CREDIT RISK ANALYSIS

# PROBLEM STATEMENT

Loan Providing company faces challenges for giving loans to people due to there insufficient or nonexistent credit histories.

When evaluating a loan application, the company faces two types of risks:

- If the applicant is likely to repay the loan, rejecting the application results in a loss of business.
- If the applicant is likely to default, approving the loan may lead to a financial loss for the company.

The provided data includes the information about loan application.

- Clients with payment difficulties: Individuals with late payments exceeding a specified number of days on at least one of the initial loan installments.
- Applications where payments are made on time.

For each loan application, there are four potential decisions:

- Approved: The company approves the loan application.
- Cancelled: The client cancels the application during approval, either due to a change of mind or receiving unfavorable pricing.
- Refused: The company rejects the loan application based on client requirements, among other factors.
- Unused offer: The client cancels the loan after approval at different stages of the process.

**Result Expected:**

- Understanding the data.
- Identify the missing data and handle the missing data.
- Identify the potential outlier.
- Identifying an imbalance in the data.

- Sanity checks.
- Univariate  analysis .
- Bivariate analysis.
- Mutivariate analysis.
- Suggestion
- Merging the two datasets.
- Analysis on combined dataset.
- Final Reccomendation.

**Understanding Data** = understanding every column of the data , identifying the target variable in *application_data.csv and* target variable in *previous_application.csv.*

- *The target variable in application_data.csv is* **TARGET.**

- *The target variable in previous_application.csv is* **NAME_CONTRACT_STATUS.**

**Identifying the missing data and handle missing =** We can identify the missing data by using syntax.

**(data.isnull().value_counts()/data.shape[0])*100**

- **data** here represent the name of the dataframe, data.shape give the length of row and then multiplying by 100 give percentage of null value in each column.

- We have the option to remove columns with 40% or more null values, and approximately 0.567% of the data can also be deleted.

**Checking for Outlier  =** We can check the outlier by using boxplot. Identifying the potential outlier.

- We can drop certain percentage of data point with extreme values from the extreme ends
- We can bin the data into group or intervals
- If outlier due to missing value, we can impute or replacing them with appropriate values.

**Identifying imbalance in data =**  We can identify the imbalance of data by examining the distribution of the target variable.

- For calculating  imbalance in ratio we can divide the number of instances in the minority class by the majority class

- We can use visualizations such as bar plots, pie charts, or histograms to represent the distribution of classes.

**Sanity Checks =** We can do some basic sanity checks whether is there any irregularities  in the data.

**Univariate Analysis =** univariate analysis means  analysis  on single variable.

- We can use Visualization  like line chart, boxplot, countplot, histogram.

**Bivariate Analysis =** Bivariate analysis means analysis on two variable

- We can use visualization like Scatter plots, Line charts, Bubble charts , Heatmaps,  Stacked Bar Chart, Grouped Bar charts, Pair Plots

**Multivariate Analysis = Multivariate Analysis means analysis on more than two variable**

- Visualization we can use Heatmap,pairplot etc.
-  Suggestion Any suggestion after analysing the partical data
- Combining Both the dataframe we can join the two dataframe and analyze the Data
- Final recommendation

# Graphs and Insight

- Checking the data imbalance we can defaulter is less then the non

- defaulter and the imbalance ratio is 11.6%

- Females are more prone to applying loan then female
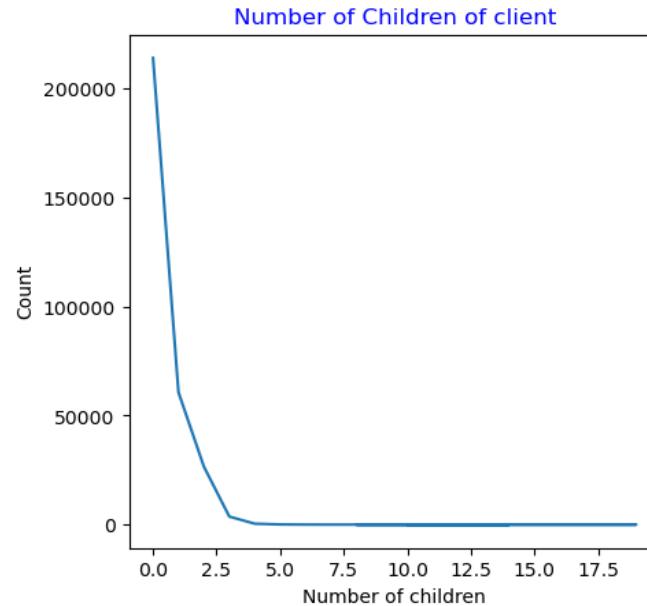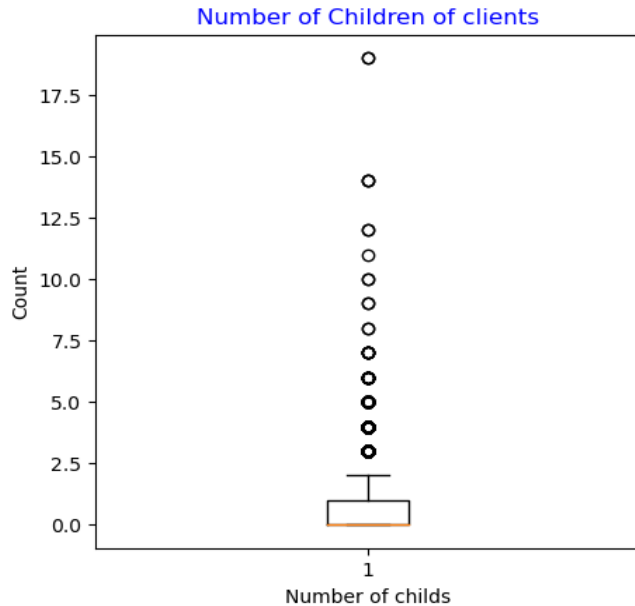
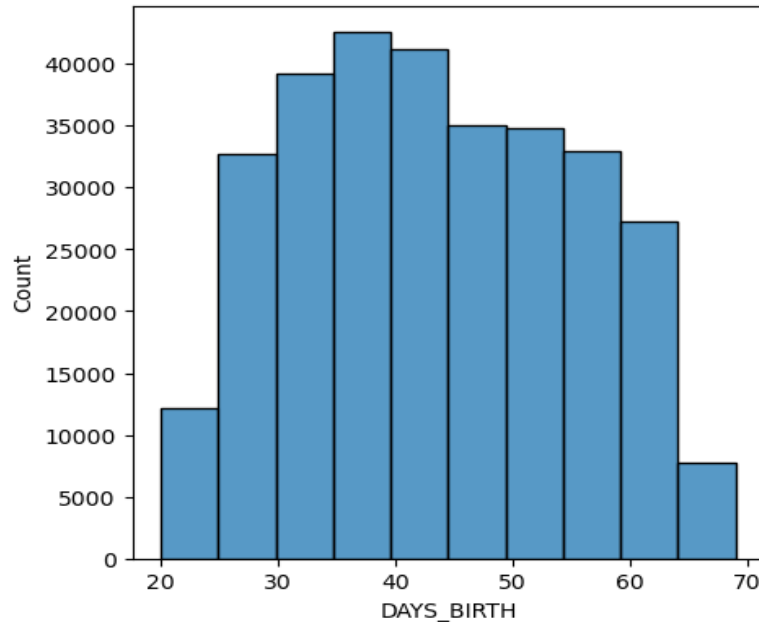Data Imbalance

Distribution of Gender for applying Loan

# Children Counts

- Most of clients having 0 children



- Most of clients belong to age group between 30 to 40 yrs, this age group is more to default in loan payement.
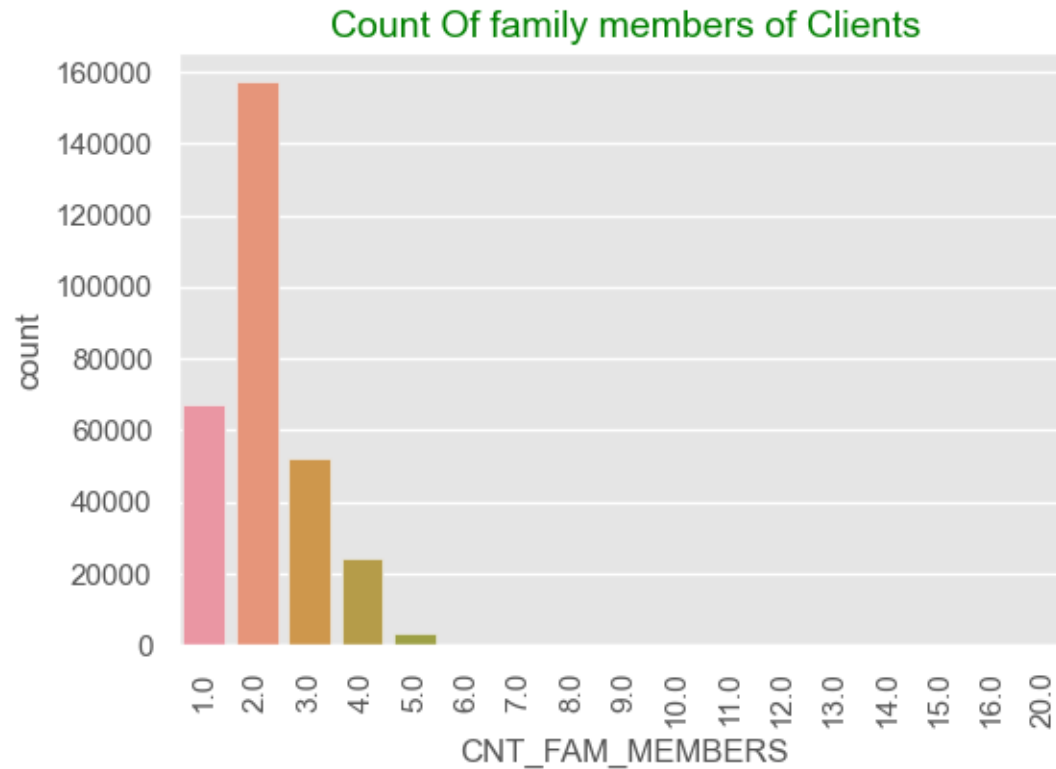
# Distribution of data

- Most of the client have occupaton type laborer.



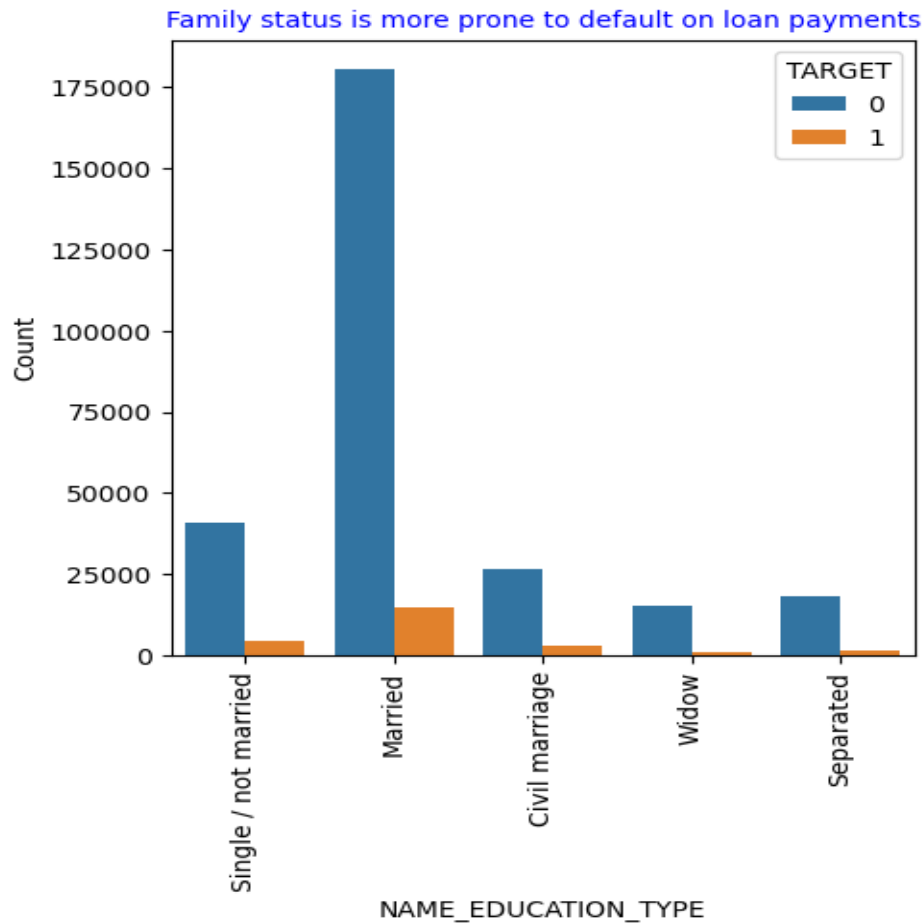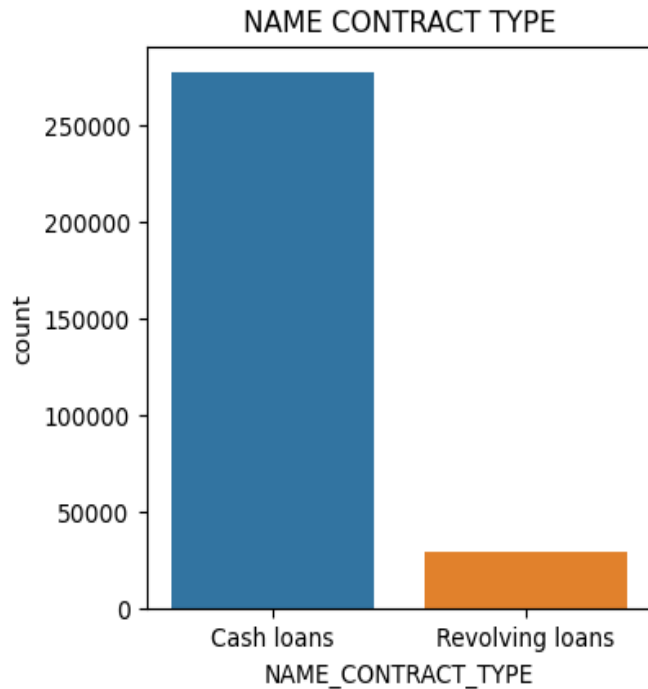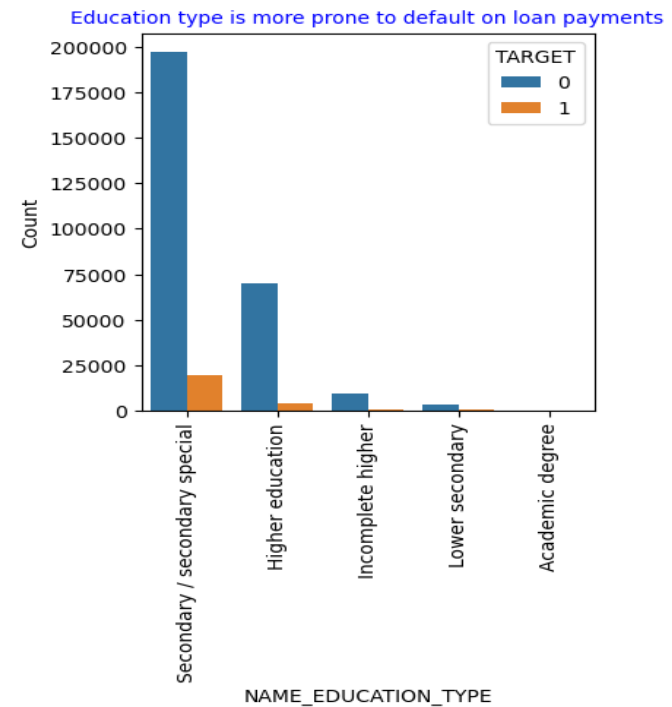Count Of family members of Clients

- Count of the family mostly clients are having 2 members

- Most of the client have family status ids married and the married peoples are the one who are more prone to take the loan.



Family status is more prone to default on loan payments

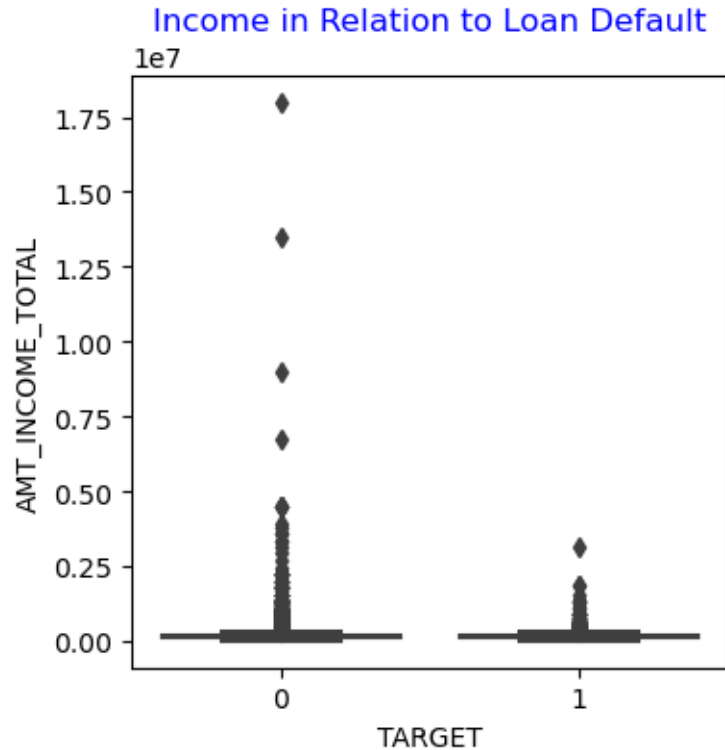- Most of client have contract type is cash loans.

- Most of the client have education type is Secondary
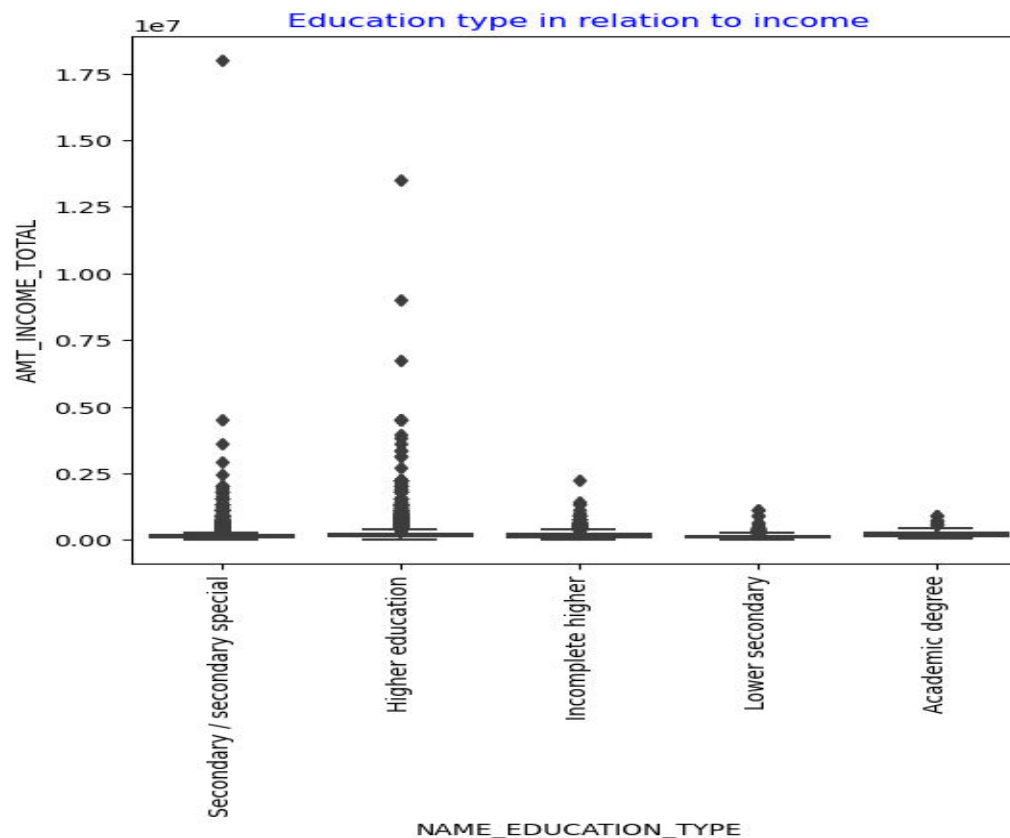
# Income and Credit in relation to default

- Low income people are face more to pay there loan.
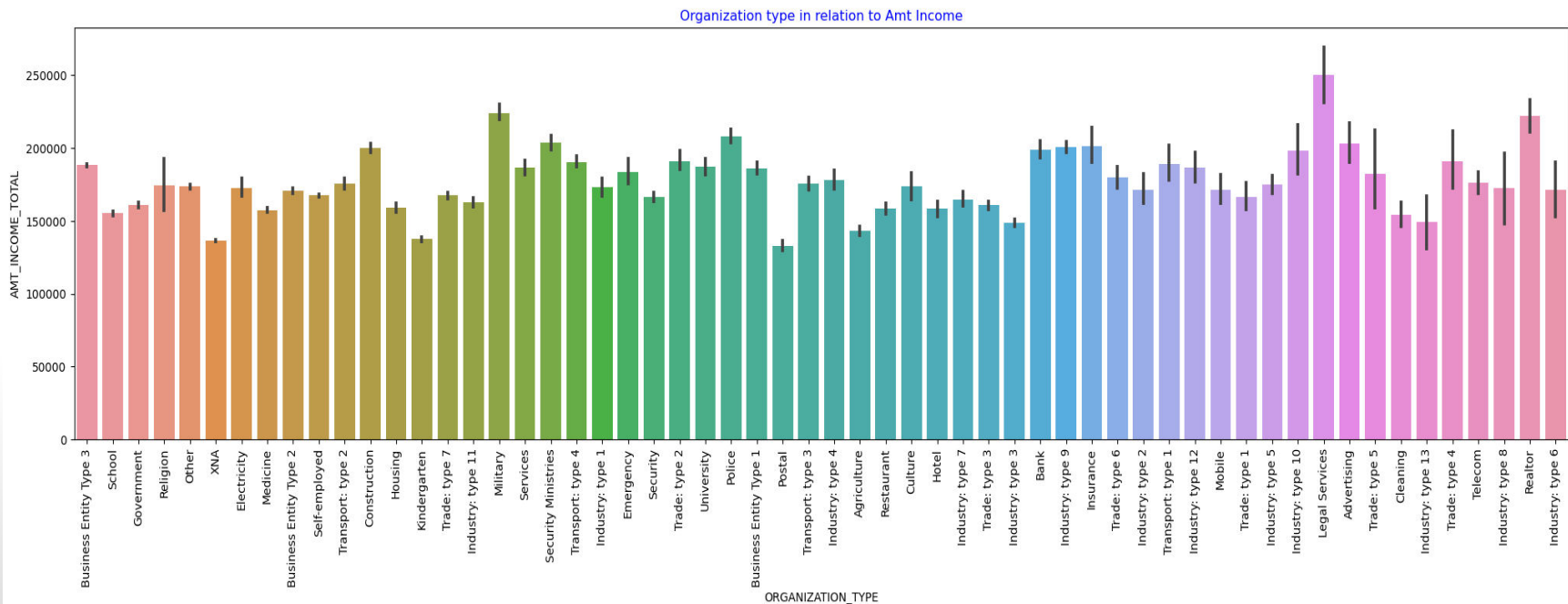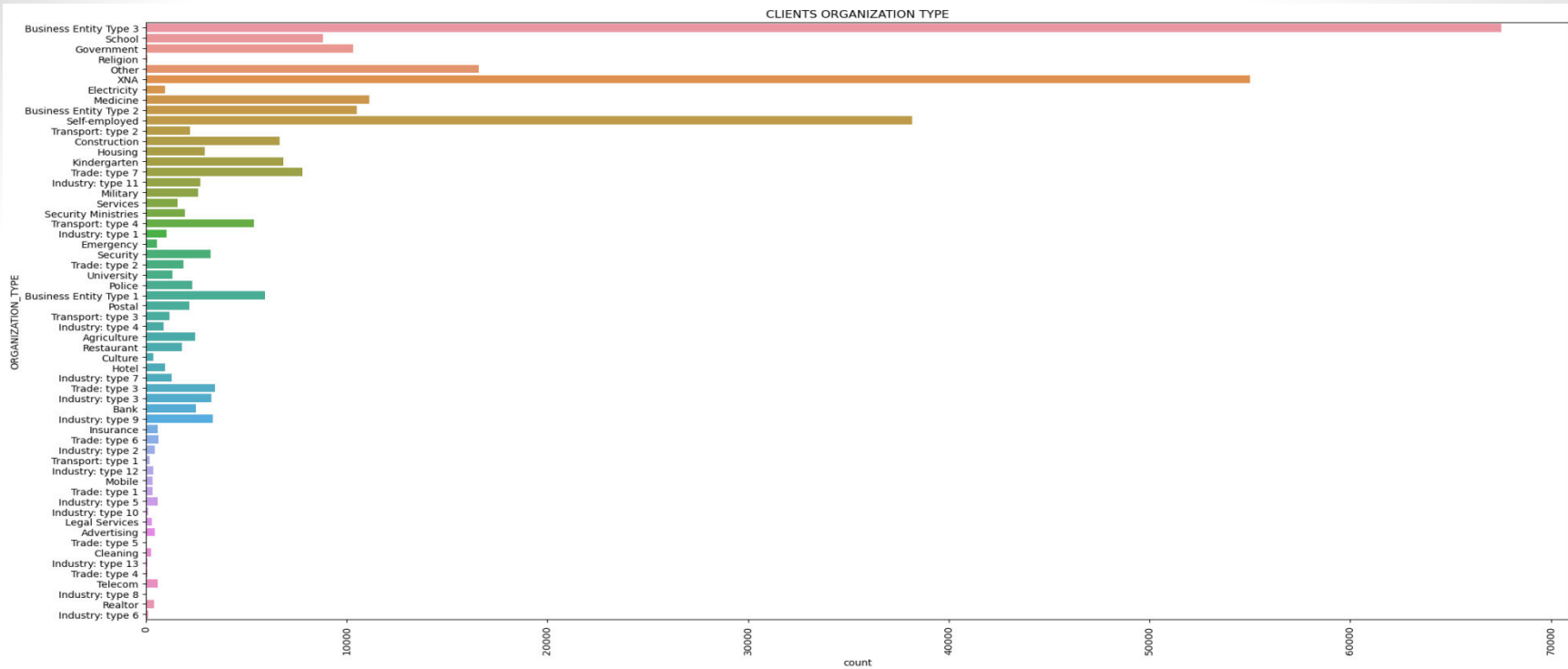
Income in Relation to Loan Default

- Credit amount in relation to default of payment

- Relationship between education and income secondary and higher educated clients are earning more.
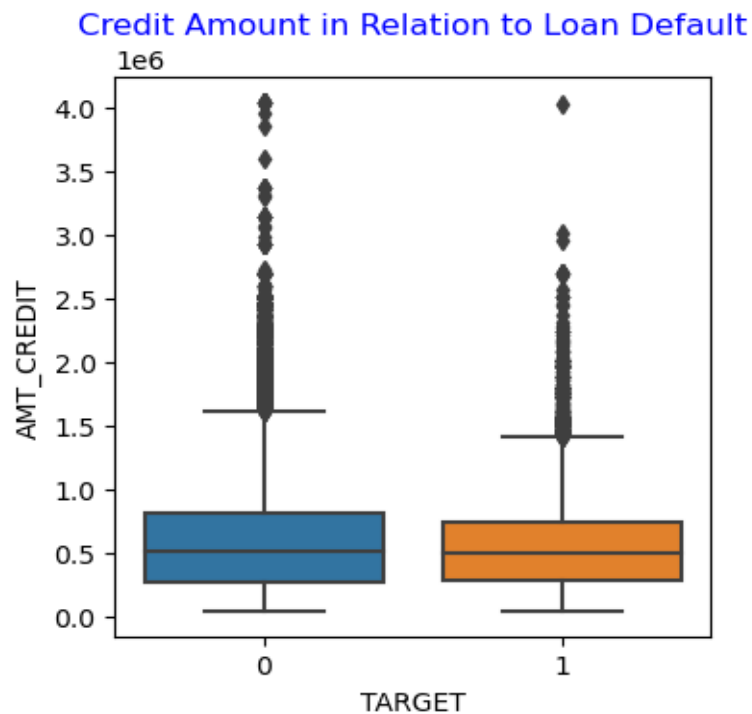


- Relationship between organization type and income and which organization type are more prone to take loan

CLIENTS ORGANIZATION TYPE

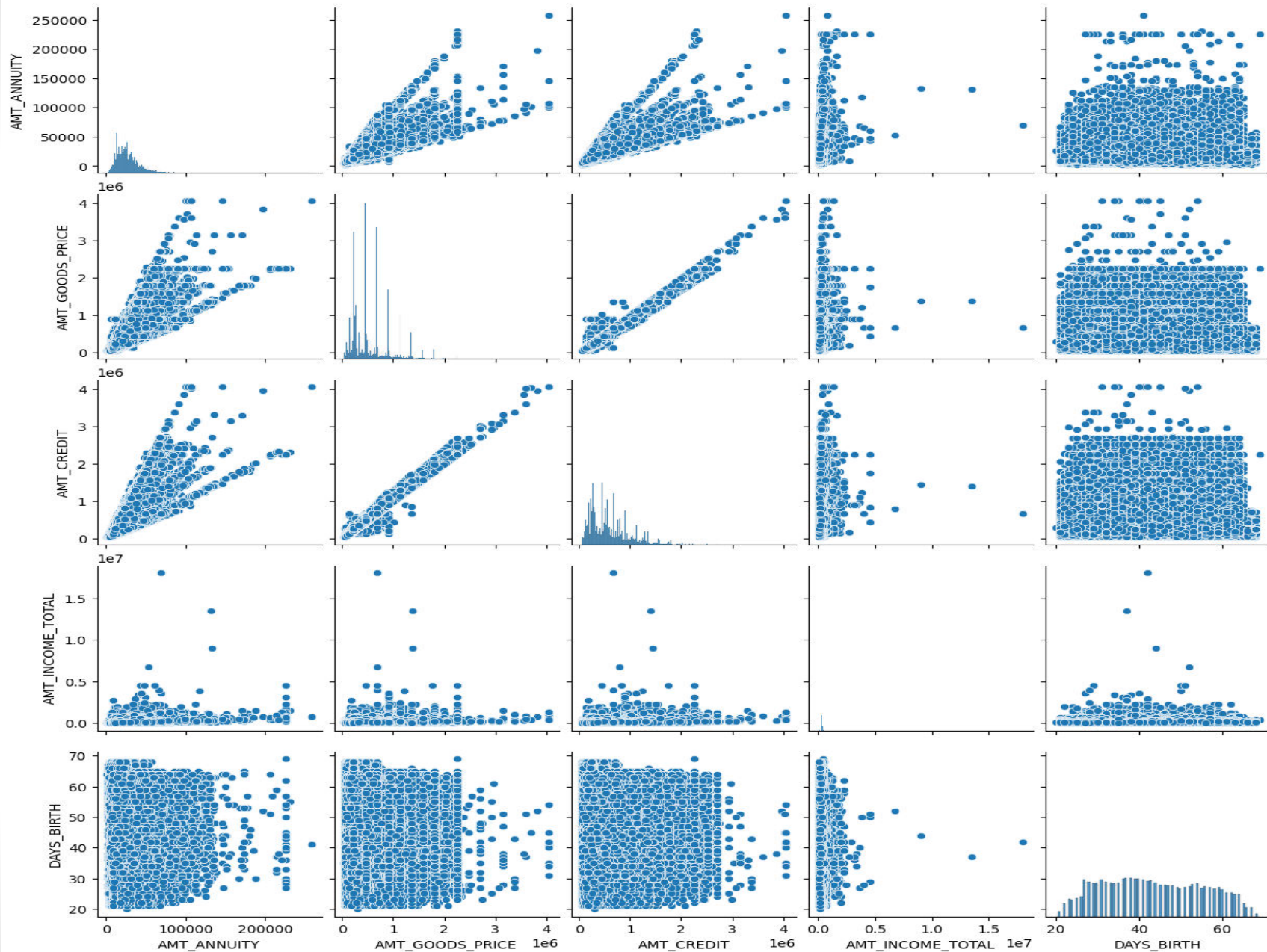Organization type in relation to Amt Income

- Relationship between credit and amount to loan default



- Comparision of all the numerical variable of the application data.

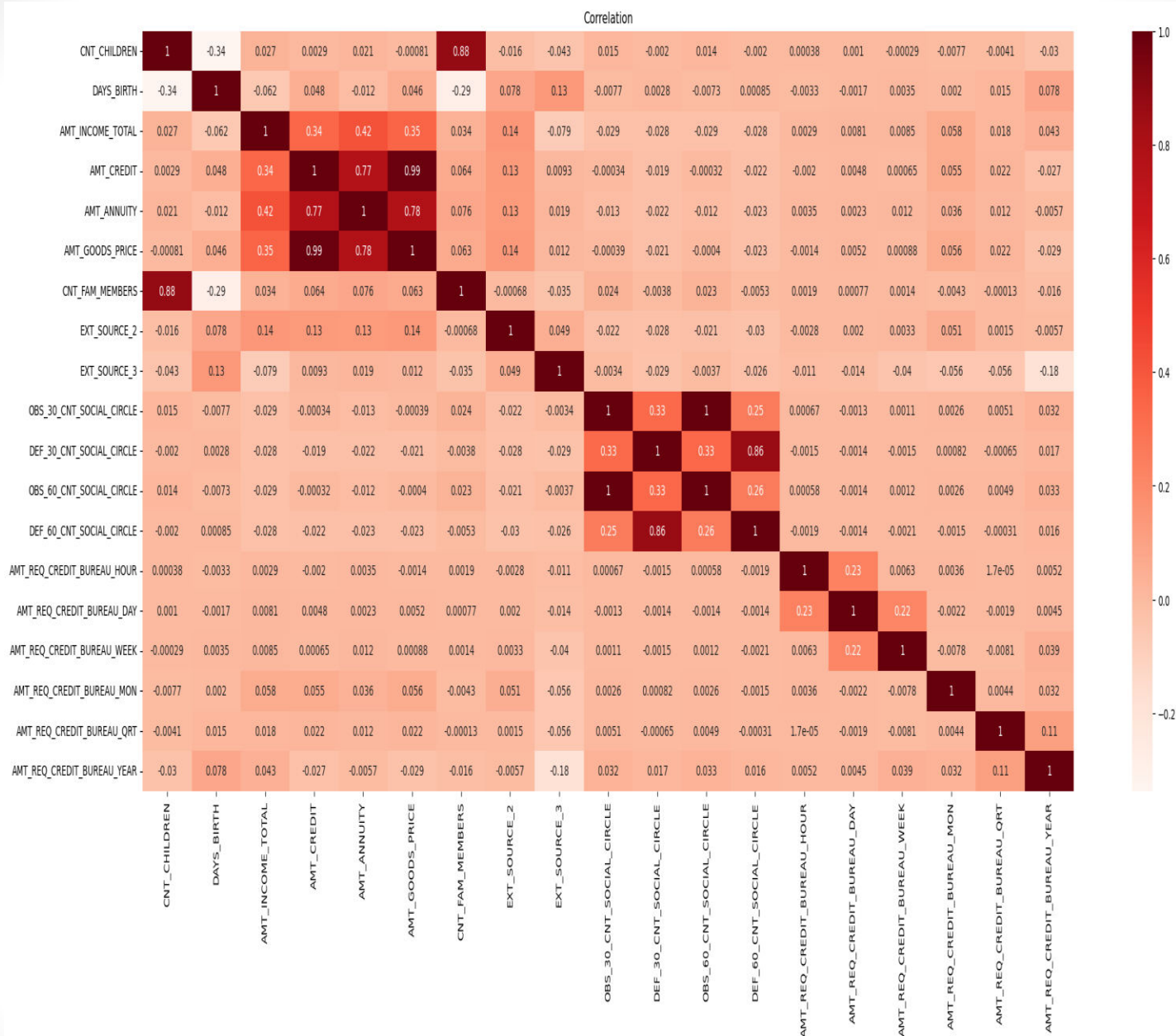Analysis of amt annuity,amt goods price, amt credit, amt income total
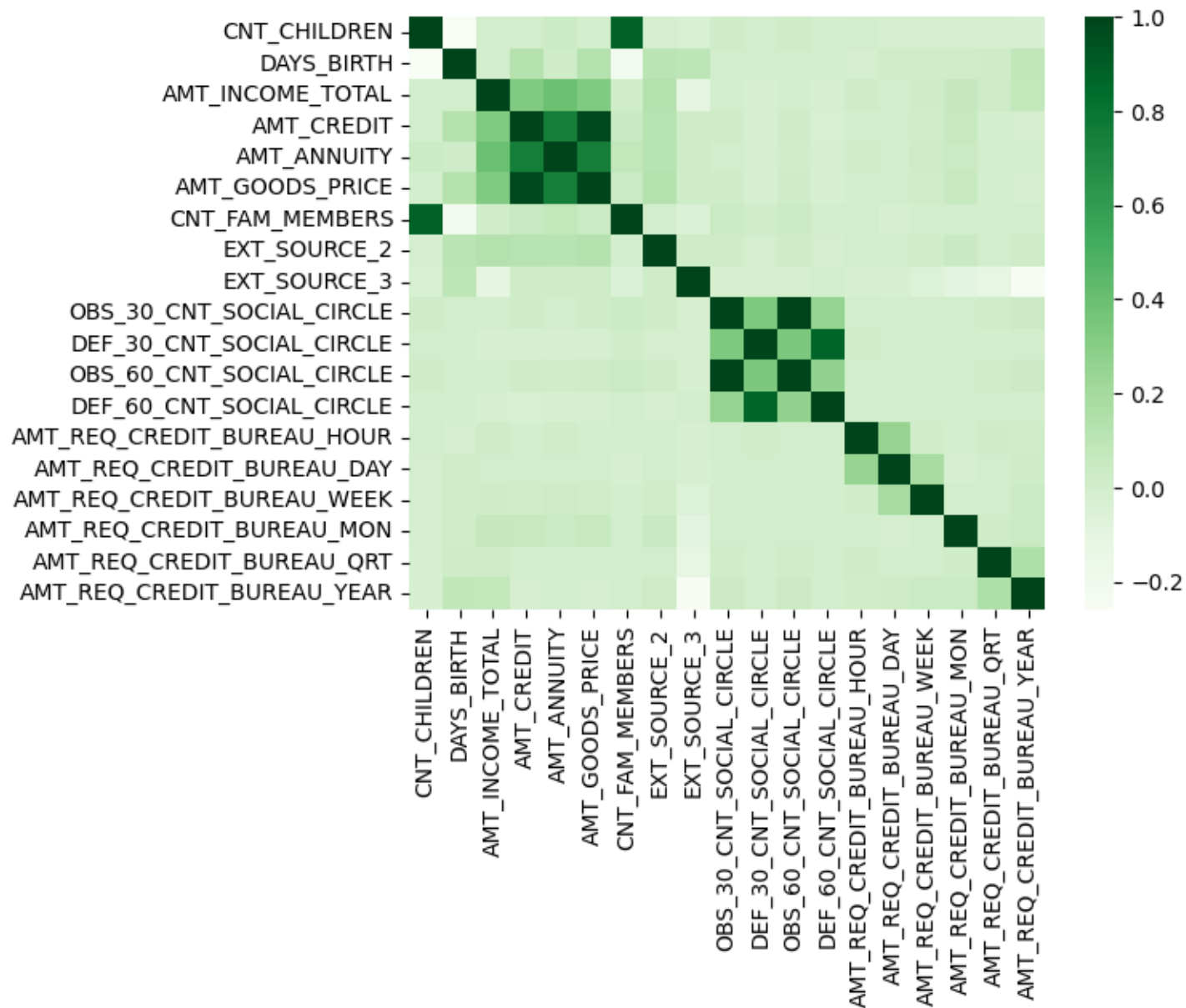
- What we infer from above pair plot

  - The correlation between an increase in the annuity amount and an increase in age is not very apparent in other words it is not   evident higher the amount annuity higher the age.

  - AMT_ANNUITY and AMT_INCOME have very weak correlation but AMT_CREDIT(Credit amount of the loan) and AMT_GOODS_PRICE(it is the   price of the goods for which the loan is given) have strong linear correlation
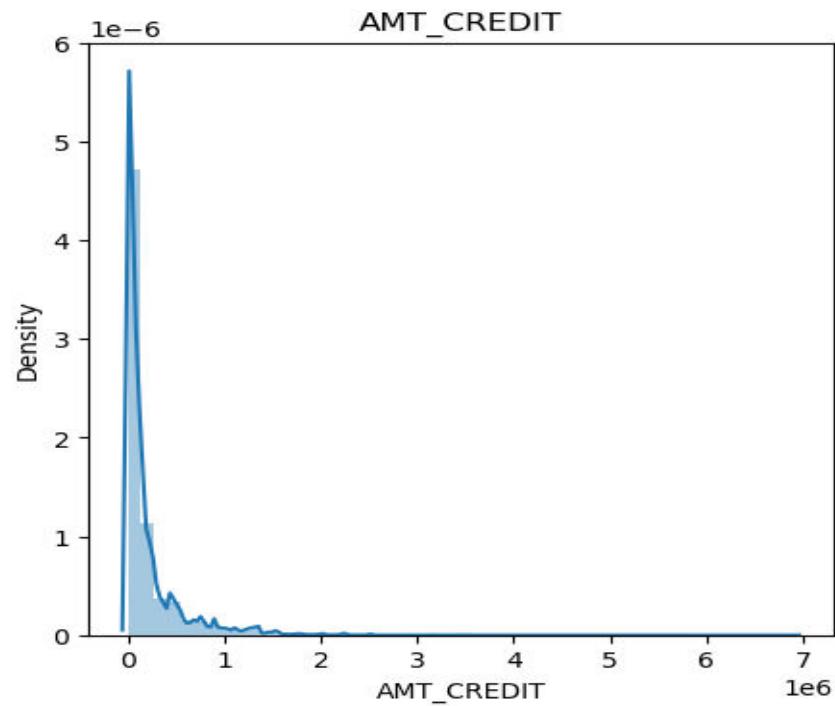
  - Increase in good price and increase in credit amount increase in amount of annuity.

- Finding the correlation for on the basic of target variable
- Divided the targeted variable into two parts Defaulters and no defaulters
- Defaulters represent by 1.
- Non Defaulters Represent by 0.
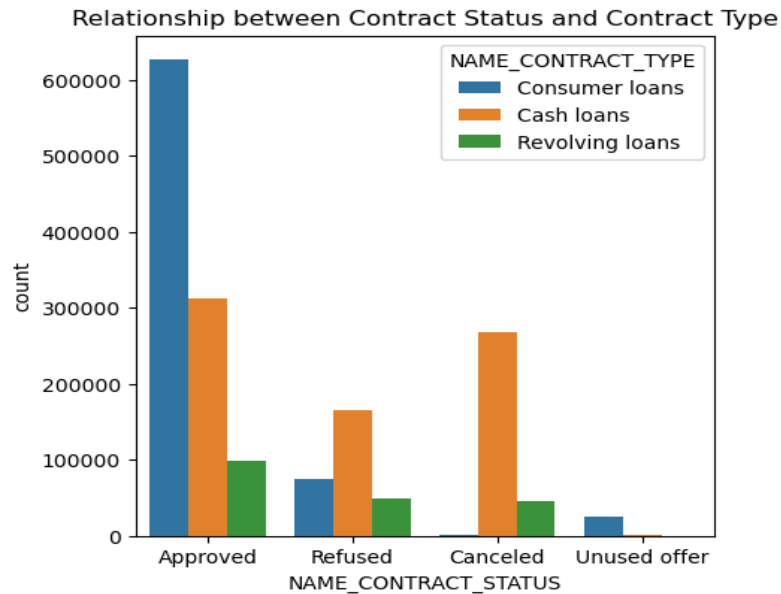
Correlation

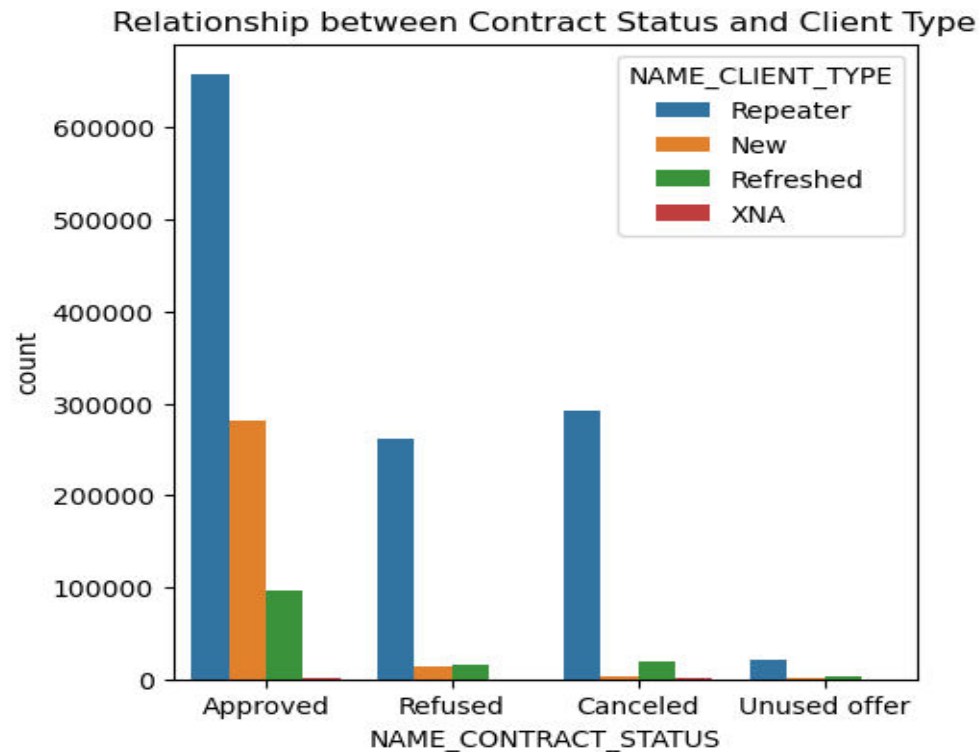- Most of the client amount credit between 0 to 100000.

- Highest number of loan approved is consumer loan and highest numer of loan is rejested and canceled is cash loan.



Relationship between Contract Status and Contract Type

- Most of the client type is Repeater. Repeater loans are most likely to approved.

Relationship between Contract Status and Client Type



- Considerable number of cash loans face rejection.
- Significant portion of cash loans are canceled.

Amount Credit Across NAME_PORTFOLIO

# Correlation of numerical variable of previous data

- All the columns AMT_ANNUITY, AMT_GOODS_PRICE, AMT_CREDIT, AMT_APPLICATION] have strong positive linear correlation with each other which is evident in below heatmap.

# Final Observation

- 54.99% of loan are approved in case of defaulters
- Education types secondary special have higher rate of loan approved
- Working professional have higher rate of loan approved
- Age group 30 -40 higher rate of loan approved
- Managers tend to request larger credit amounts but have a higher incidence of defaulting on loan payments.

# Final Conclusion

- females are applying more for loans in compared to men
- age group between 30 to 40 are applying more for loans.
- married peoples are applying for more loan
- Individuals who apply for multiple loans typically have a family size of two and do not have any children.

# From above all the conclusion ***

Probably Newly Married people age group between 30 to 40    are applying more for loan and having no child.

# THANK YOU

# DEEPANJANA ROY