# EDA – CREDIT ASSIGNMENT

NAME : PRATEEK JAIN

BATCH ID: 1975

**Data Science Program - March 2022**

# INTRODUCTION

| Given dataset of loan providing company | 3 .csv files as dataset | 1. Application<br><br>2. Previous Application<br><br>3. Column description |

# About Dataset

**Application**
- Contain information at the time of applying
  - **The client with payment difficulties:**
  - **All other cases:**

**Previous Application**
- Contain information about four types of decisions that could be taken
  - **Approved**
  - **Cancelled**
  - **Refused**
  - **Unused offer**

**Column Description**
- Contains information
  - Description/meaning of the columns so as to get better understanding of dataset

# Business Understanding

The loan providing companies find it hard to give loans to the people due to their insufficient or non-existent credit history

When the company receives a loan application, the company has to decide for loan approval based on the applicant's profile

# Objective

## Understanding

- How consumer attributes and loan attributes influence the tendency of default.

## Identifying

- Patterns indicating if a client has difficulty paying their installments
- Top 10 correlations

## Predicting

- Consumers capable of repaying the loan are not rejected.
- Understand the driving factors (or driver variables) behind loan default

# APPLICATION DATAFRAME

1. • Importing Libraries
2. • Reading the data set and finding percentage of null values
3. • Dropping columns with missing values >45%
4. • Identifying continuous and categorical columns/variable
5. • Continuous column - Columns containing unique values > 58
6. • Categorical column – Columns containing unique values < 58
7. • Imputation for missing value<45 (categorical – mode, continuous – median)
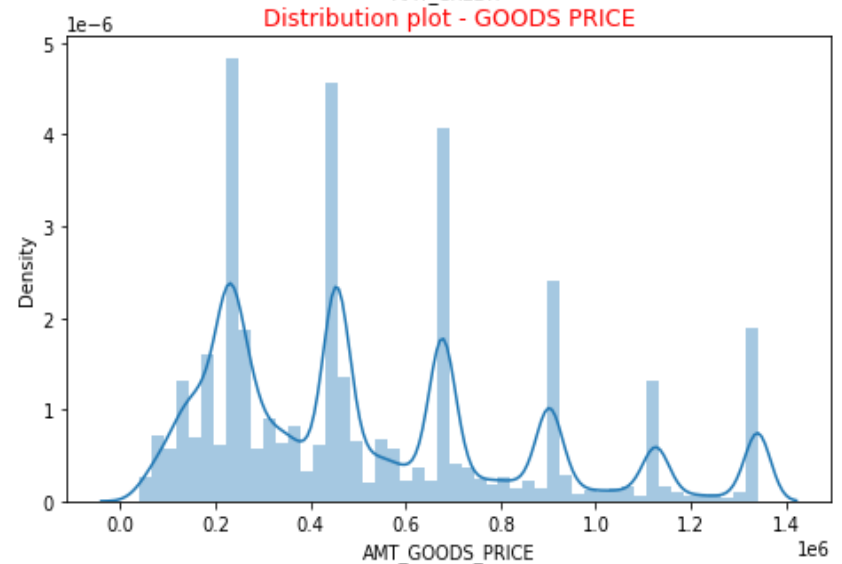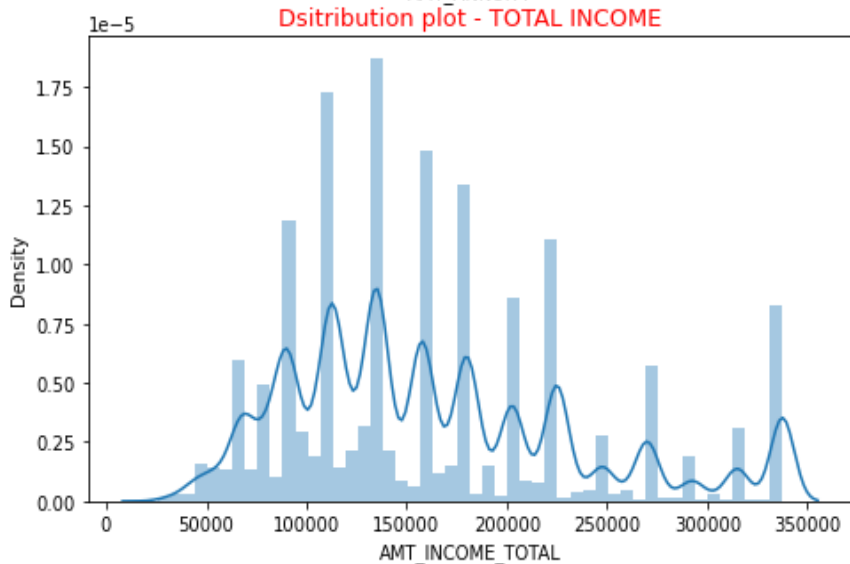8. • Dropping unnecessary columns
9. • Detecting Outliers – Using Subplots
10. • Handling outliers by flooring and capping

HEATMAP

# UNIVARIATE ANALYSIS

# SUBPLOTS



**Most of the people taking loans have annuity between 20000 and 30000**

# DEFAULTER AND NON DEFAULTERS

## APPLICATION LOAN DATASET – TWO MAIN VARIABLES
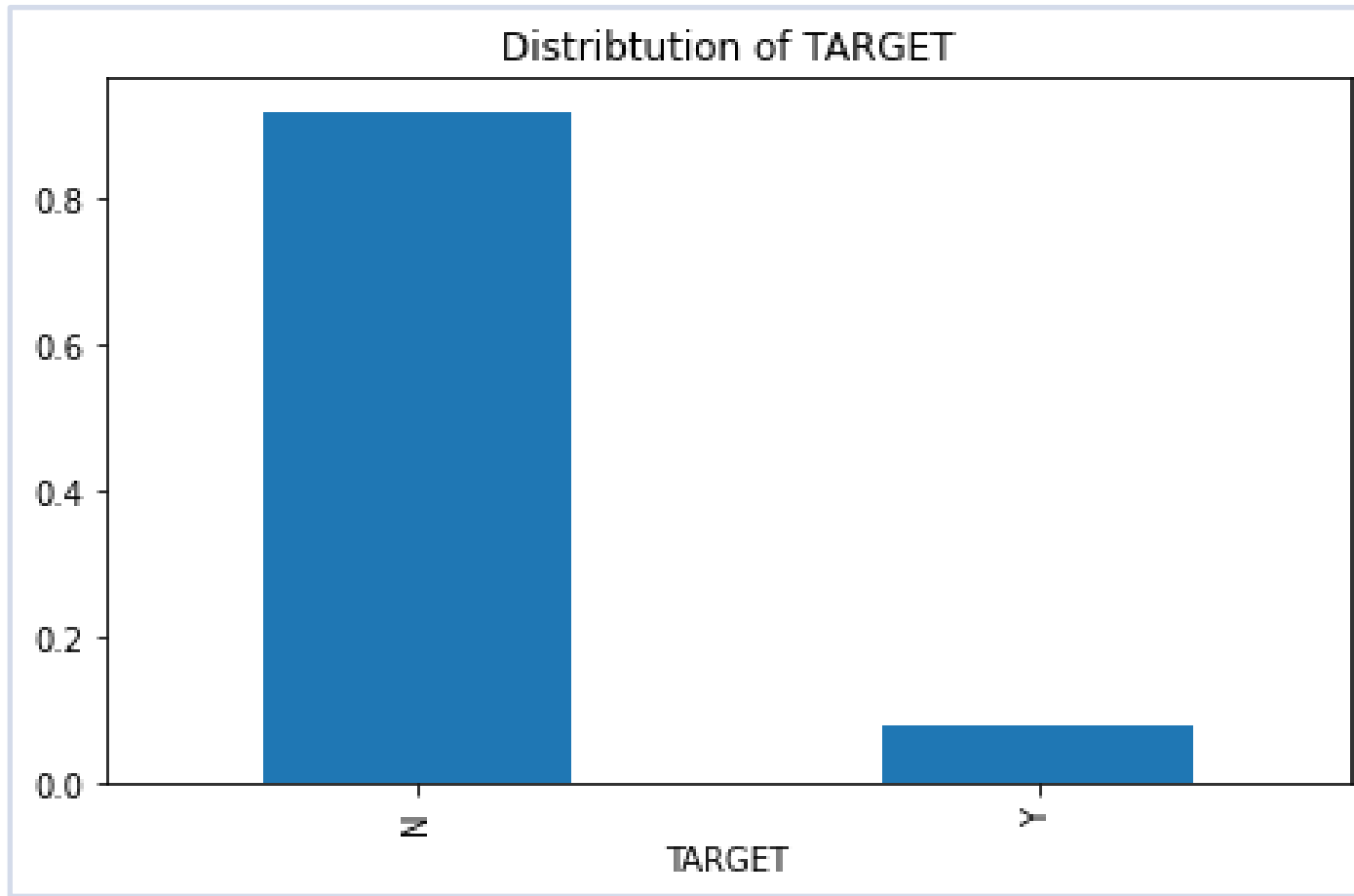
- DEFAULTER
- NON DEFAULTER

## DEFAULTERS

- The client with payment difficulties
- Response is stored as 1 ('Y'- Yes)

## NON DEFAULTERS

- **All other cases**
- Response is stored as 0 ('N'- No)

## BARPLOT – DEFAULTER AND NON DEFAULTERS



Distribtution of TARGET

- **Large number of people applying for loans are non defaulters**

# COUNTPLOT – LOAN TYPE



**Number of people whether defaulters or non defaulters prefer to take 'CASH LOAN' as compared to 'REVOLVING LOAN'**
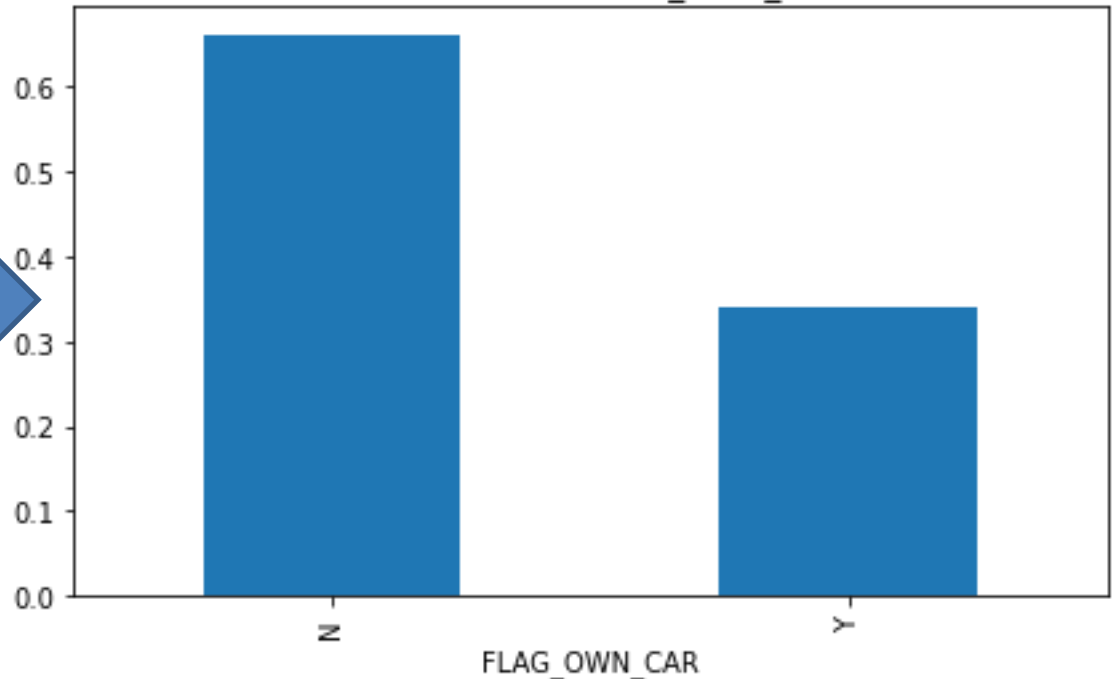
# COUNTPLOT - GENDER



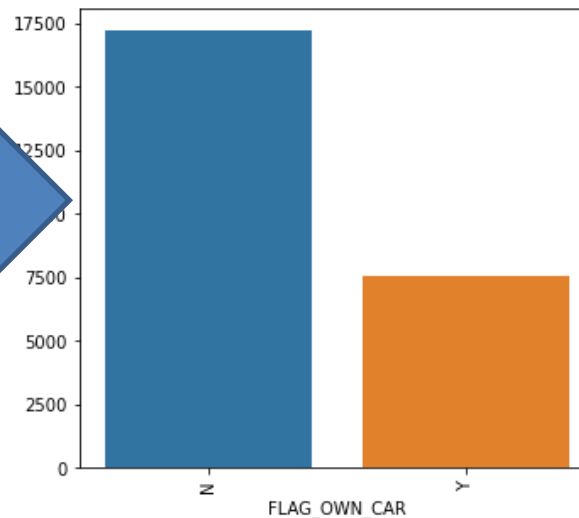Number of "FEMALE" taking loans is much higher than the number of "MALE" for both the target variables

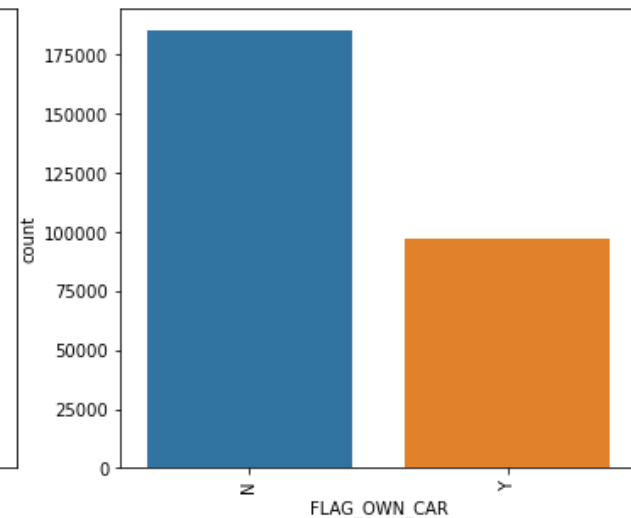# COUNTPLOT - CAR



Distribtution of FLAG_OWN_CAR

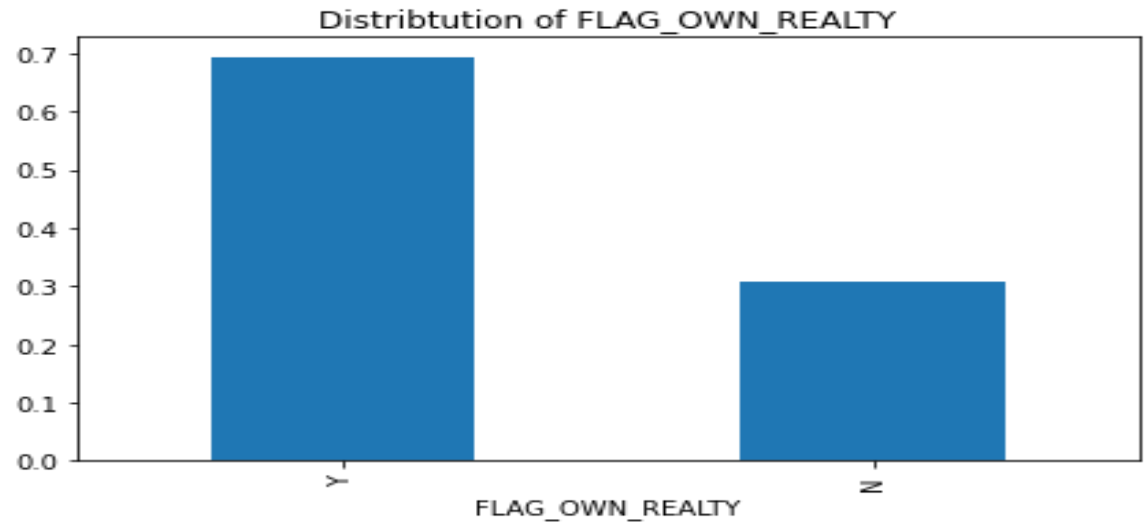People who "DON'T OWN CAR" are applying more for loan

COUNTPLOT - DEFAULTERS

COUNTPLOT - NON DEFAULTERS

Most people applying for the loan and have car are "NON DEFUALTER"

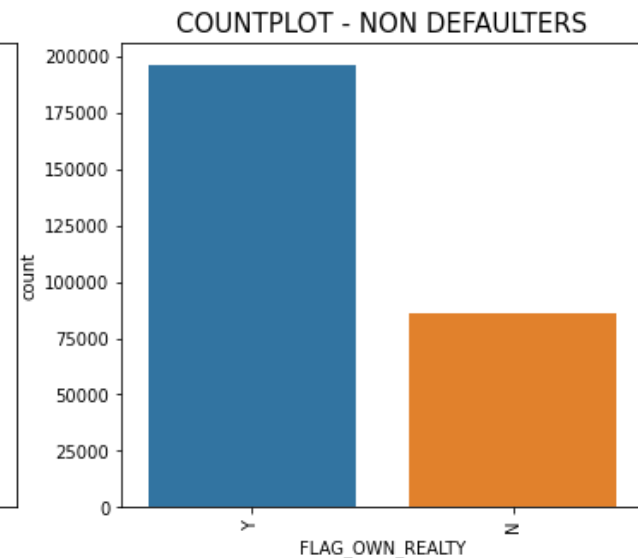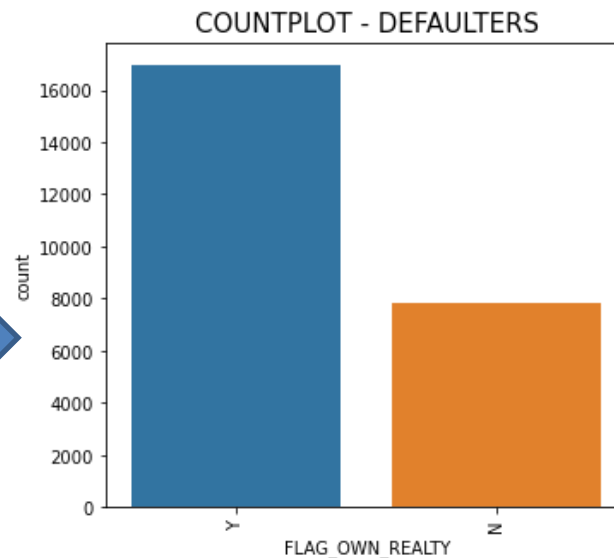# COUNTPLOT - HOUSE



Most people applying for the loan have own house

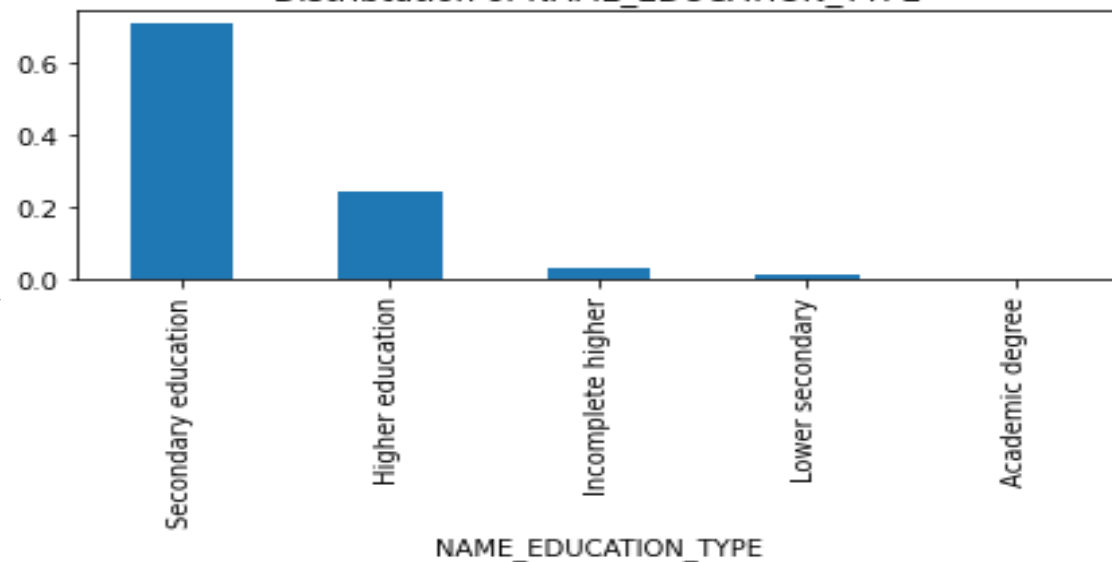People having 'OWN HOUSE' and non defaulters are more

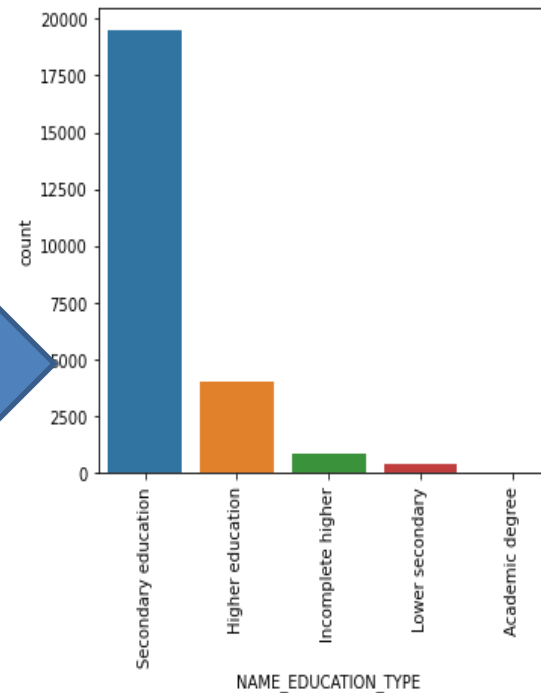Distributution of NAME_EDUCATION_TYPE

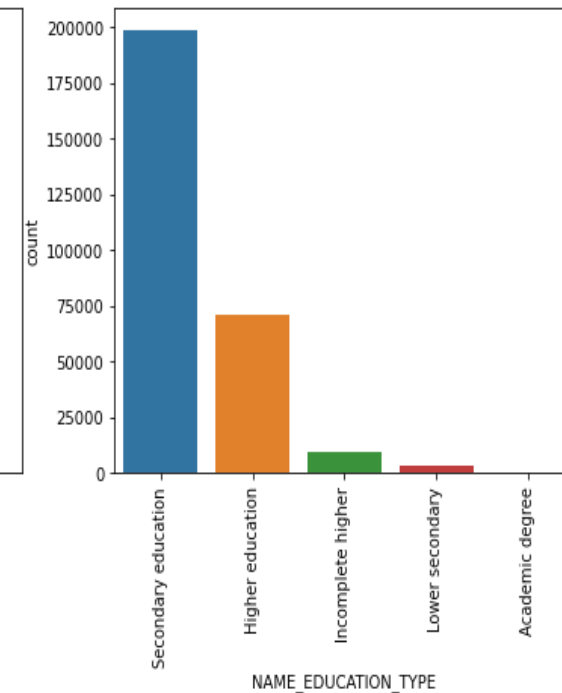People applying for the loan, mostly have secondary education

•People having secondary & higher education have less payment difficulties

•People have academic degrees can also be targeted as they apply less for the loans
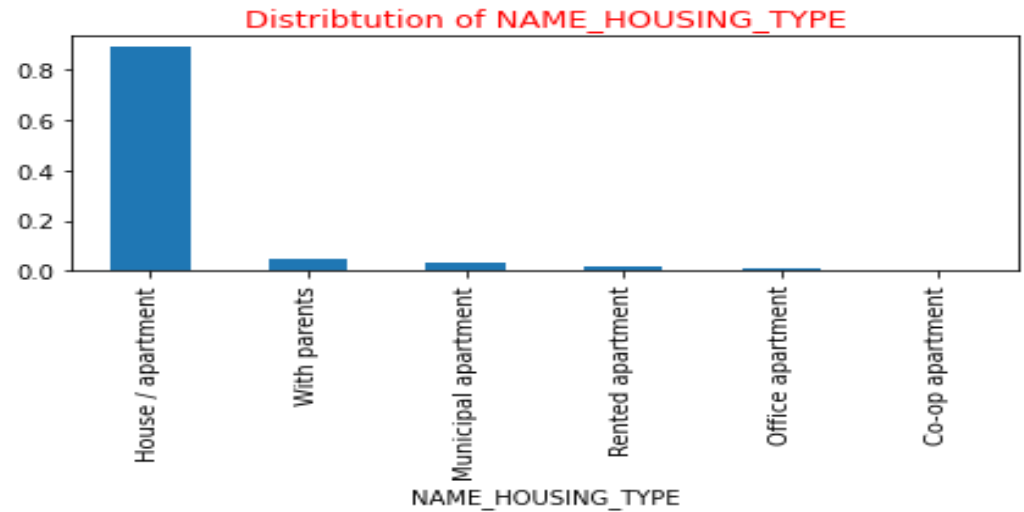
Distribution for DEFAULTERS
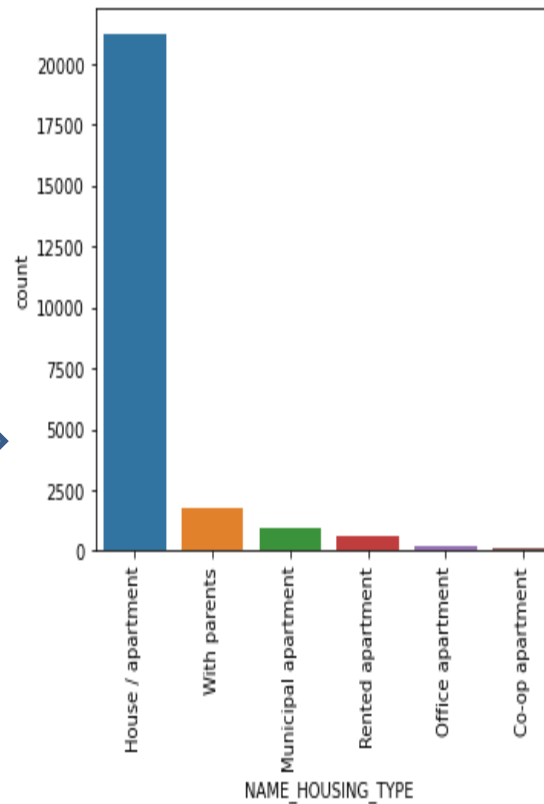
Distribution for NON DEFAULTERS

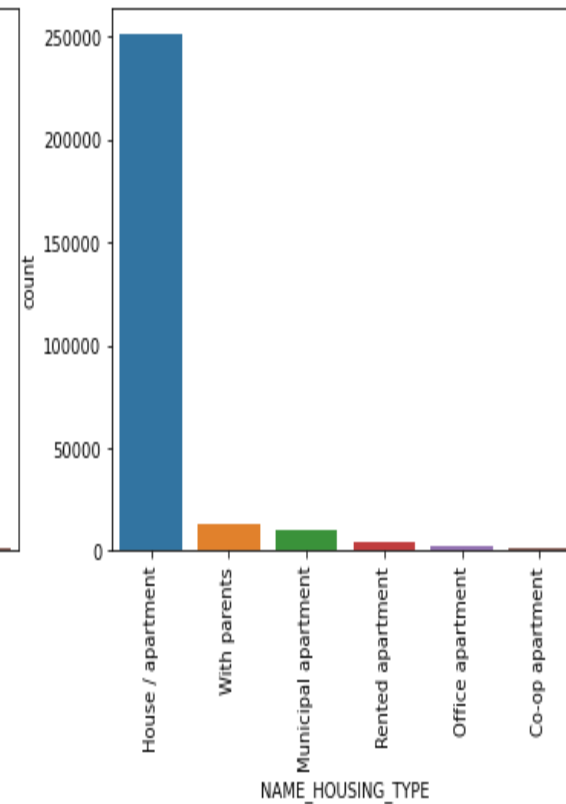People applying for the loan are mostly living in a house/apartment

Distribtution of NAME_HOUSING_TYPE

•Large number of people applying for loans have there own 'HOUSE/APARTMENT'

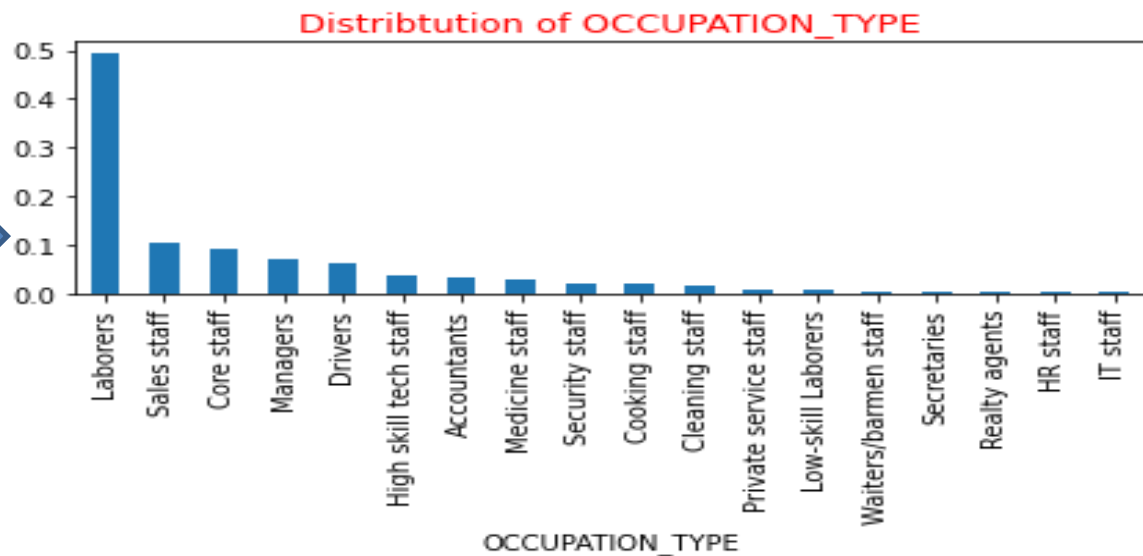•People living 'WITH PARENTS' are more likely to have payment difficulties.
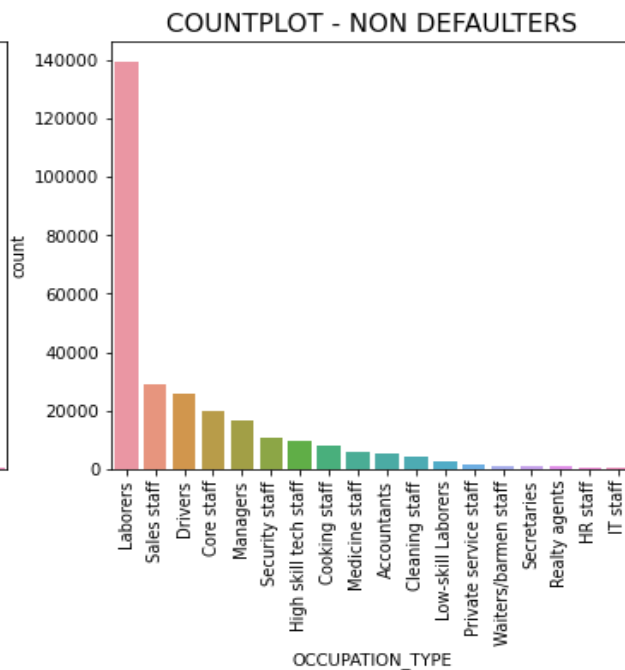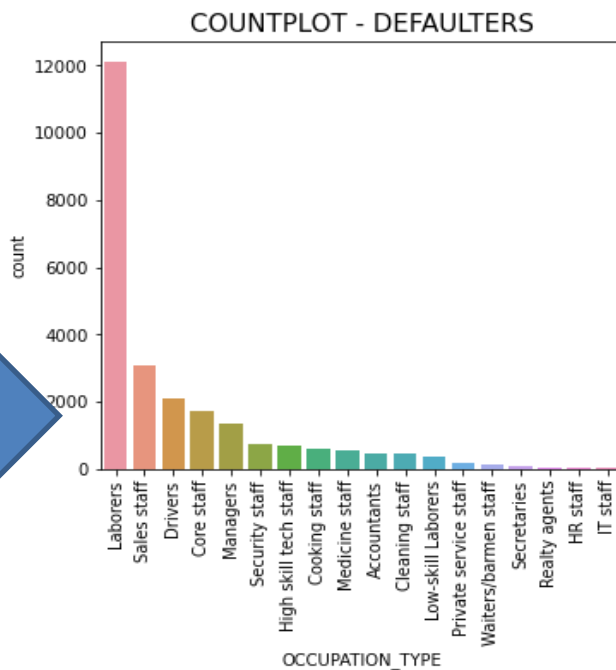
COUNTPLOT - DEFAULTERS

COUNTPLOT - NON DEFAULTERS

Labourers are applying more for loan as compared to other category

Distribtution of OCCUPATION_TYPE
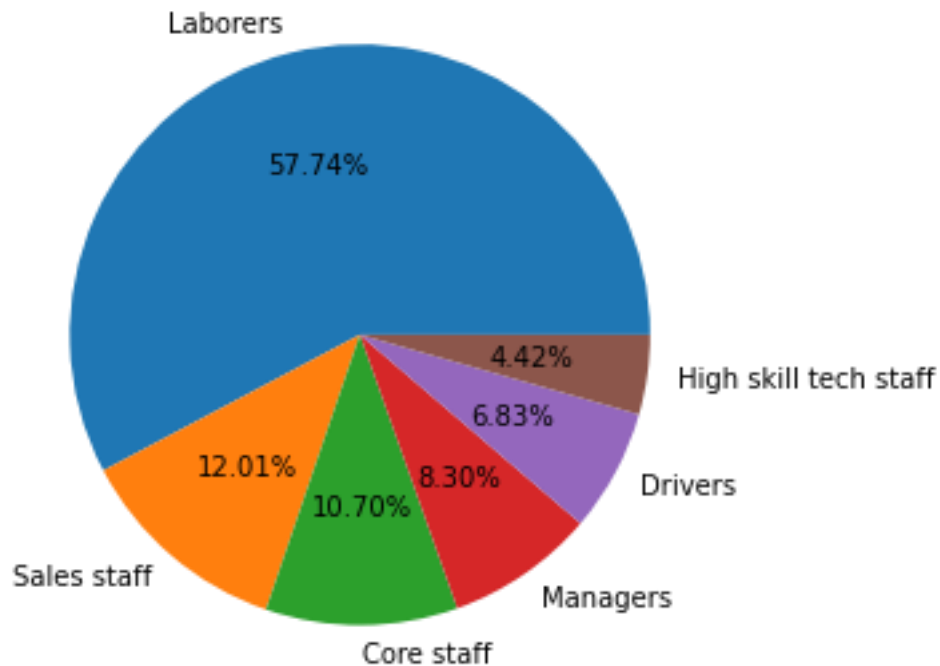
Labourers are the maximum number of defaulters and non defaulters

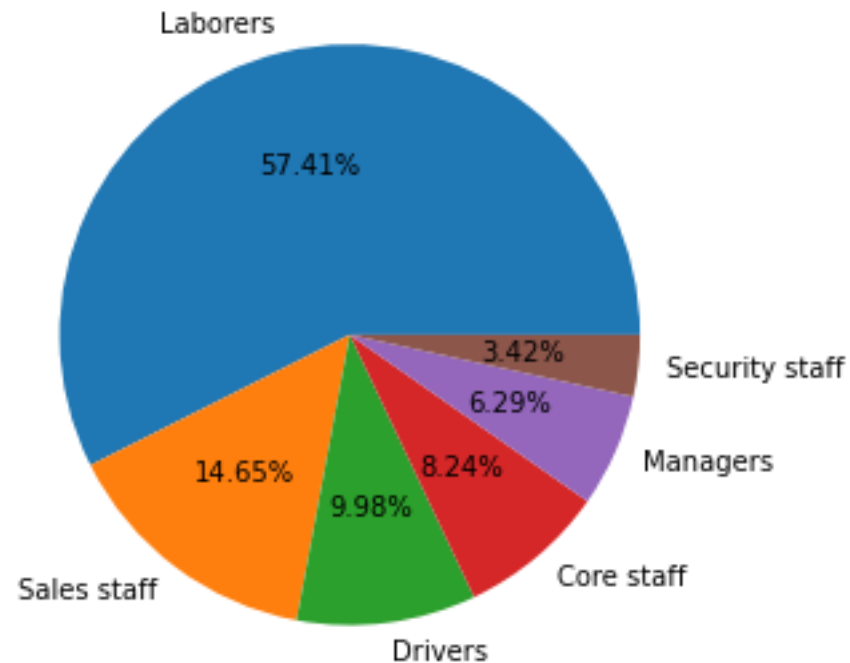COUNTPLOT - DEFAULTERS

COUNTPLOT - NON DEFAULTERS

# From the previous graph we took only 5-6 categories so as to more insight



Distribution for NON DEFAULTERS

Distribution for DEFAULTERS

# INFERENCES FROM PREIVOIUS SLIDE accept the applications of Non Defaulter where % is higher.

*Core Staff - 5.34% higher than defaulters*

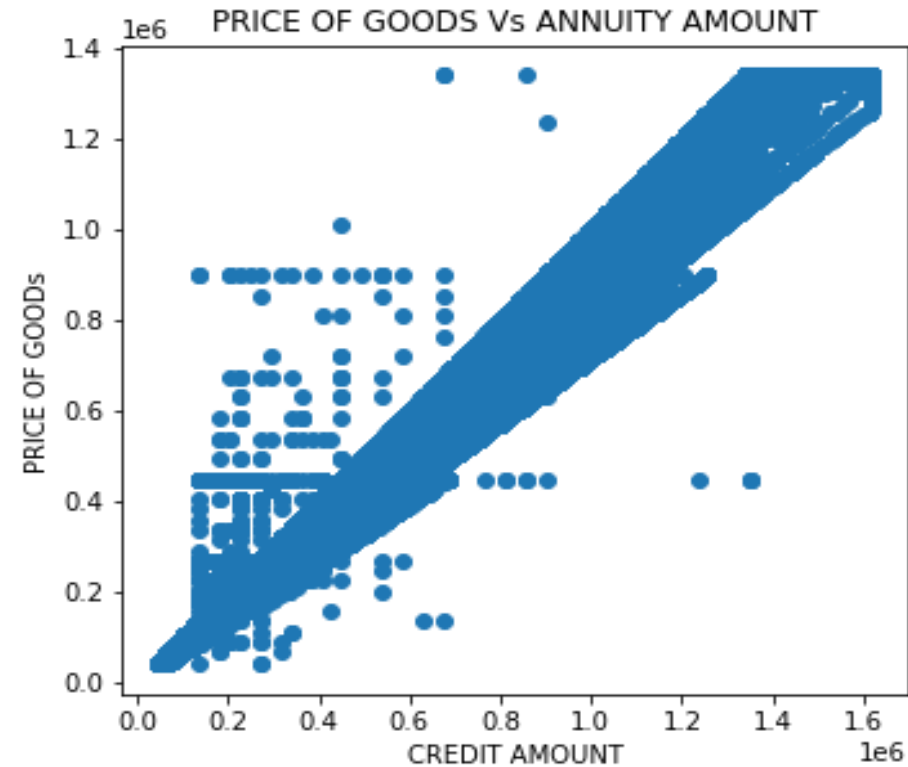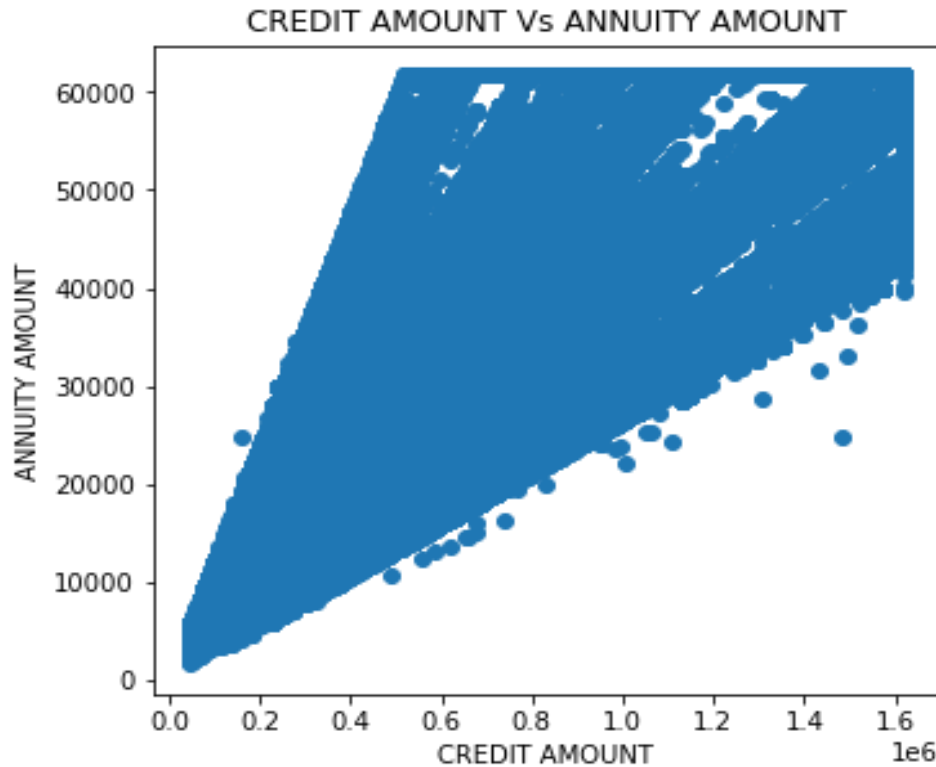*Managers - 1.94% higher than defaulters*

*High skill tech staff - Not listed in Defaulters list as we have taken only 6 maximum values*
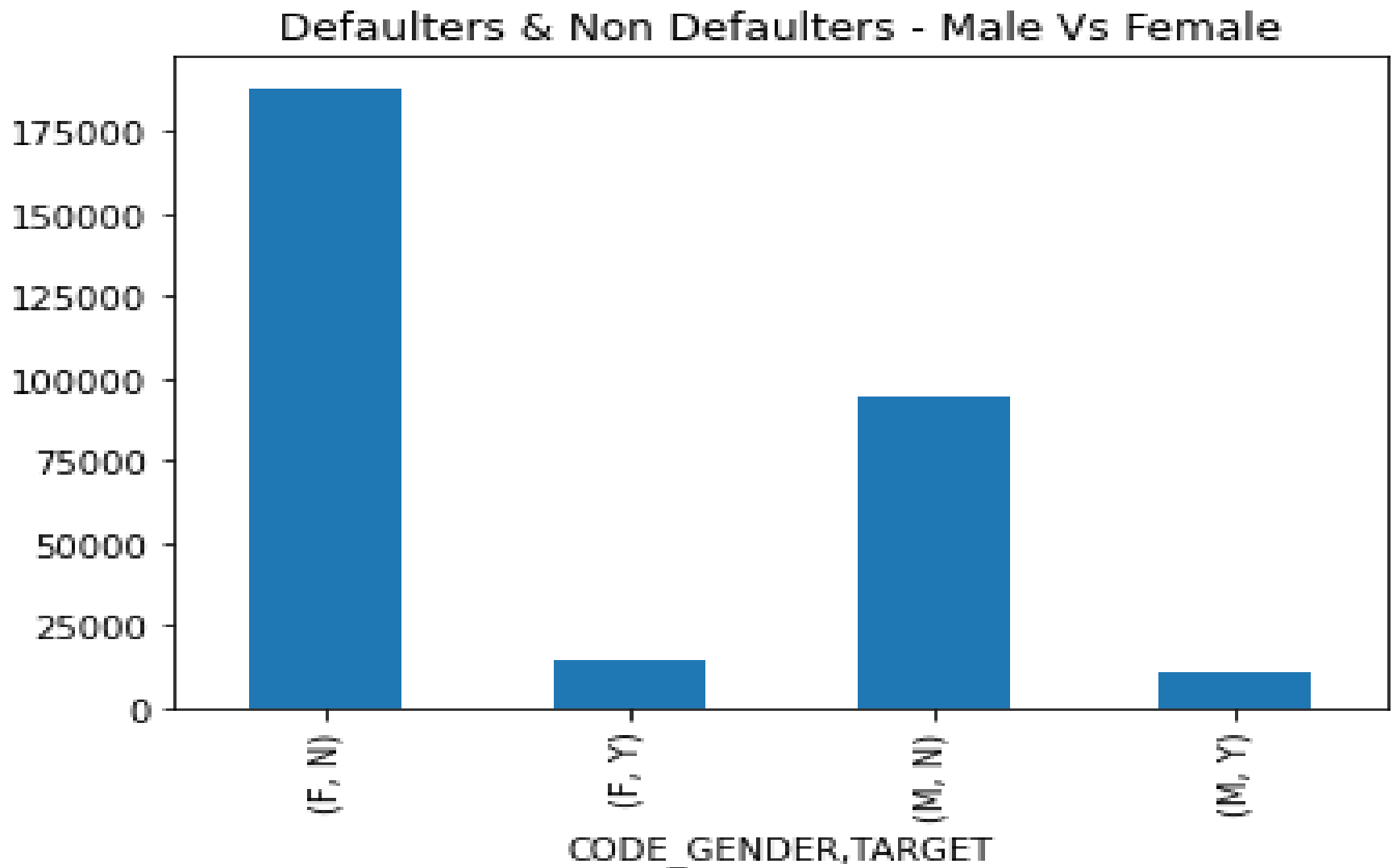
# BIVARIATE & MULTIVARIATE ANALYSIS

- Now we will find out the defaulters
- We will work upon the separate dataset for deaulters (TARGET = 1)
- We will try to plot relationships between different variables toget more insight

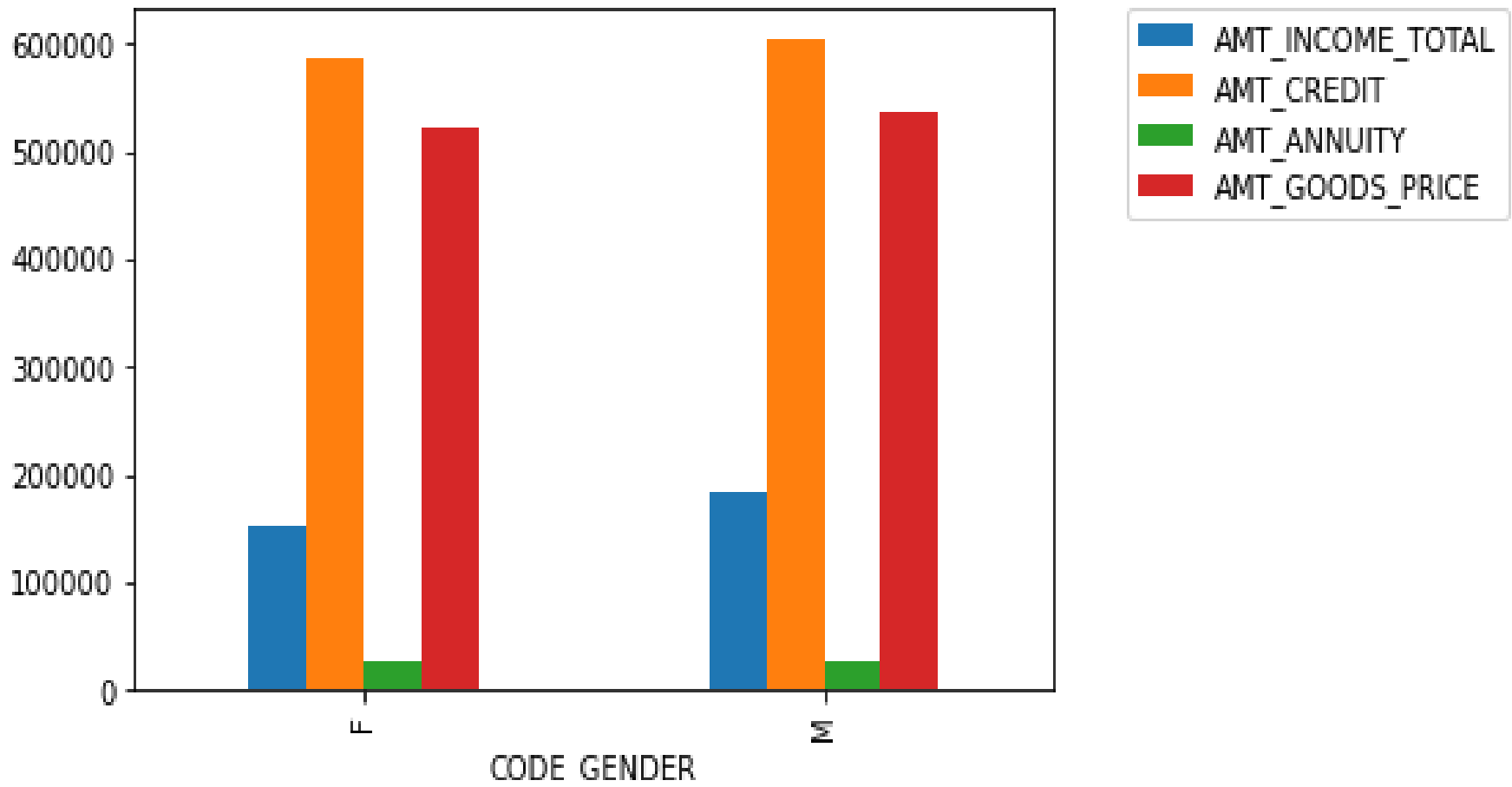# SCATTER SUBPLOTPLOT – CREDIT Vs ANNUITY AMOUNT, CREDIT Vs PRICE OF GOODS



➢**Positive correlation between ANNUITY & CREDIT AMOUNT**
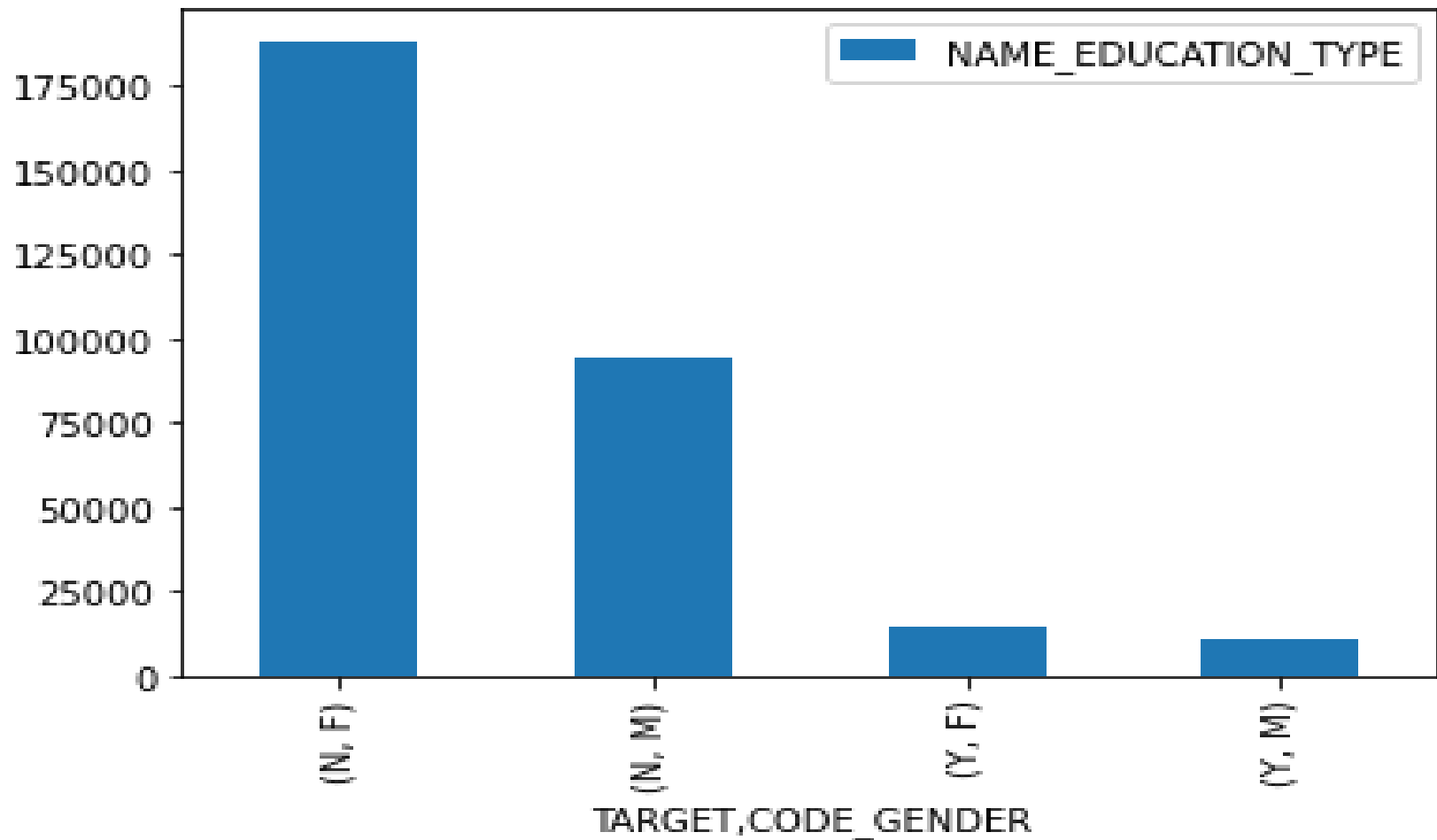➢**Positive correlation between PRICE OF GOODS & CREDIT AMOUNT**

Defaulters & Non Defaulters - Male Vs Female

➢**Number of female defaulter and non defaulter are more than males**

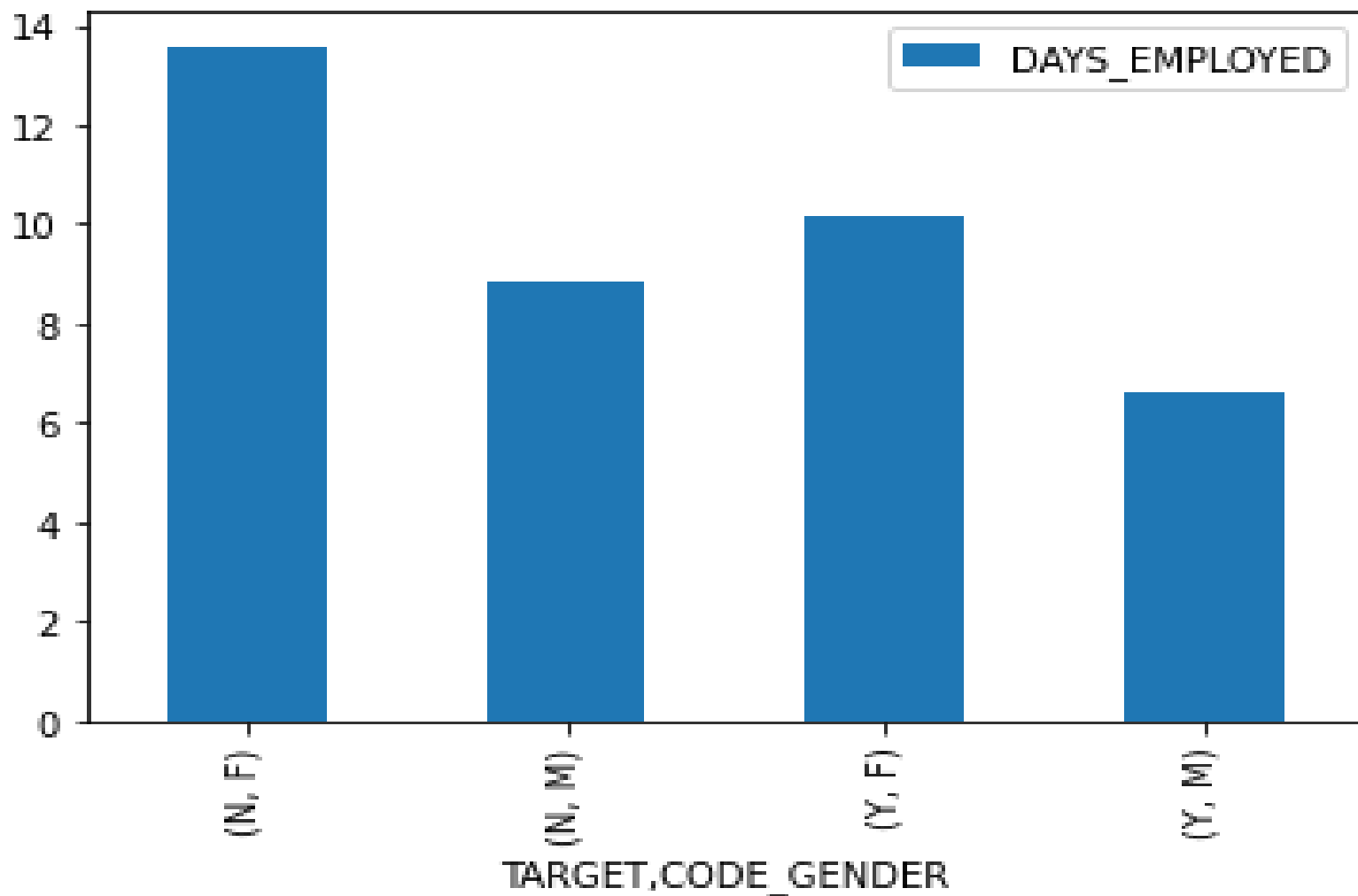GENDER VS INCOME, CREDIT AMOUNT, ANNUITY, GOODS PRICE

➤**Income of male is pretty higher than of female so we can say that males have less loan payment difficulty.**

➤**Males can get loan easily as compared to female.**

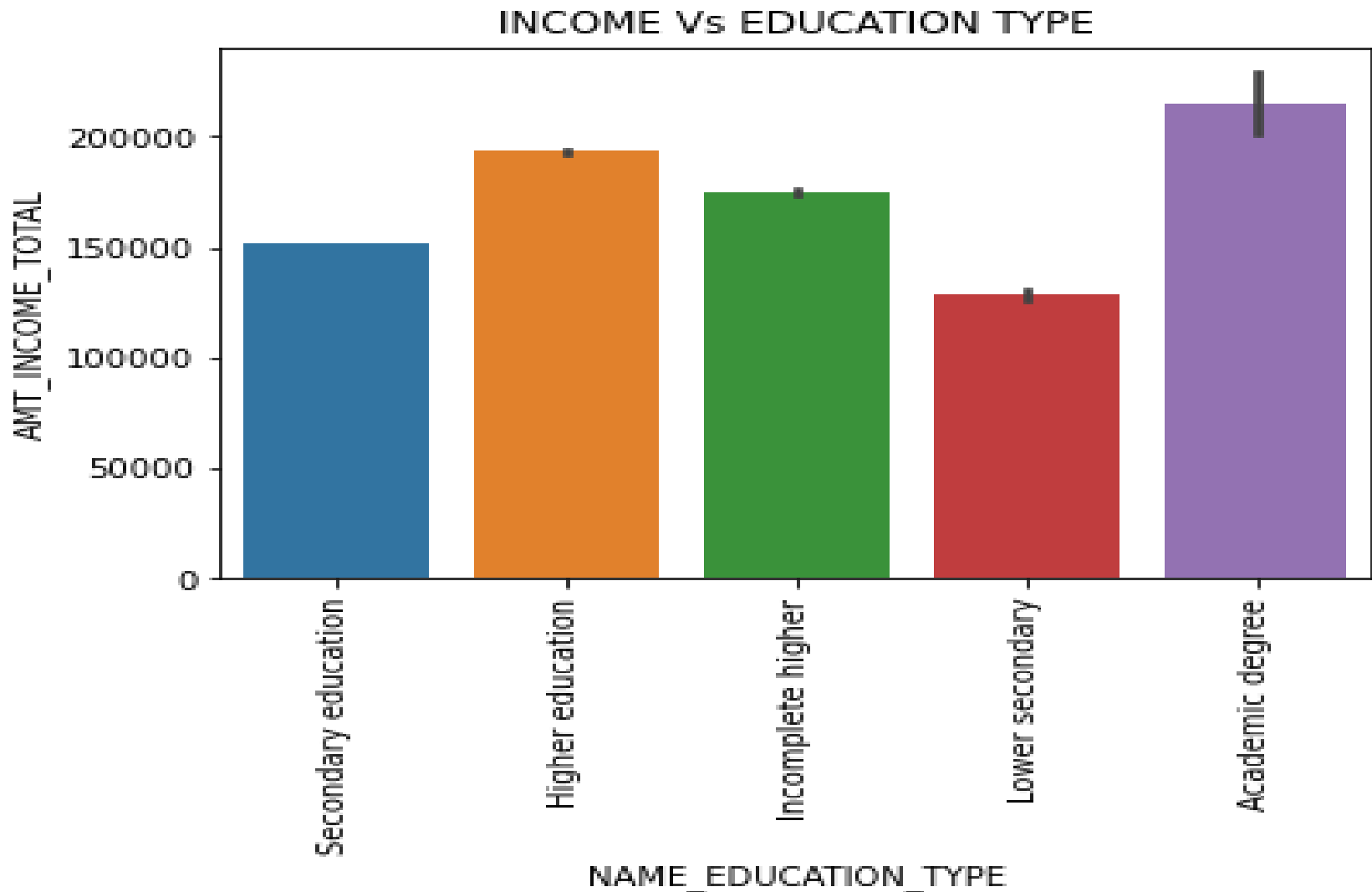DEFAULTERS - NON DEFAULTERS, GENDER Vs EDUCATION TYPE

➢ **Most of the Non defaulter women are educated as compared to males.**
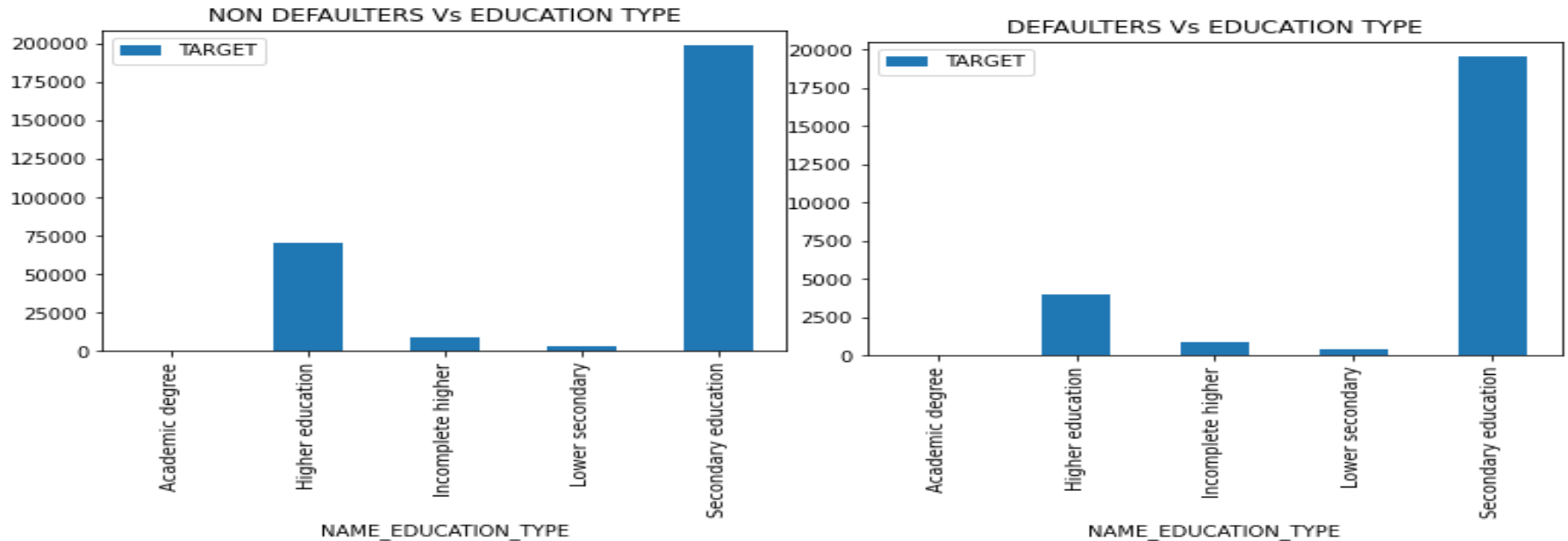
DEFAULTERS - NON DEFAULTERS, GENDER Vs DAYS EMPLOYED

➢ **Females employed days are higher than males**
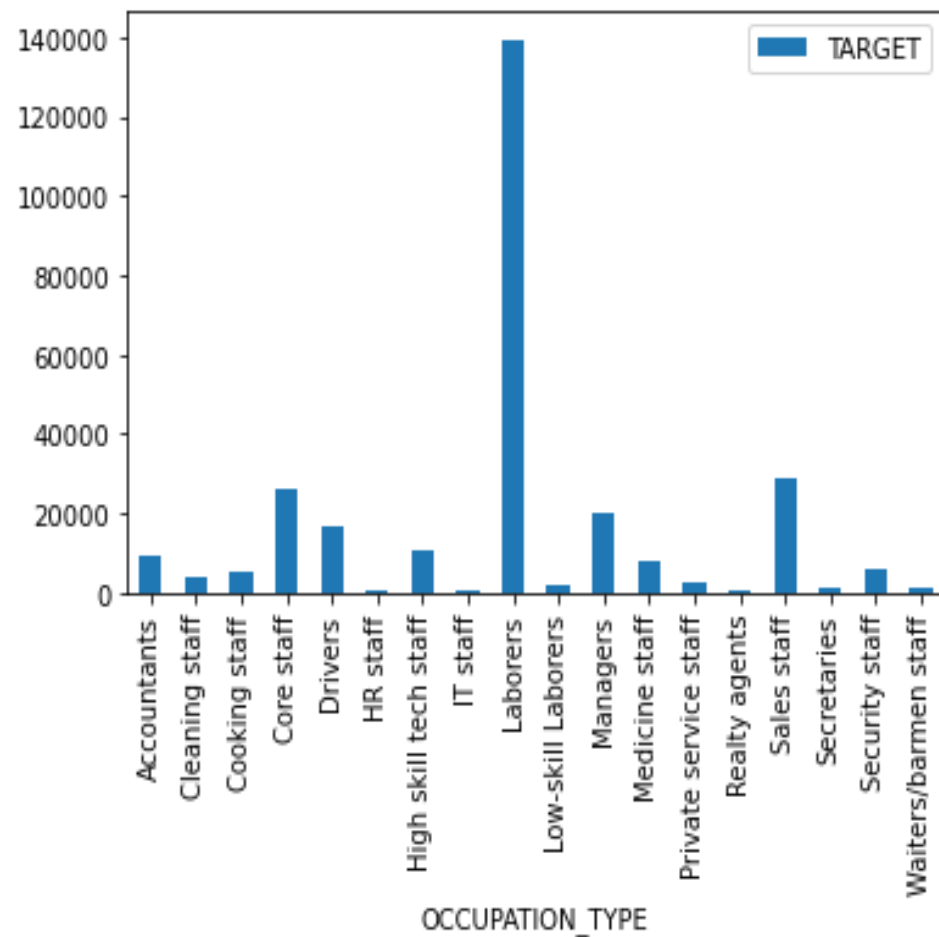
**INCOME Vs EDUCATION TYPE**

➢**People having academic degree have high median salary. So people with Academic Degree can be targeted**

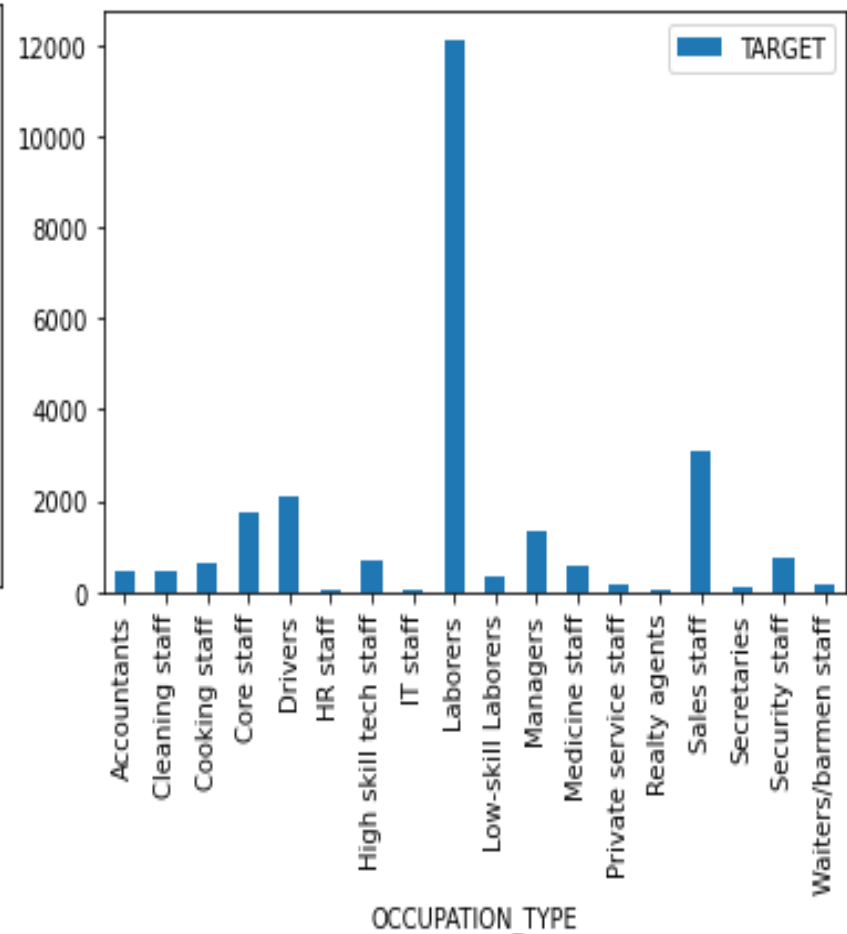# DEFAULTER – NON DEFAULTERS VS EDUCATION TYPE



➢**People having secondary education are highest number of defaulter and non defaulters**

➢**People having Higher education are more in non defaulter bucket so they can be targeted**
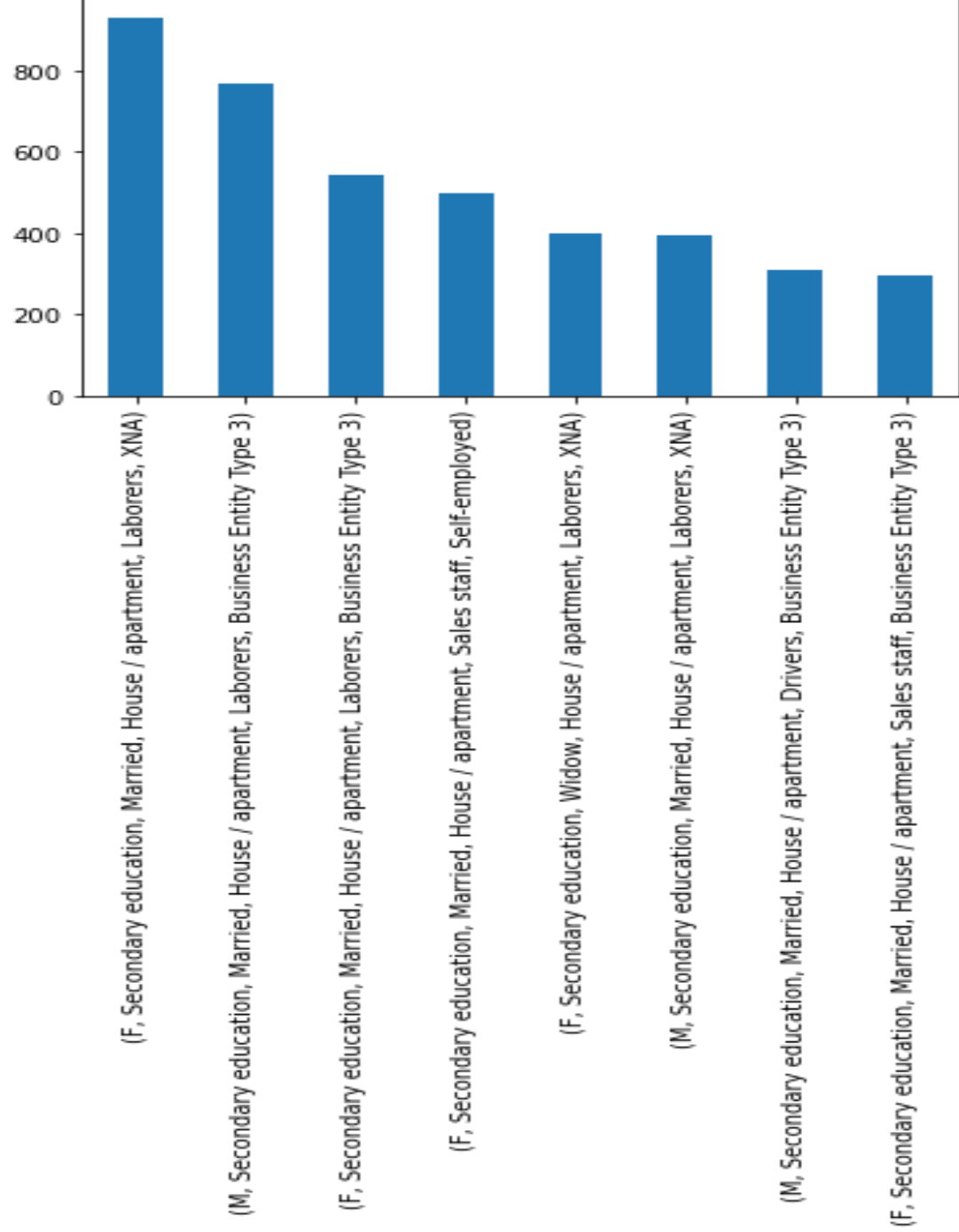
NON DEFAULTERS Vs OCCUPATION TYPE / DEFAULTERS Vs OCCUPATION TYPE

➢ **Maximum number of Labourers are the defaulters and non defaulters both**

CODE_GENDER,NAME_EDUCATION_TYPE,NAME_FAMILY_STATUS,NAME_HOUSING_TYPE,OCCUPATION_TYPE,ORGANIZATION_TYPE
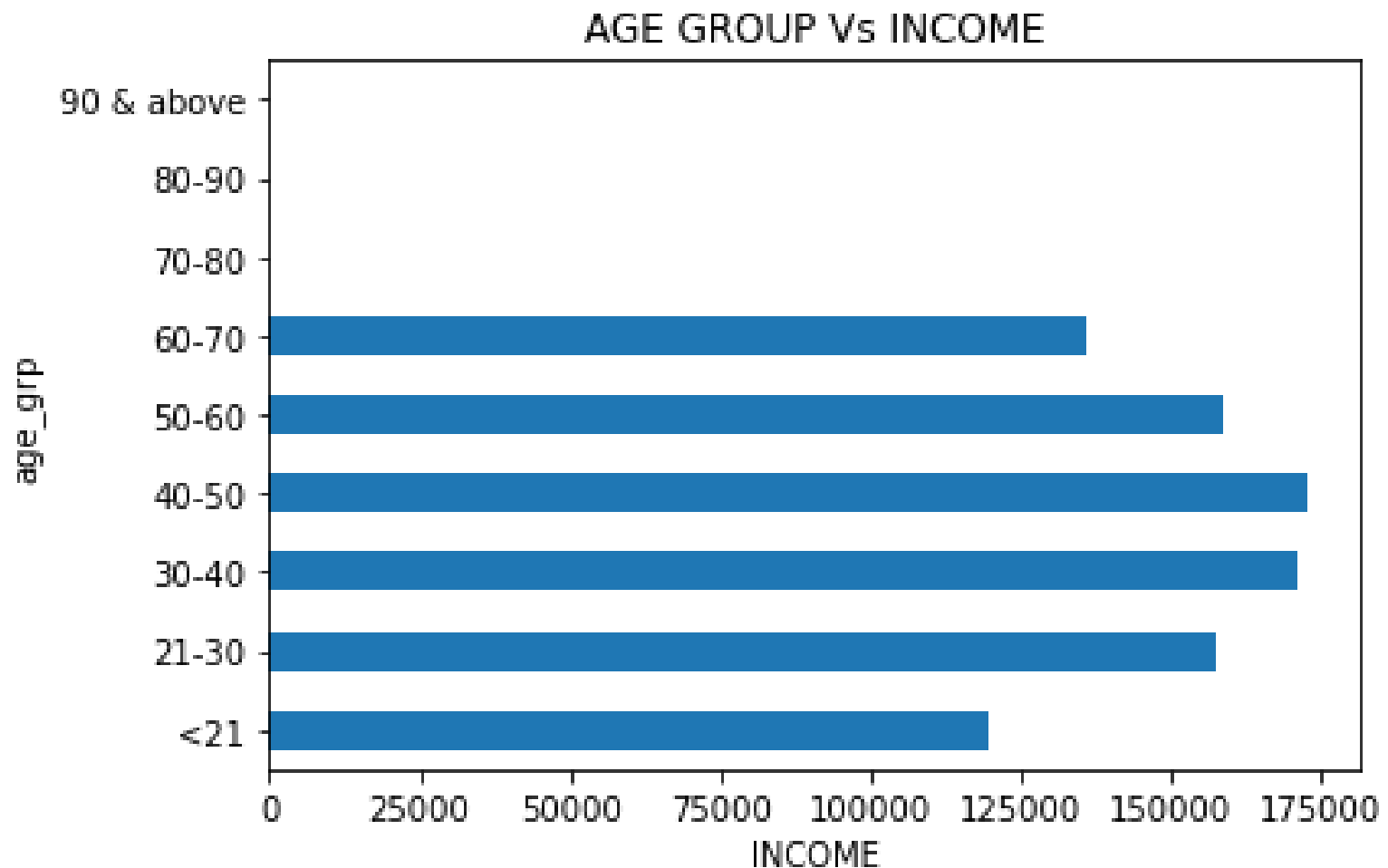
# INSIGHT

MAJOR DEFAULTERS

MALE & FEMALE

MARRIED

LIVING IN APARTMENT

WORKING IN A BUSINESS ENTITY TYPE 3

**Maximum number of the people applying for the loan are between 40-50 age and have high mean salary**

# PREVIOUS APPLICATION DATASET

Data cleaning approach is same as in application dataset

1. • Importing Libraries

2. • Reading the data set and finding percentage of null values

3. • Dropping columns with missing values >45%

4. • Identifying continuous and categorical columns/variable

5. • Continuous column - Columns containing unique values > 58

6. • Categorical column – Columns containing unique values < 58

7. • Imputation for missing value<45 (categorical – mode, continuous – median)
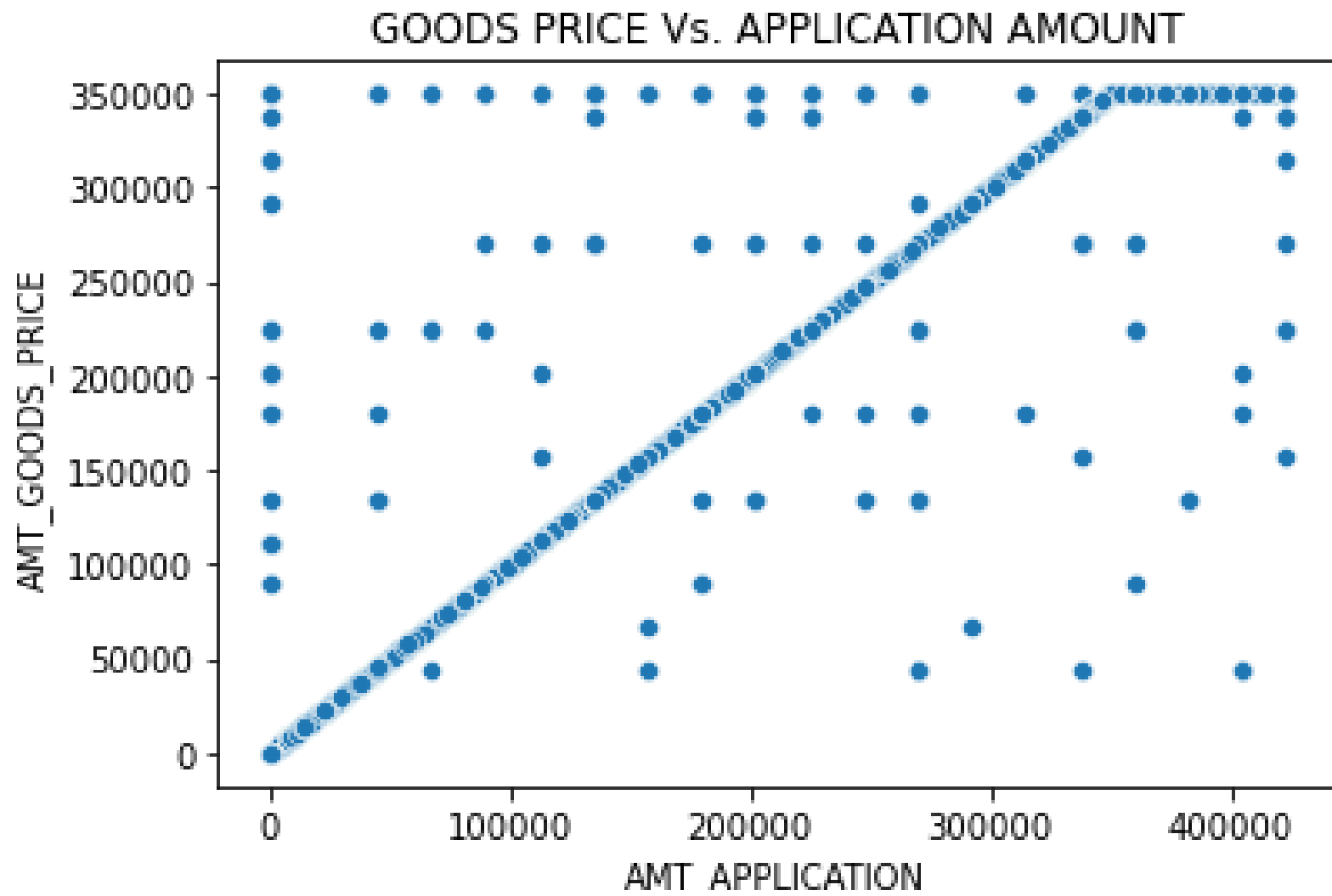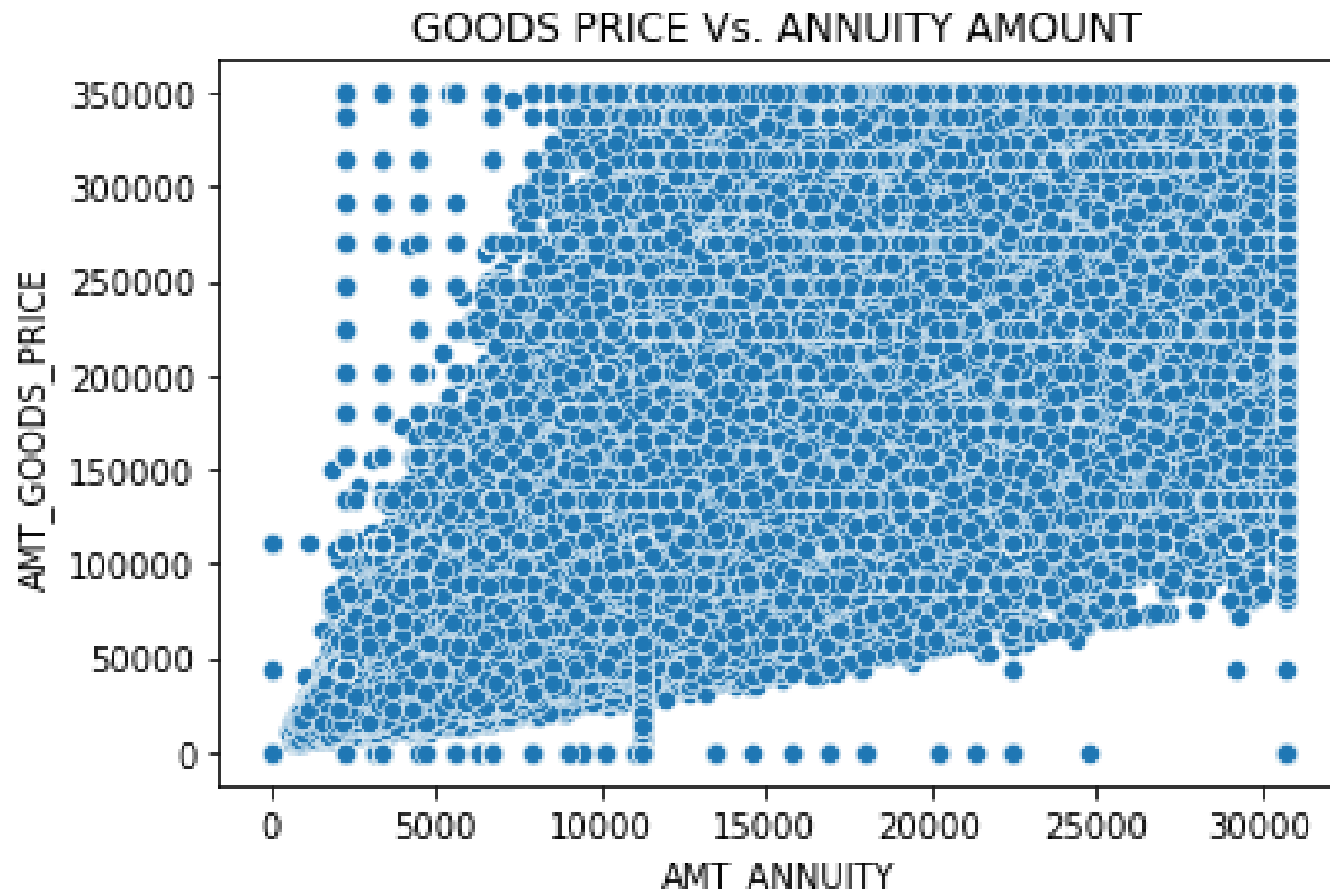
8. • Dropping unnecessary columns

9. • Detecting Outliers – Using Subplots

10. • Handling outliers by flooring and capping

GOODS PRICE Vs. APPLICATION AMOUNT
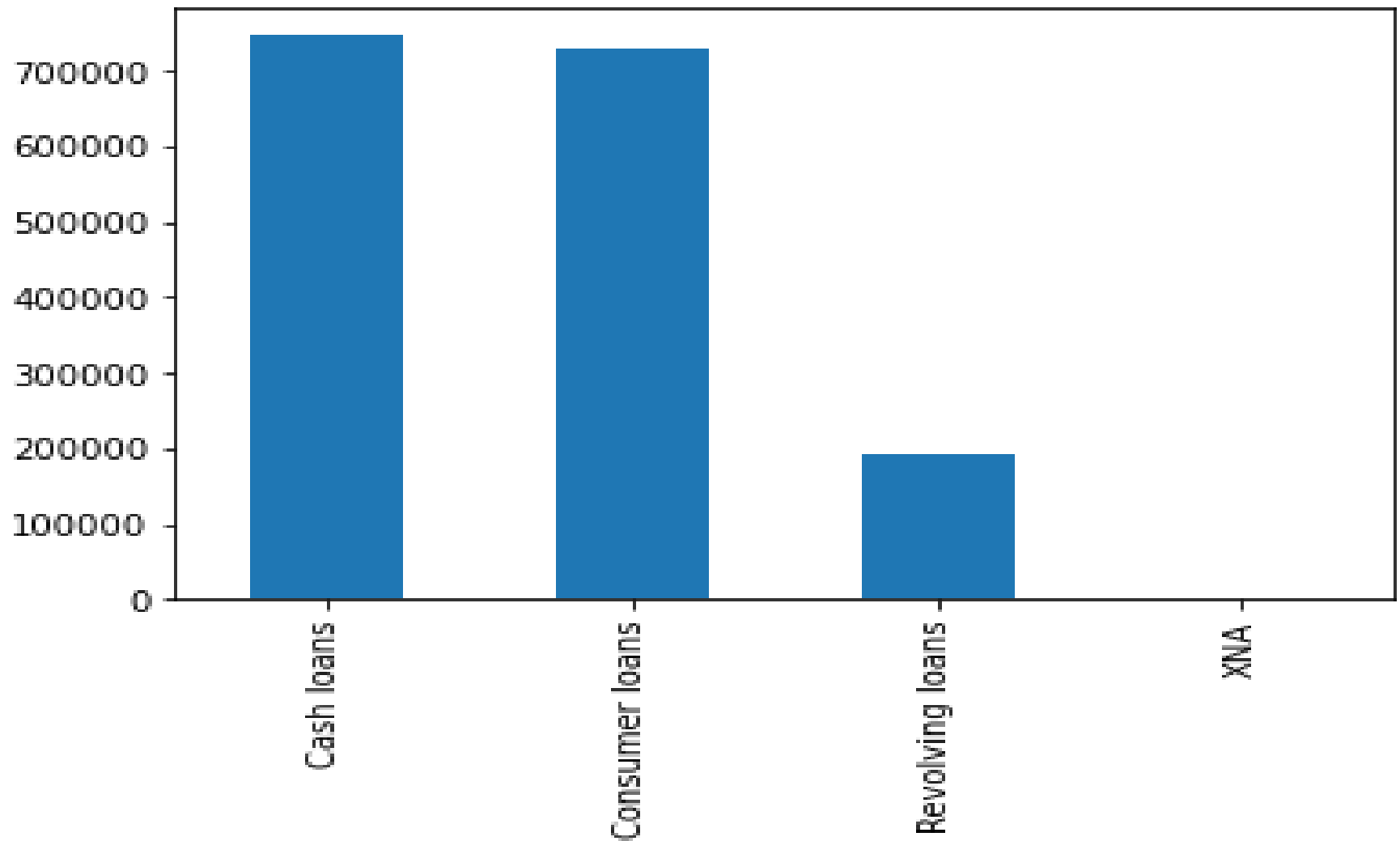
➢**As the goods price increases amount of loan application also increases**
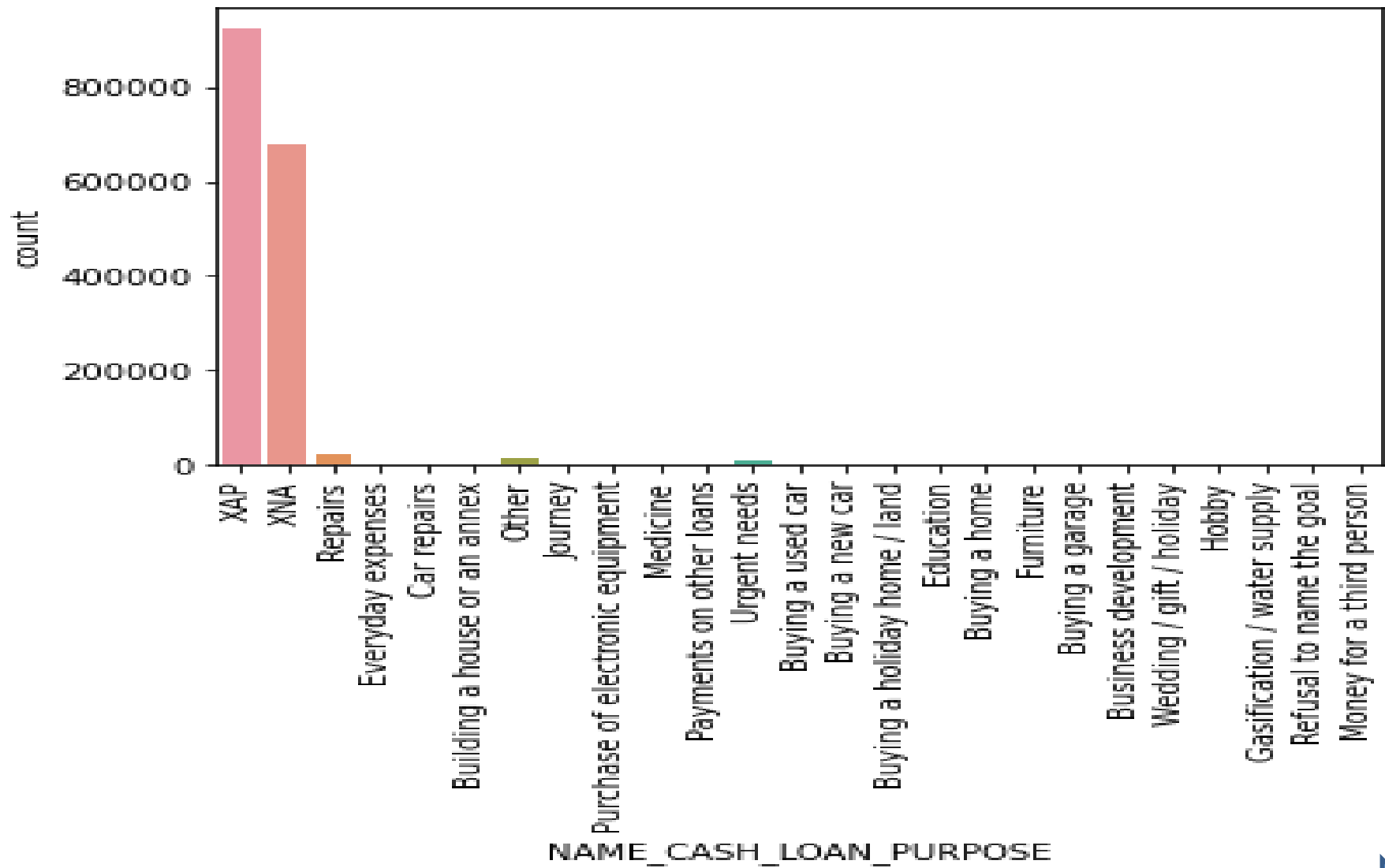
GOODS PRICE Vs. ANNUITY AMOUNT

➢**As the goods price increases annuity amount also increases**
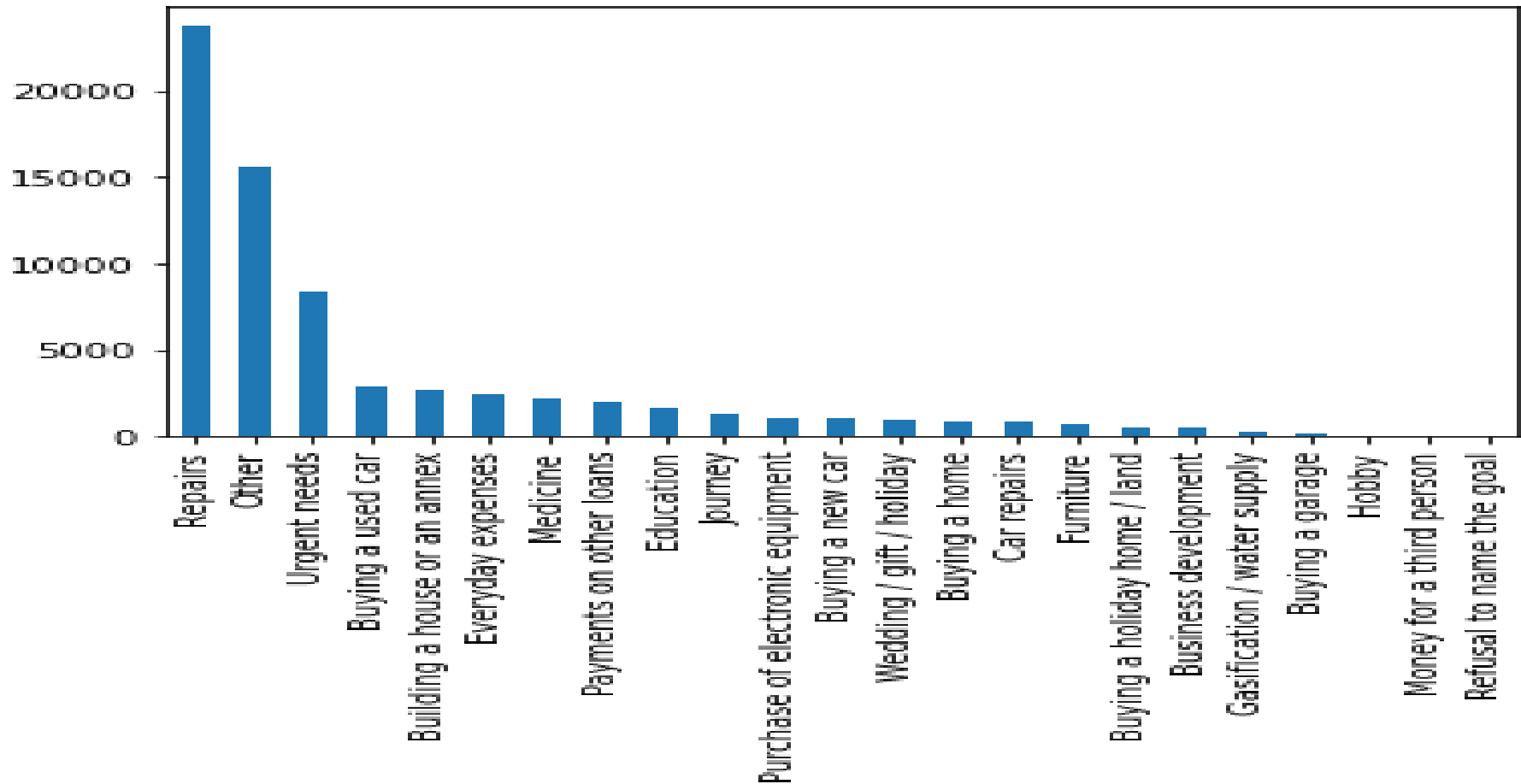
# COUNT Vs CONTRACT TYPE

➢**People previously applied for cash loans more**
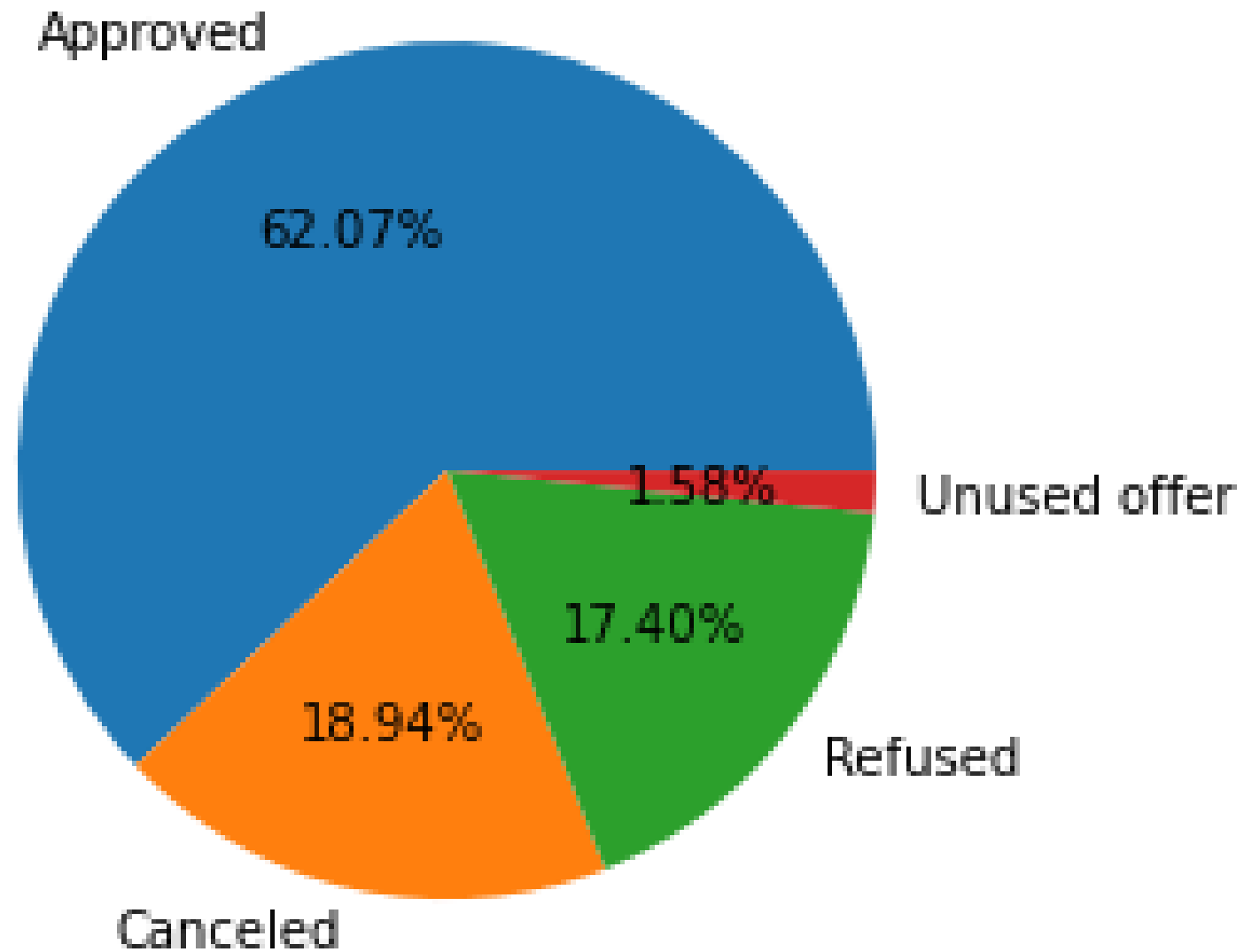
COUNT Vs CASH LOAN PURPOSE

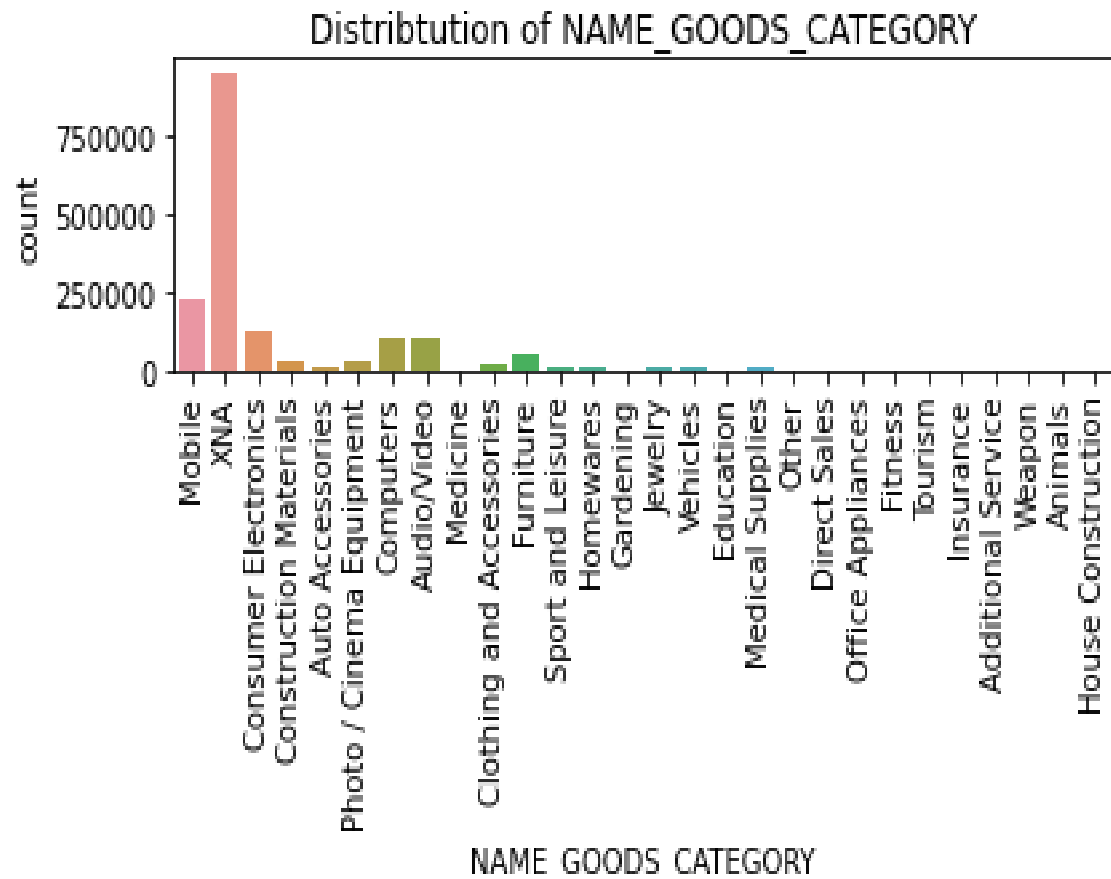IGNORING XAP, XNA. GOING DEEPER FOR OTHER CASH LOAN PURPOSE
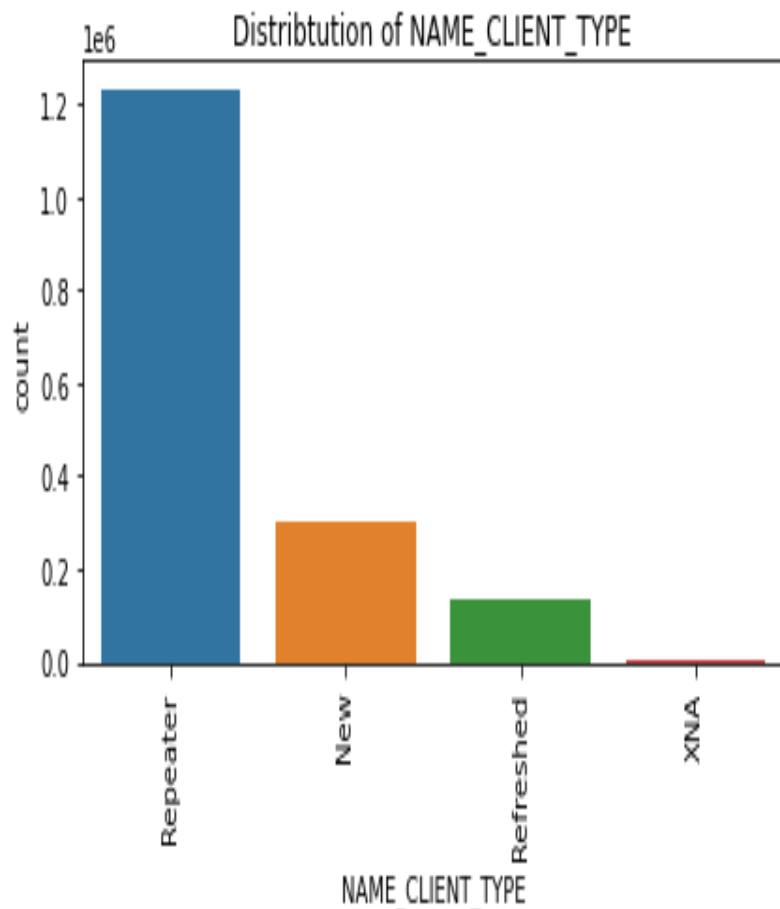
COUNT Vs. CASH LOAN PURPOSE - FILTERED DATASET

➤**Maximum number of people are applied for loan for repair purpose**
➤**Secondly people applies, for other reasons**
➤**People also require loan for urgent needs which are not specified**
➤**People applies loan for buying a used car, building a house.**

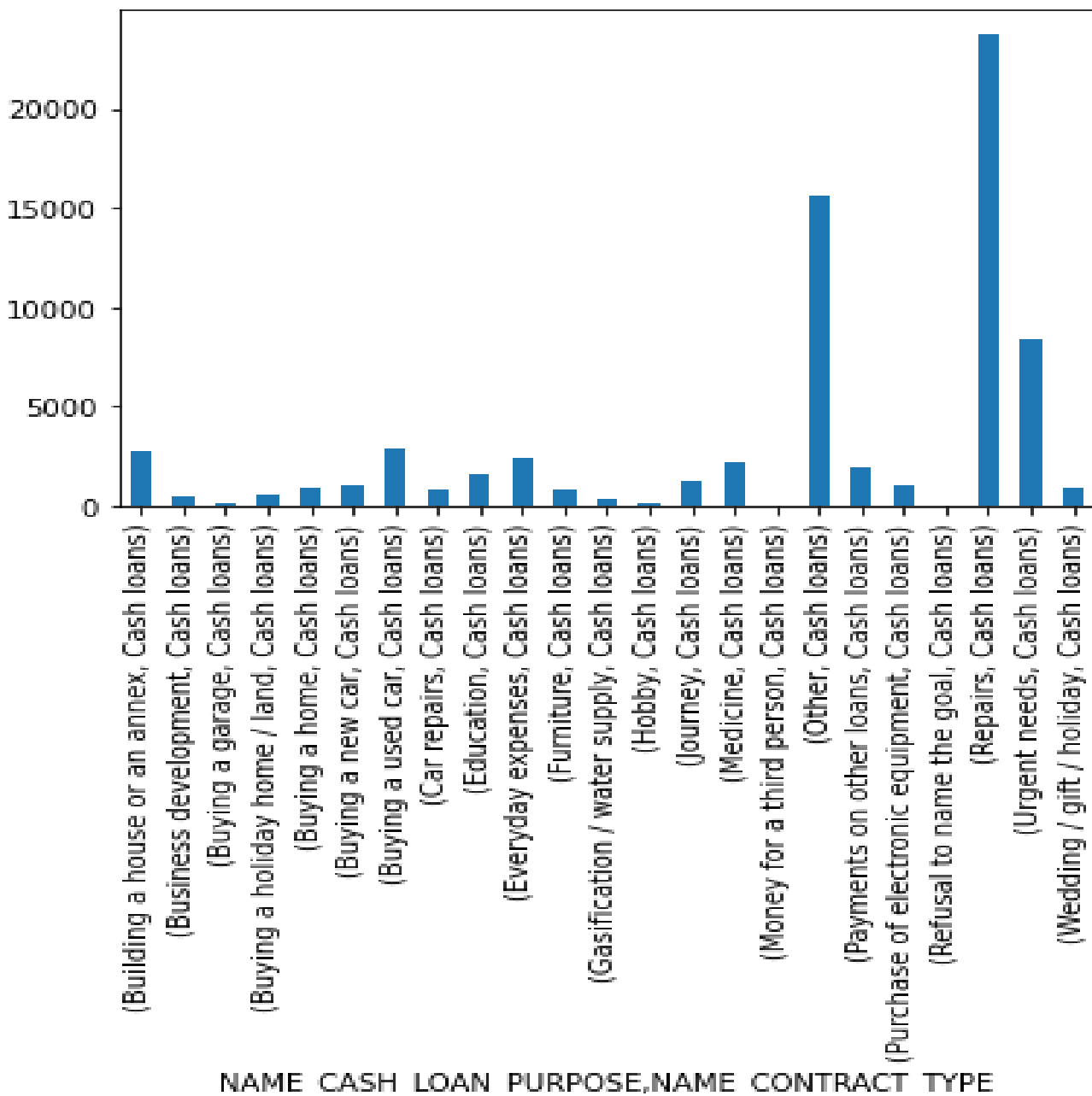# Pie chart for Application Status

Distribtution of NAME_CLIENT_TYPE

Distribtution of NAME_GOODS_CATEGORY

➤**Maximum number of old clients are applying for the loan again**

➤**Other than XNA, people are applying loan for mostly mobile phones.**

**Maximum number of people are taking cash loans for Repair**

LOAN PURPOSE Vs LOAN TYPE without XNA, XAP

NAME_CASH_LOAN_PURPOSE,NAME_CONTRACT_TYPE

**Maximum number of loans that are approved is Consumer loans**

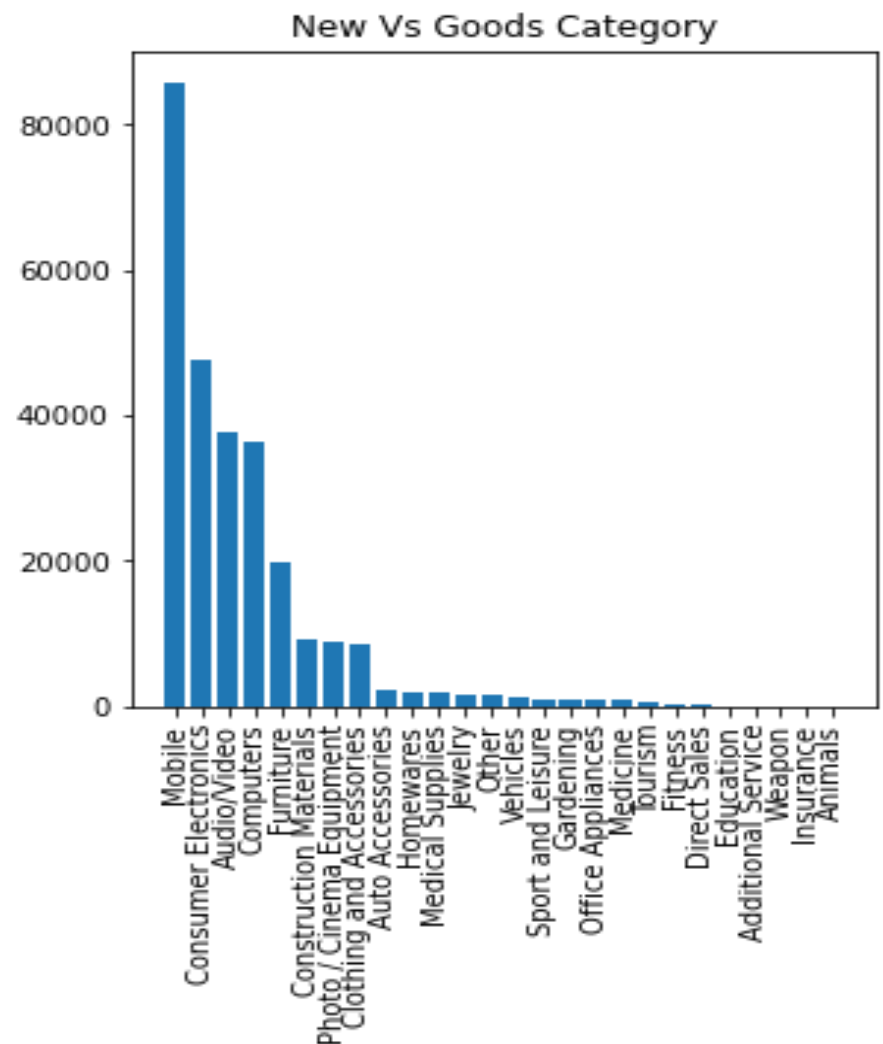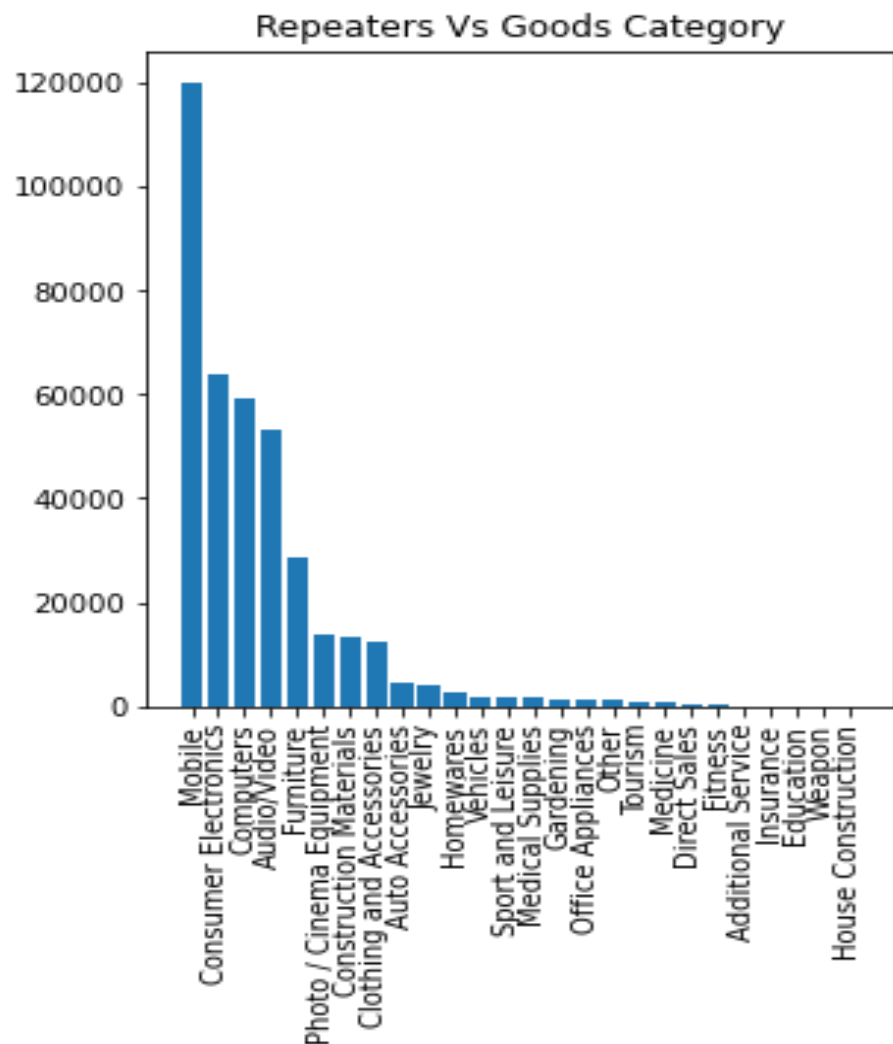**Cancellation of Consumer loans is very less as compared to cash and revolving loans**

LOAN TYPE Vs LOAN STATUS

NAME_CONTRACT_TYPE,NAME_CONTRACT_STATUS

**Company approves the loan for maximum number of people who are applying for the new loan again**

**For People applying for the first time, company tries not to cancel the loan application that is why cancellation for new people is very less as compared to refreshed and repeaters**
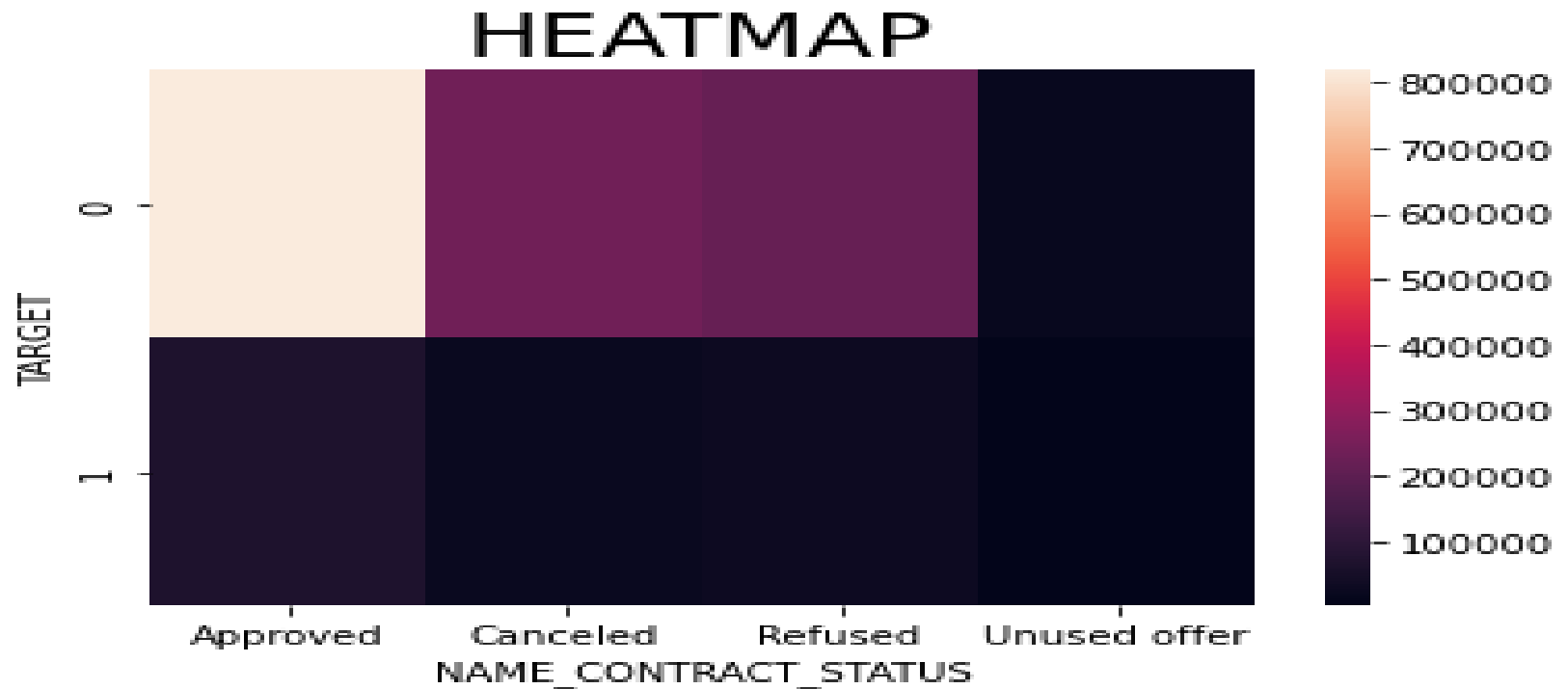


CLIENT TYPE Vs LOAN STATUS

**Maximum number of people ( Repeaters and New ) are applying loan for mobile phones.**

Merging two data sets → Application & Previous Application → Drawing Inferences

# HEATMAP



➢ **Maximum number of loan is approved for non defaulters**
➢ **Loans are approved for defaulters also.**
➢ **Loans are refused for non defaulters also.**

# MAJOR – INSIGHT

By taking mode of categorical columns with **Target = 0**, **CONTRACT_STATUS = Refused**
By taking mode of categorical columns with **Target = 1**, **CONTRACT_STATUS = Approved**

LOAN REFUSED PREVOUSLY FOR DEFAULTERS
LOAN APPROVED PREVIOUSLY FOR NON DEFAULTERS

APPLYING FOR CASH LOANS

FEMALE, MARRIED, HAVE SECONDARY EDUCATION

DOESNOT OWNS CAR, BUT HAVE HOUSE

LABOUR BY OCCUPATION AND WORKING IN BUSINESS ENTITY 3