# EDA Case Study-Final Report

Group Members:

ANKIT KUMAR SINGH

ARJUN SUCCENA

# Business Understanding

**Business Objective:**

► This case study aims that the company wants to understand the driving factors (or driver variables) behind loan default, i.e. the variables which are strong indicators of default. The company can utilize this knowledge for its portfolio and risk assessment.

**Goal of Analysis:**

► To find out the relation between the different variables and their impact on loan default. And suggest which attributes contributes a significant difference in Loan Default.

# Data to be analysed

► **'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.

► **'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.

► **'columns_description.csv'** is data dictionary which describes the meaning of the variables

# Data Understanding

The company has come across some important attributes in order to understand behavior of their approved loan customers w.r.t. loan default. Thus, the lending company has decided to work only on these variables to mitigate the future risk. The driver variables consider for this case study are:
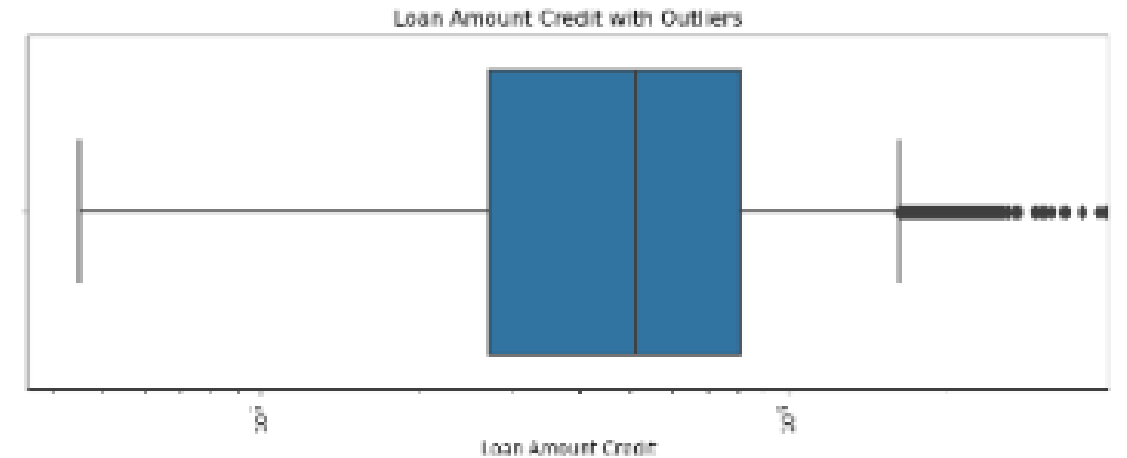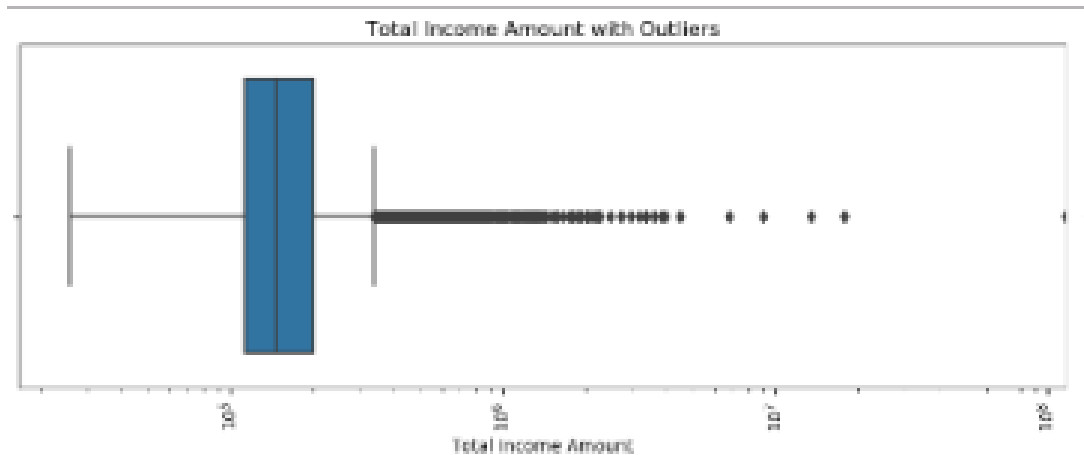
- NAME_CONTRACT_TYPE
- CODE_GENDER
- FLAG_OWN_CAR
- FLAG_OWN_REALTY
- CNT_CHILDREN
- AMT_INCOME_TOTAL
- AMT_CREDIT
- AMT_ANNUITY
- AMT_GOODS_PRICE
- NAME_TYPE_SUITE
- NAME_INCOME_TYPE
- NAME_EDUCATION_TYPE
- NAME_FAMILY_STATUS
- NAME_HOUSING_TYPE
- DAYS_BIRTH

# Data Cleaning and Preparation

- We loaded the file application.csv in Jupyter Notebook and found that there are 307511 observation and 122 variables, since in our data understanding stage we have already identified our driver variable we will be focusing our analysis on those variables.

- We found that SK_ID_CURR is the unique variable, so we have check for any duplicate in our data sets.

- We will look out for columns which have 10 -13% of missing values(Suggest metrics(no imputation) using Mean or Median based on the Outliers)

- We will look out for columns which have more that 50% of missing values(Drop the columns).

- We have removed the all columns with more than 50% of missing values and have made a subset of approx. 30 columns which will be useful for the analysis.

- We have then divided the whole dataset into two categories-Target 0 (Defaulters) and Target 1 (Non Defaulters).
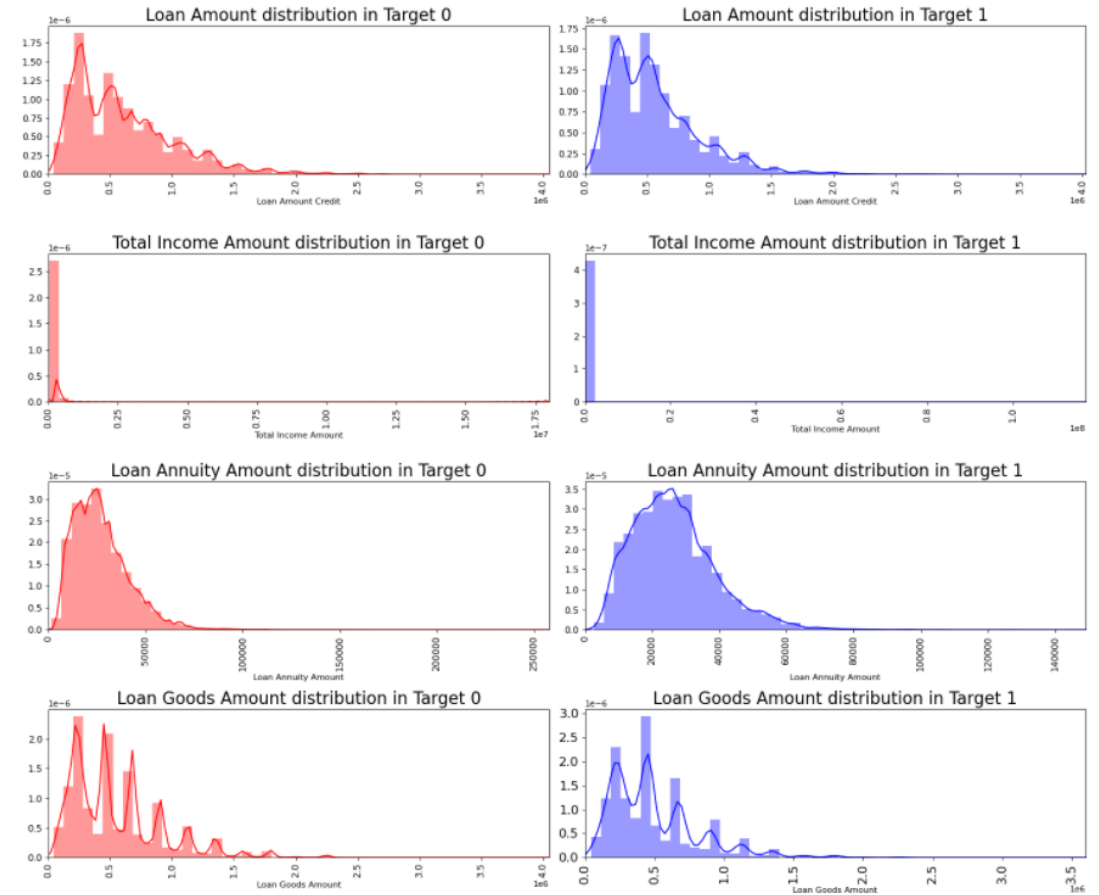
# Exploratory Data Analysis: Outlier Detection

The variables Total Income Amount, and loan amount credit have outliers present.

The plot shows the presence of outliers in the data set.

# Exploratory Data Analysis: Univariate Analysis(Numerical Variables)

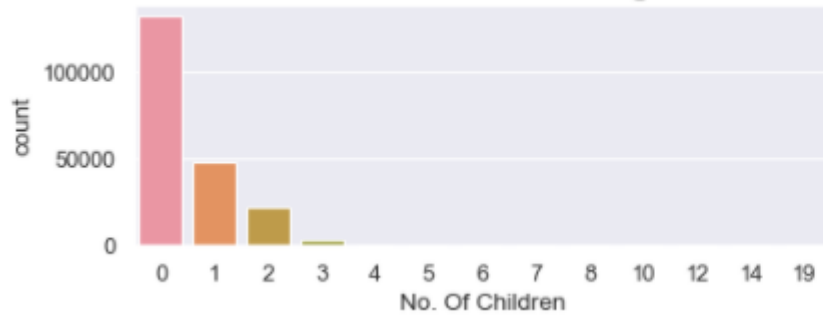We have plotted the Loan amount credit for target 0 and 1 as well as total income amount for target 0 and 1 both

As we can see in the plots, there is no significant difference in plot of loan amount and total income for both target 0 and target 1 data set.

# Exploratory Data Analysis: Univariate Analysis (Categorical Variables)

We have plotted bar graph of count of children and count of family members where we can see no significant observation which we can draw from it.
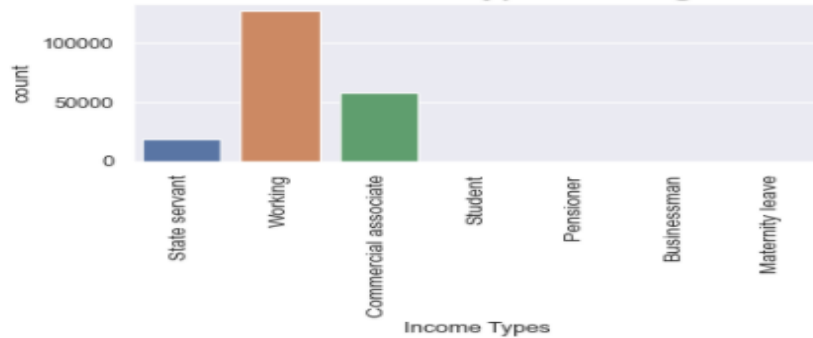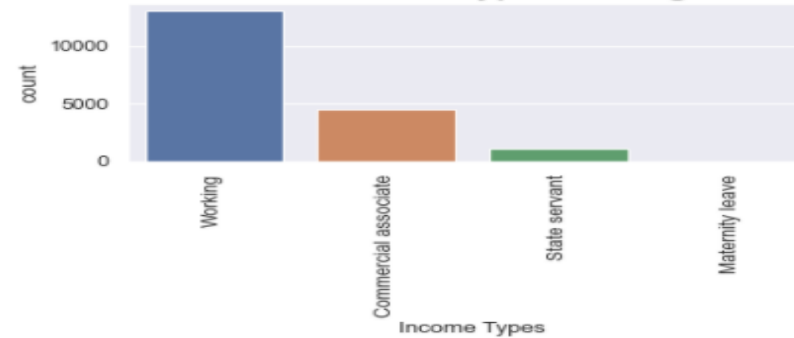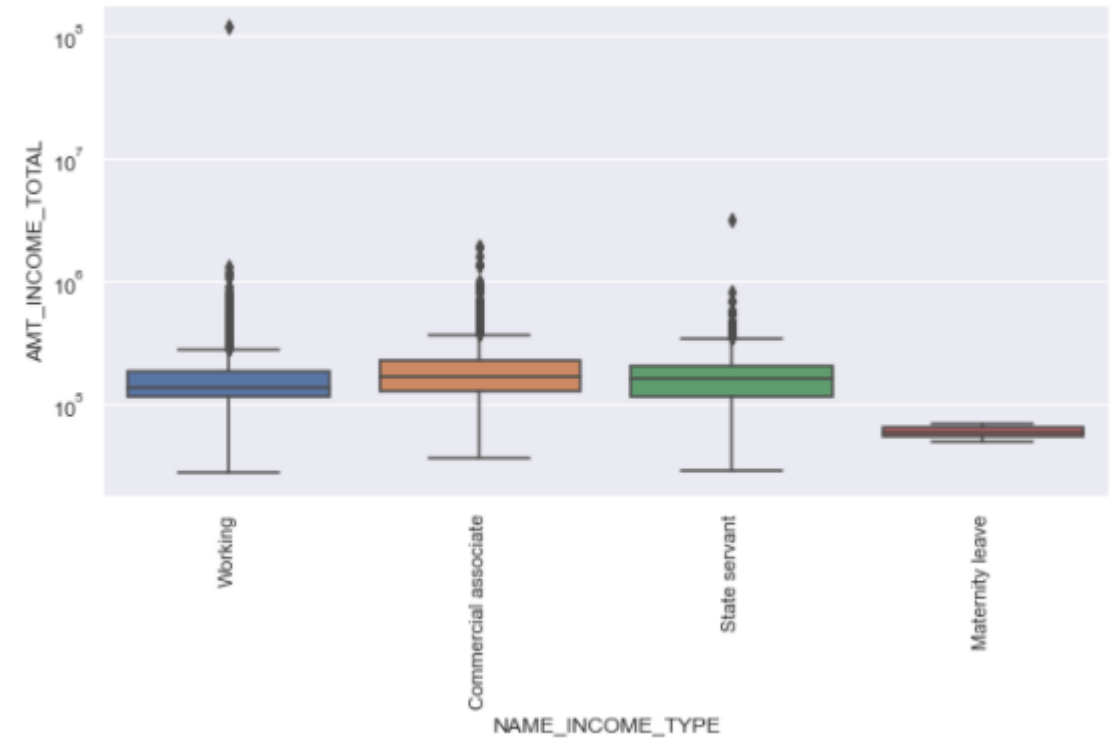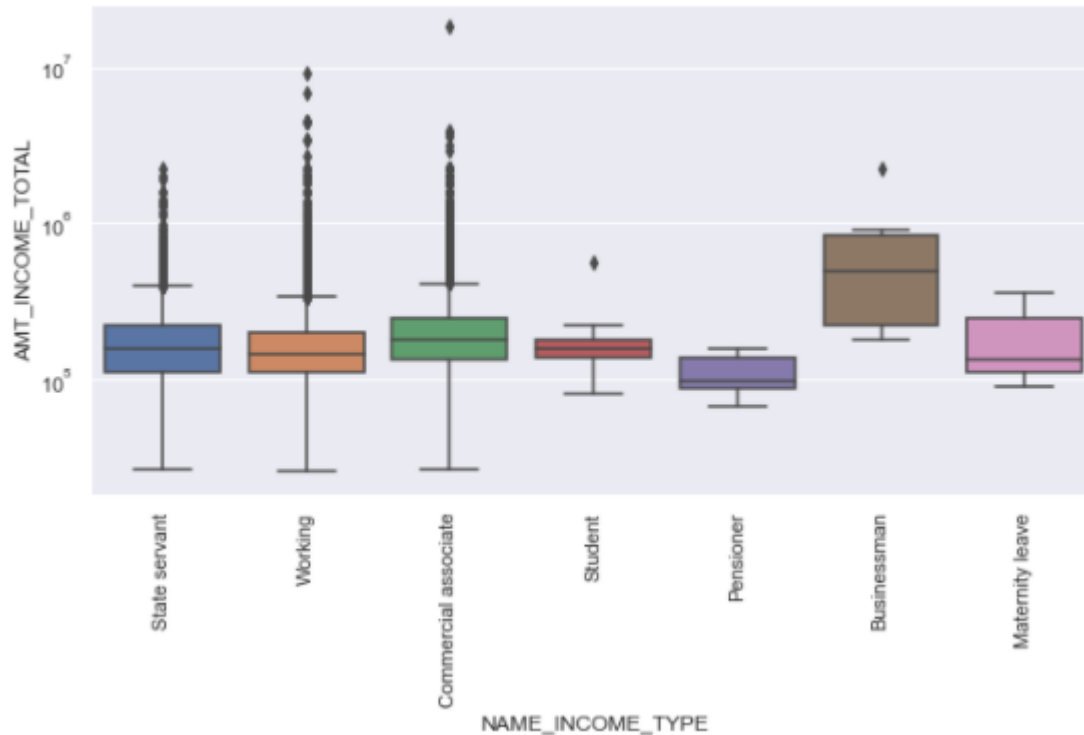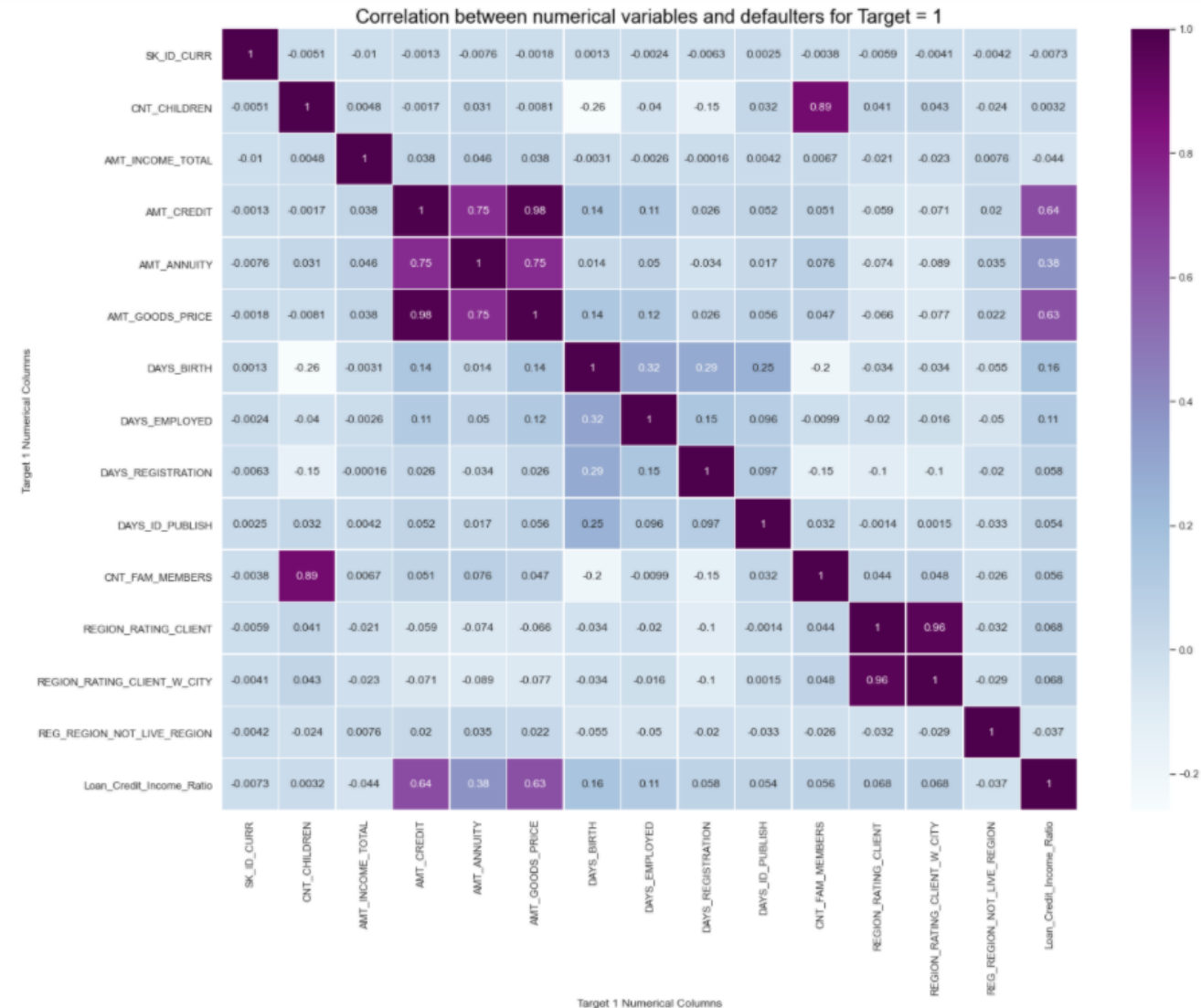
# Exploratory Data Analysis: Bivariate Analysis

Comparing the Defaulters list and the Non-Defaulters list, we can see that the Pensioners and Students are less likely to default.
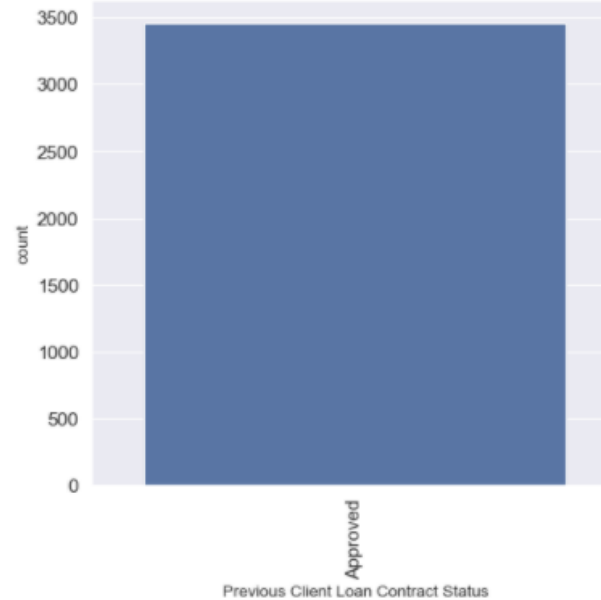
# Exploratory Data Analysis: Correlation

- ▶ The maroon boxes show the most correlated sub-sectors. Specifically, (0.98), (0.96), (089), (0.75) and (0.74) are the most correlated pairs.

- ▶ As inferred from the Correlation: Amount Credit and Amount of Goods are highly correlated.

- ▶ Also, Amount credit and Amount Annuity are highly correlated pairs as inferred from the naming convention



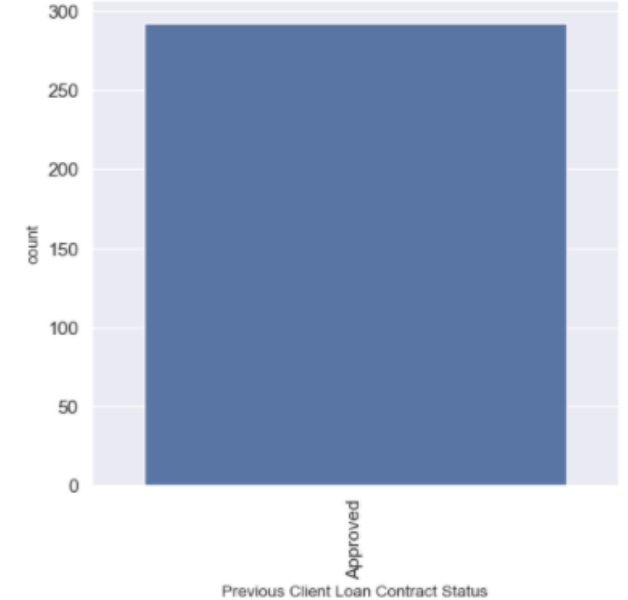Correlation between numerical variables and defaulters for Target = 1

# Exploratory Data Analysis: Univariate analysis on Previous Application Data

Based on the graph, we cannot see any significant impact on defaulters and non-defaulters list when we consider the previous variables for analysis.

# Conclusion

We can see that the Credit Amount and Amount of Goods are highly correlated.

We see a similar correlation between Credit Amount and Annuity Amount.

After analysis, we can conclude that pensioners and students are less likely to default on loan repayment.

We also see that the females are relatively higher in the defaulters list than males.

And we can also see that the defaulters are usually lying between 0-30 years of age.

Thank you