# CREDIT
# EDA - ASSIGNMENT

By:- Rashmi Singh

# Introduction

This assignment aims to give an idea of applying EDA in a real business scenario. In this assignment, apart from applying the EDA techniques we also develop a basic understanding of risk analytics in banking and financial services and understand how data is used to minimize the risk of losing money while lending to customers.

The loan providing companies find it hard to give loans to the people due to their insufficient or non-existent credit history. Because of that, some consumers use it to their advantage by becoming a defaulter. Suppose you work for a consumer finance company which specializes in lending various types of loans to urban customers. You have to use EDA to analyze the patterns present in the data. This will ensure that the applicants capable of repaying the loan are not rejected

# Business Objective

This case study aims to identify patterns which indicate if a client has difficulty paying their instalments 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, etc. 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. In other words, 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

# Data Cleaning Approach

➢ Due to Threshold theory for any dataset if we have more than 40% missing value we can drop the column because that column don't give accurate result.

➢ I have useless columns like flag and in this values are 0 or 1 and  which is not giving any useful information so I have removed all the unnecessary columns

➢ **checked days column where negative values are there so we need to convert it to positive by using absolute**

➢ **Also by analysis the data which I get is Female applicants are more as compare to male applicants**
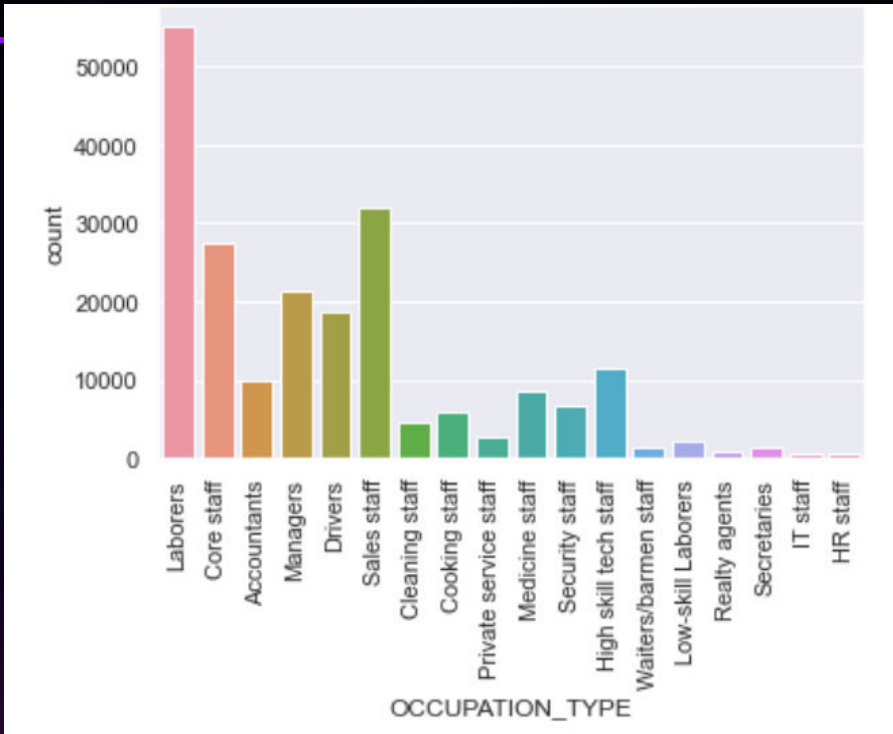
# Missing Value

## Application Data

| | |
|---|---|
| COMMONAREA_MEDI | 69.872297 |
| COMMONAREA_AVG | 69.872297 |
| COMMONAREA_MODE | 69.872297 |
| NONLIVINGAPARTMENTS_MODE | 69.432963 |
| NONLIVINGAPARTMENTS_AVG | 69.432963 |
| NONLIVINGAPARTMENTS_MEDI | 69.432963 |
| FONDKAPREMONT_MODE | 68.386172 |
| LIVINGAPARTMENTS_MODE | 68.354953 |
| LIVINGAPARTMENTS_AVG | 68.354953 |
| LIVINGAPARTMENTS_MEDI | 68.354953 |
| FLOORSMIN_AVG | 67.848630 |
| FLOORSMIN_MODE | 67.848630 |
| FLOORSMIN_MEDI | 67.848630 |
| YEARS_BUILD_MEDI | 66.497784 |
| YEARS_BUILD_MODE | 66.497784 |
| YEARS_BUILD_AVG | 66.497784 |
| OWN_CAR_AGE | 65.990810 |
| LANDAREA_MEDI | 59.376738 |
| LANDAREA_MODE | 59.376738 |

| | |
|---|---|
| NONLIVINGAREA_MODE | 55.179164 |
| NONLIVINGAREA_AVG | 55.179164 |
| NONLIVINGAREA_MEDI | 55.179164 |
| ELEVATORS_MEDI | 53.295980 |
| ELEVATORS_AVG | 53.295980 |
| ELEVATORS_MODE | 53.295980 |
| WALLSMATERIAL_MODE | 50.840783 |
| APARTMENTS_MEDI | 50.749729 |
| APARTMENTS_AVG | 50.749729 |
| APARTMENTS_MODE | 50.749729 |
| ENTRANCES_MEDI | 50.348768 |
| ENTRANCES_AVG | 50.348768 |
| ENTRANCES_MODE | 50.348768 |
| LIVINGAREA_AVG | 50.193326 |
| LIVINGAREA_MODE | 50.193326 |
| LIVINGAREA_MEDI | 50.193326 |
| HOUSETYPE_MODE | 50.176091 |
| FLOORSMAX_MODE | 49.760822 |
| FLOORSMAX_MEDI | 49.760822 |

## Previous Data

| | |
|---|---|
| NAME_TYPE_SUITE | 49.119754 |
| NFLAG_INSURED_ON_APPROVAL | 40.298129 |
| DAYS_TERMINATION | 40.298129 |
| DAYS_LAST_DUE | 40.298129 |
| DAYS_LAST_DUE_1ST_VERSION | 40.298129 |
| DAYS_FIRST_DUE | 40.298129 |
| DAYS_FIRST_DRAWING | 40.298129 |
| AMT_GOODS_PRICE | 23.081773 |
| AMT_ANNUITY | 22.286665 |
| CNT_PAYMENT | 22.286366 |
| PRODUCT_COMBINATION | 0.020716 |
| AMT_CREDIT | 0.000060 |
| CHANNEL_TYPE | 0.000000 |
| NAME_YIELD_GROUP | 0.000000 |
| NAME_SELLER_INDUSTRY | 0.000000 |
| SELLERPLACE_AREA | 0.000000 |
| SK_ID_PREV | 0.000000 |
| NAME_PRODUCT_TYPE | 0.000000 |
| NAME_PORTFOLIO | 0.000000 |
| SK_ID_CURR | 0.000000 |
| NAME_CLIENT_TYPE | 0.000000 |
| CODE_REJECT_REASON | 0.000000 |
| NAME_PAYMENT_TYPE | 0.000000 |
| DAYS_DECISION | 0.000000 |
| NAME_CONTRACT_STATUS | 0.000000 |
| NAME_CASH_LOAN_PURPOSE | 0.000000 |
| AMT_APPLICATION | 0.000000 |
| NAME_CONTRACT_TYPE | 0.000000 |

# MISSING VALUE TREATMENT

1.if variable is **object** means categorical so will fill with **Mode**

2.if variable is **int** or **float** means numerical so will fill with **Median** or **Mean**

**If we have numerical variable so we preferred fill with median because due to outlier in data set they impact mean but not median**
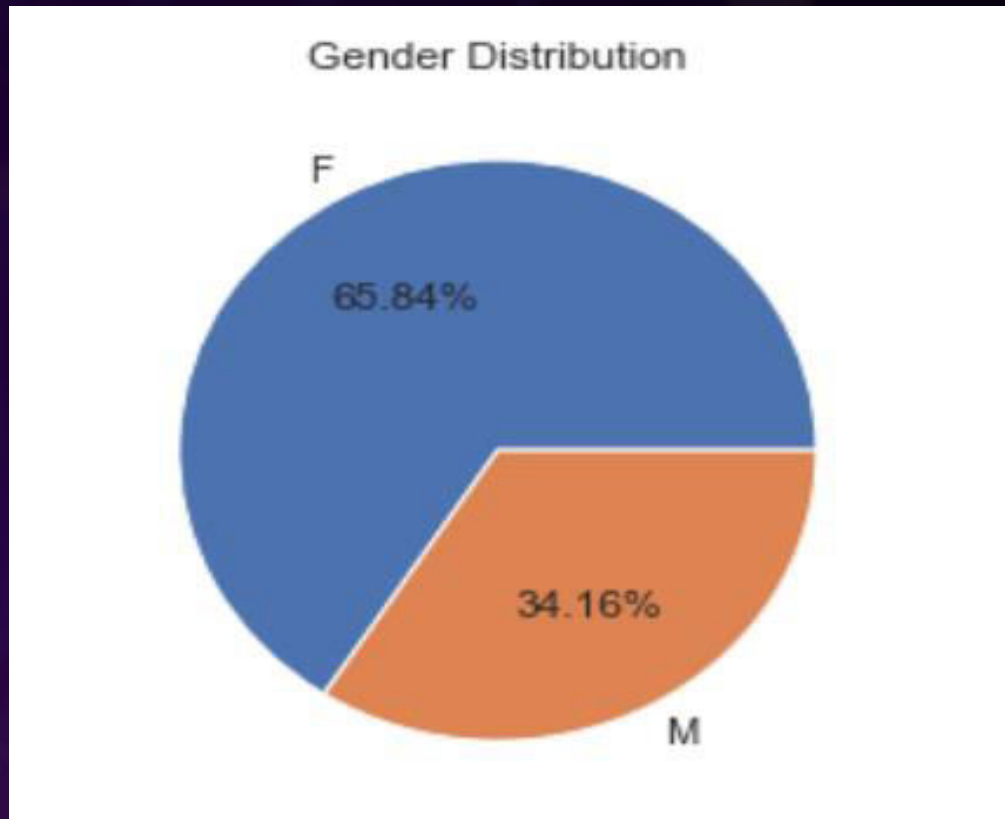
# HANDLING OUTLIER



The occupation Type is categorical so we need to fill it with mode
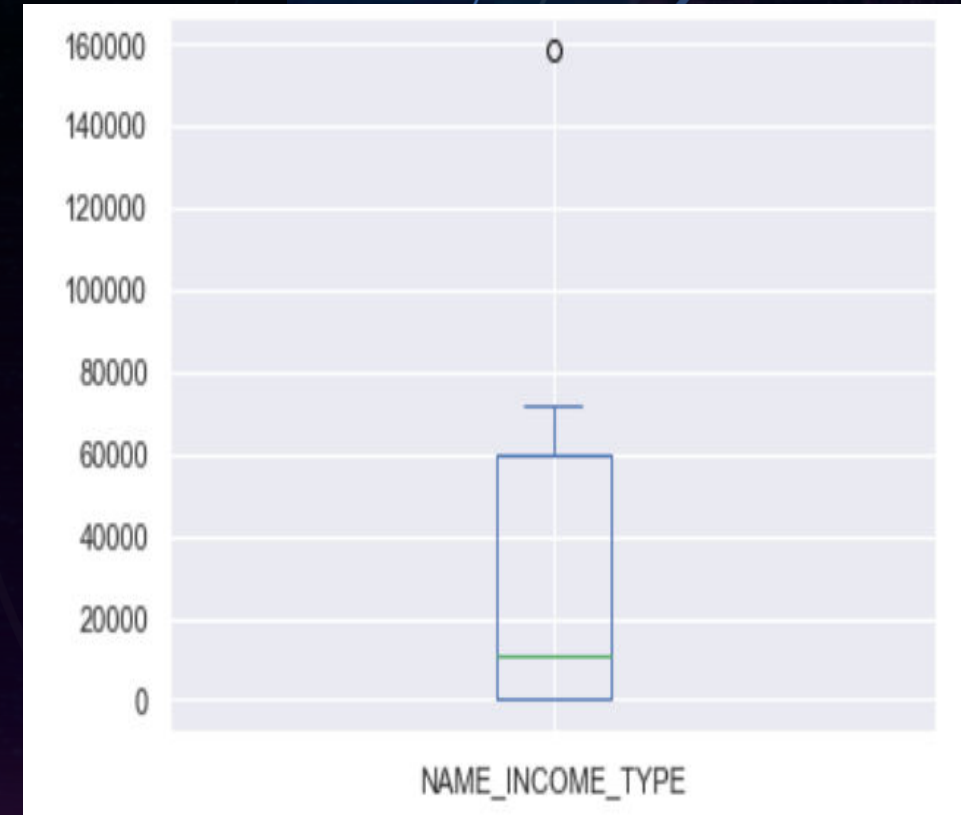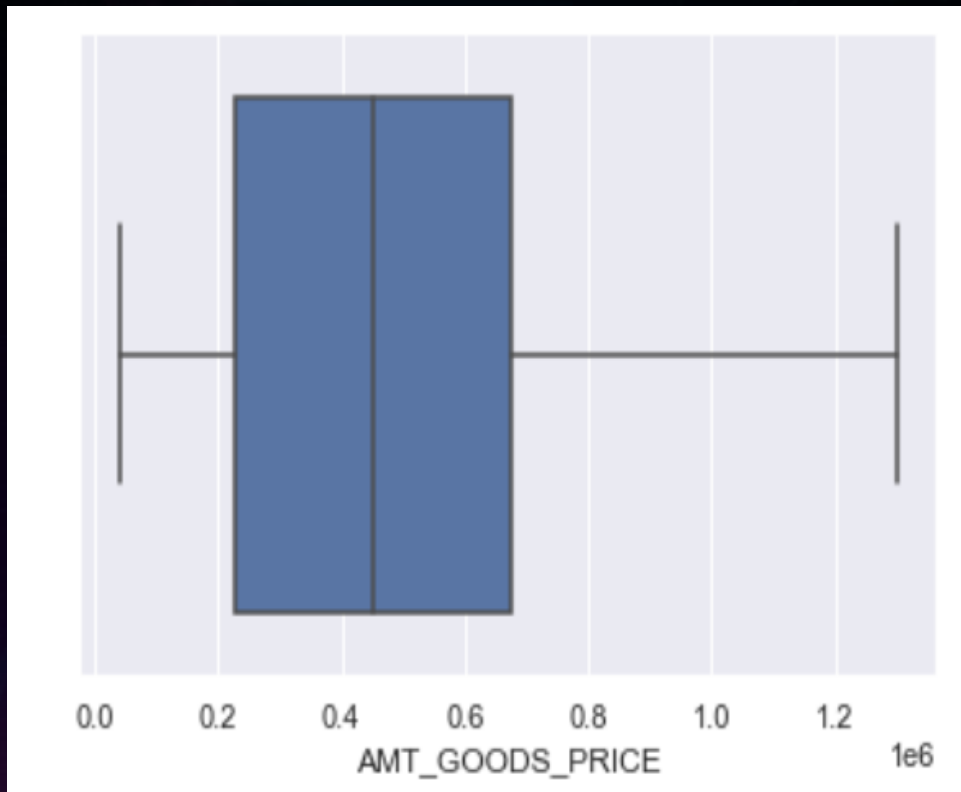
# Gender Distribution



Gender Distribution

F

65.84%

34.16%

M

Female applicants are more as compare to male applicants

# OUTERLIERS

Outliers can be removed after setting the value upto 0.95 percentile

- Outliers identified which is at the max point

# Previous Application Data

## Univariate Analysis

## Bivariate Analysis

The CASH LOAN ARE MORE AS COMPARE TO OTHER TYPE

OF LOANS



Repeaters have the highest AMT_GOODS_PRICE Cash Loans are also more
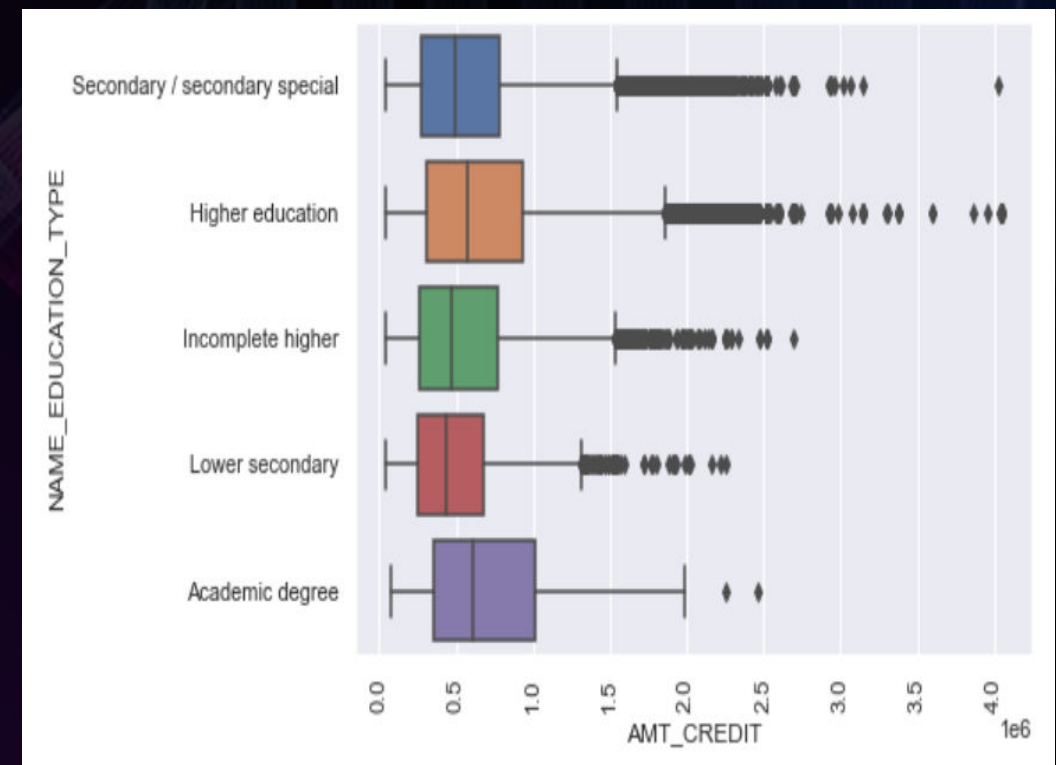as compare to other loans¶



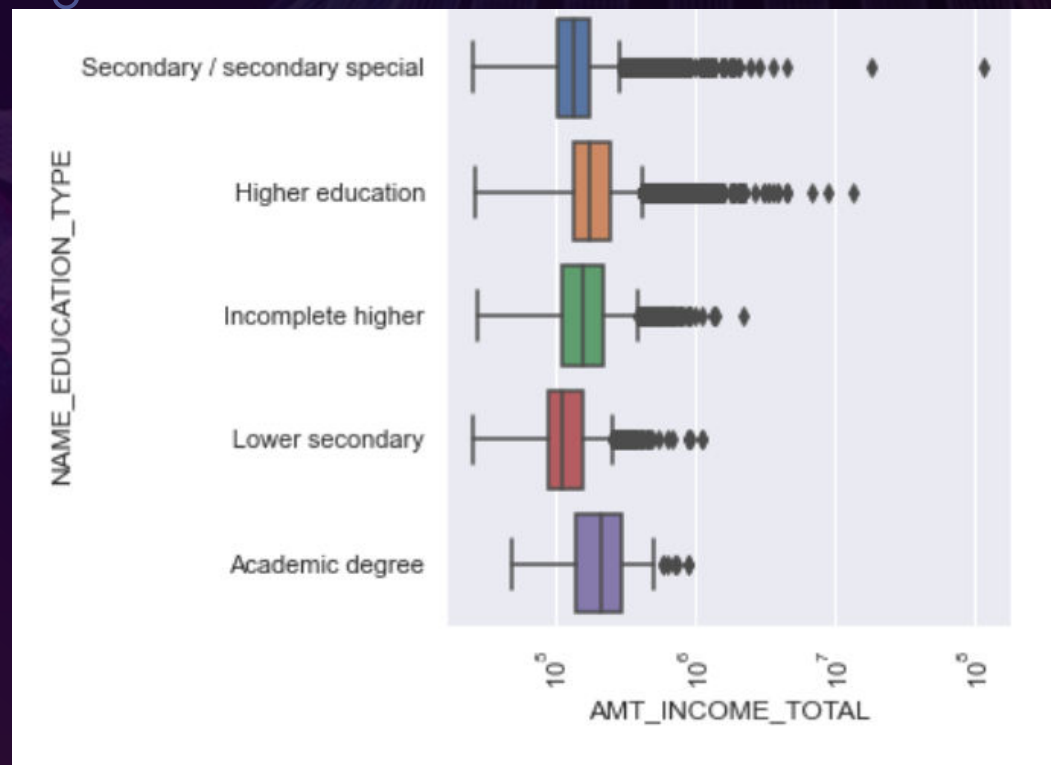Repeater clients are more as compare to others



THE NUMBER OF REFUSED LOANS ARE for LOANS AND REPAIR PURPOSE
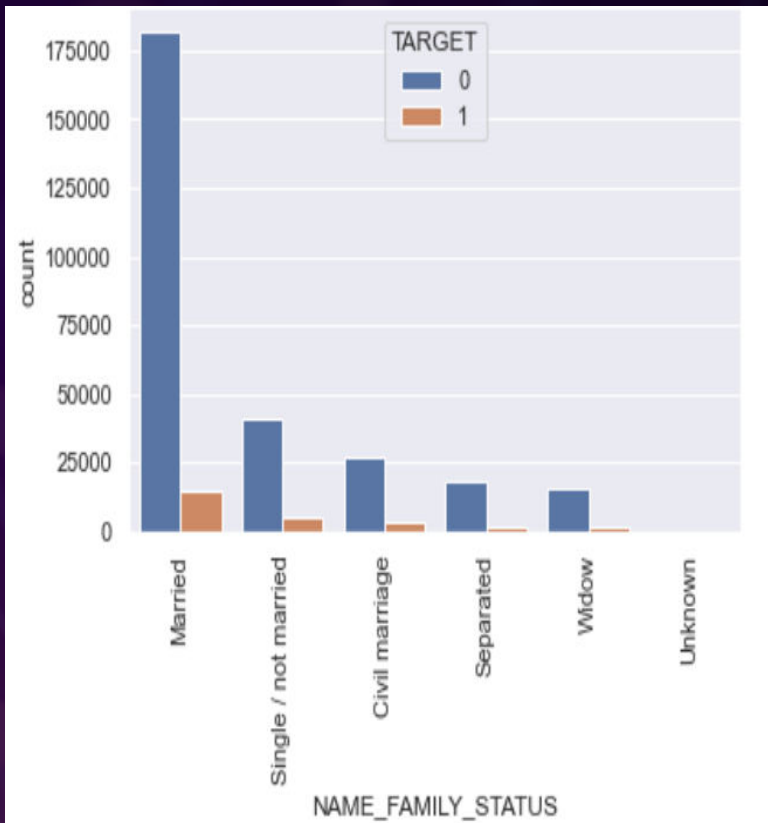UNKNOWN AND UNKNOWN 1 ARE ALSO VERY POPULAR

# Amount Income Total/Amount Credit

Clearly the business Income type is domaining the other business types with working which as a verity of range
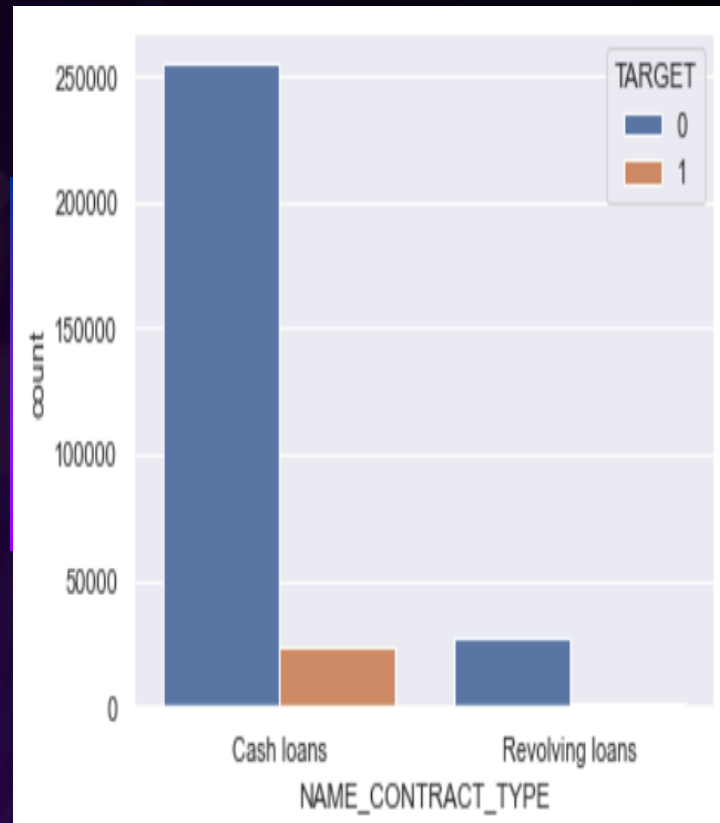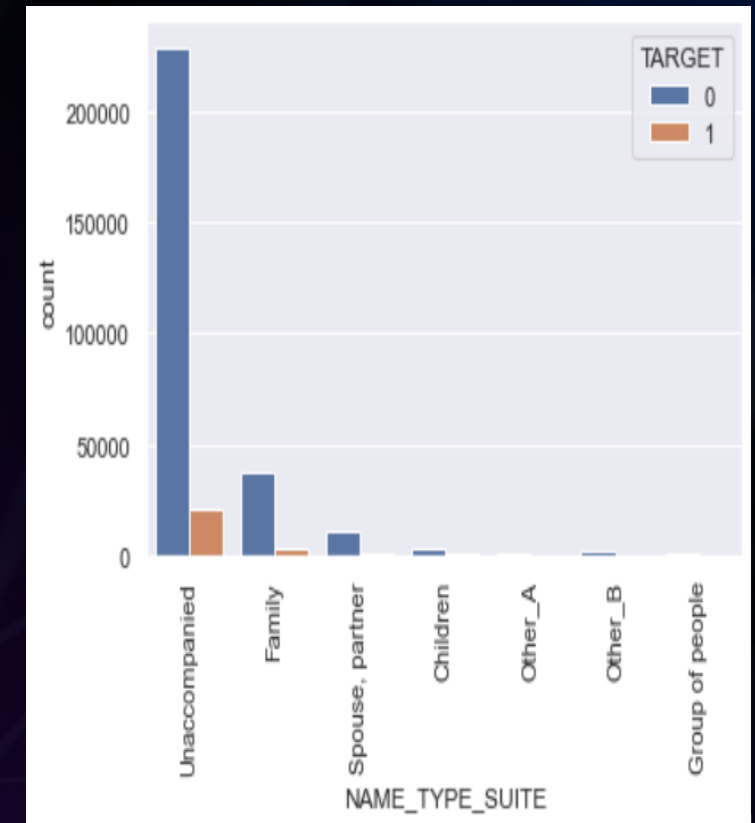
# BIVARIATE ANALYSIS

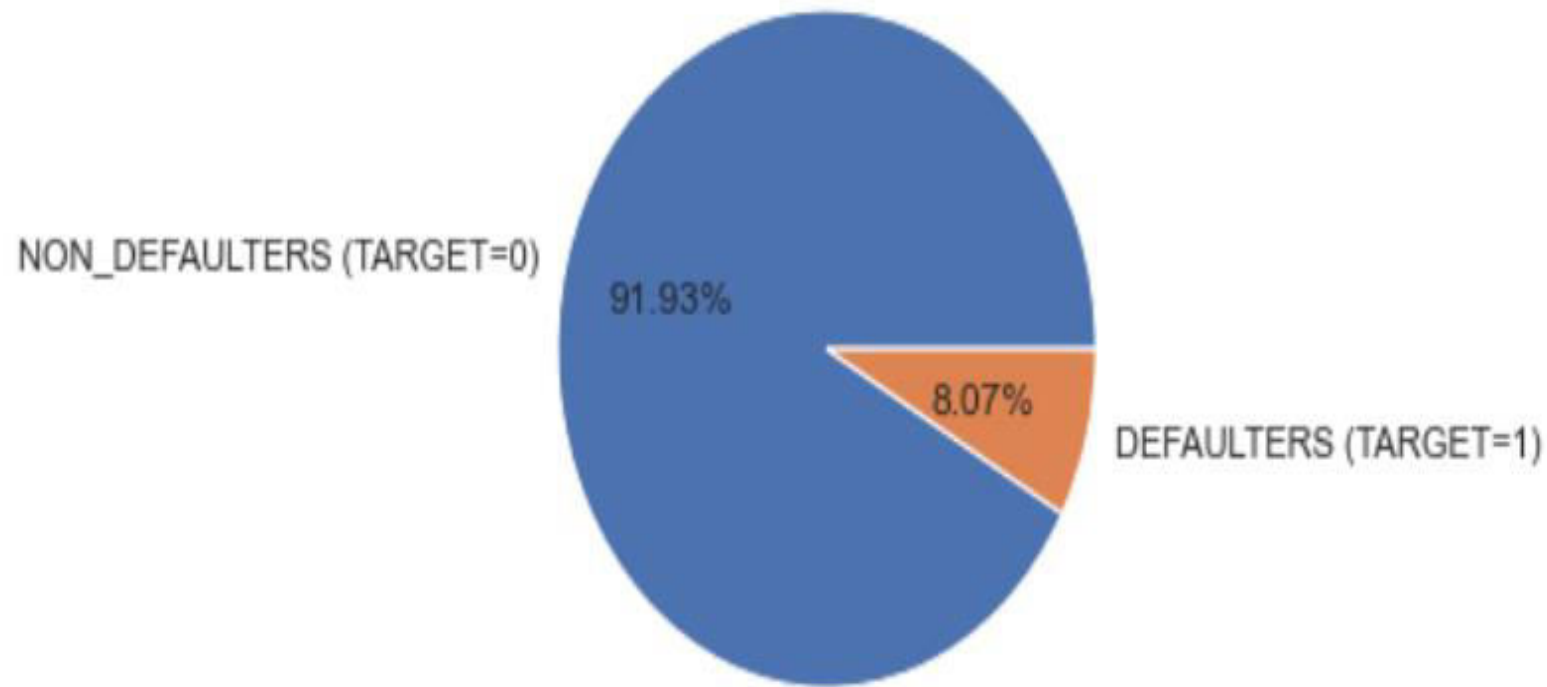Married and single are top 2 category to target which has highest no. of non defaulters

Cash loans have more non defaulters than revolving loans

Unaccompanied and family are the people who have less defaulters than other category
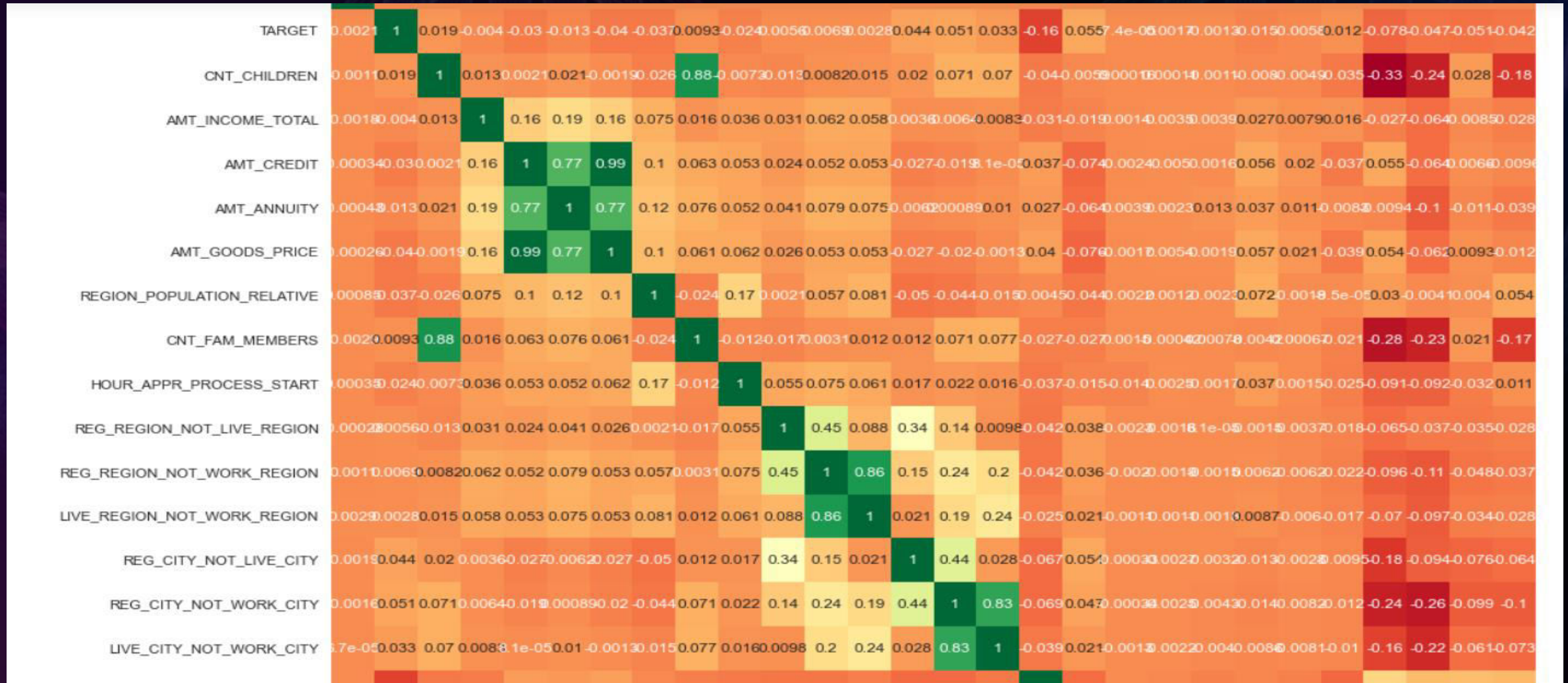






13

# SEGMENTED VARIABLE



We can see there is a huge difference in the data
The Defaulter rate is less 8.07 % and Non-Defaulter rate is high 91.93%

# MULTIVARIATE ANALYSIS

CHECKING For Co-relation ALL THE NUMERIC VARIABLE AT ONES

# FINAL OBSERVATION

i.    Target variable for Application dataset - "TARGET"

ii.    Target variable for Previous dataset - "NAME_CONTRACT_STATUS"

iii.   The rate of defaulters are less in the range of 20-40 & 40-60 are good target audience.

iv.   Laborers , Core and Sales Staff is the occupation type that has the loan approved and has the highest non defaulter rate.

v.    Married people are more likely to get loan approved in comparison to any other Marital Status of the people so this is also a good target audience .

vi. Secondary Education has the Highest Approval rate ,although the Income of Academic degree holder are more as

Compare to Secondary education still the approval rate is more than Academic Degree holders.