

# Bank Loan Case Study

Arcot Navya Sai ([navyasaipatnaik@gmail.com](mailto:navyasaipatnaik@gmail.com)) +91-8074851394

**Project Description:** This Project is about handling the unclean data in the industry where we will be focusing on Exploratory Data Analysis on the given Dataset. As the case study we are dealing with is the Bank Loan Case Study we will be developing 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.

In the current Dataset we are having three datasets namely

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 whether the previous application had been Approved, Cancelled, Refused or Unused offer.
3. ``columns_description.csv`` is data dictionary which describes the meaning of the variables.

**Approach:** As we are handling the unclean data our first step is to understand the data and the columns we are having three datasets where one dataset `columns_description.csv` is totally explaining about columns of the other two datasets namely `application_data.csv`, `previous_application.csv`.

After understanding the data, we will explore the columns and description of the data along with info and later we will be handling the missing data in the columns and rows, we will also drop the irrelevant rows and columns.

Later we will do the analysis of outliers and remove them followed by Univariate and Bivariate analysis on the discrete and the continuous variables Finally we summarize the data and the analysis by visualizing the relations in excel.

**Tech-Stack Used:** In this Case Study We have used

1. Microsoft Excel Professional Plus 2016 version
2. Jupyter Notebook from Anaconda 3
3. Python 3
4. Pandas
5. Numpy etc. and all the necessary libraries in python

**Insights:** Initially we have started with exploring the columns, info, description etc. Of the dataset given in the Jupyter Notebook.

Missing values in application Dataset: <https://drive.google.com/file/d/1jI5YVqZ-AiRK7LhOvOjzX-JnKhKmJcsT/view?usp=sharing>

We have removed the columns which were irrelevant for the data Analysis

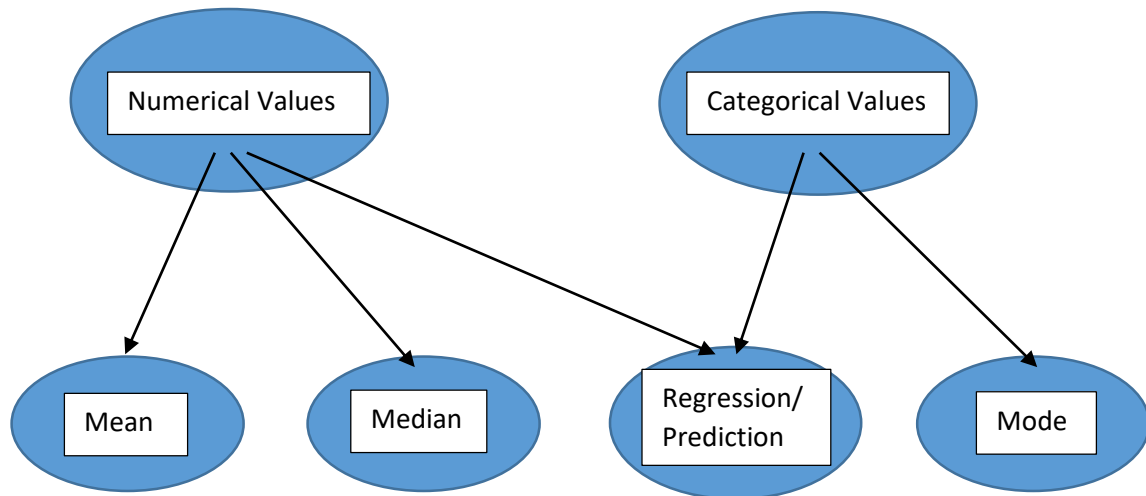
'DAYS\_REGISTRATION','FLAG\_MOBIL','FLAG\_EMP\_PHONE','FLAG\_WORK\_PHONE','FLAG\_CONT\_MOBILE',  
'FLAG\_PHONE','FLAG\_EMAIL','WEEKDAY\_APPR\_PROCESS\_START','HOUR\_APPR\_PROCESS\_START','LIVE\_REGION\_NOT\_WORK\_REGION','REG\_CITY\_NOT\_LIVE\_CITY','REG\_CITY\_NOT\_WORK\_CITY','LIVE\_CITY\_NO

T\_WORK\_CITY','DAYS\_LAST\_PHONE\_CHANGE','OBS\_30\_CNT\_SOCIAL\_CIRCLE','DEF\_30\_CNT\_SOCIAL\_CIRCLE','OBS\_60\_CNT\_SOCIAL\_CIRCLE','DEF\_60\_CNT\_SOCIAL\_CIRCLE','NAME\_TYPE\_SUITE'

After the Whole Cleaning Process, the size of the application data set is (307511, 51) nearly 58% of columns were removed many were irrelevant and others had large missing values

Then we checked for the rows having more than 50% of missing values.

How did we handle the missing values?



**Outlier Analysis:** For the Outliers we have used the Z score to Identify and remove the outliers in the numerical columns. (294063, 53) is the shape of the Dataset. After the outlier analysis 4.37% of rows got deleted now we have the data set where there are no missing values and the outliers.

### Data Imbalance:

Then we did the Data Imbalance Analysis to check the Numerical Data and Categorical Data  
We can see that there is data imbalance in below columns: -

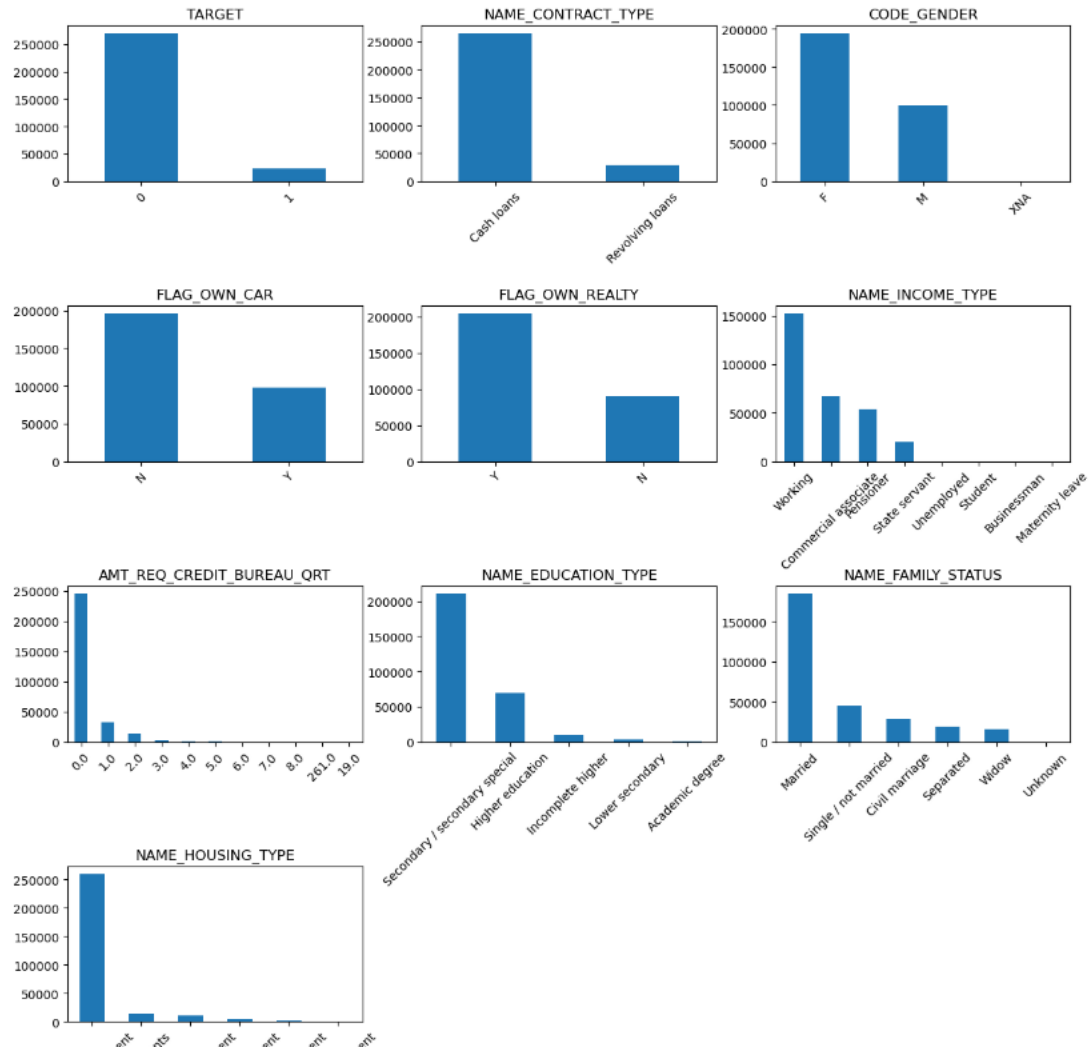
TARGET - There are very few defaulters (1) compare to non-defaulters (0)

NAME\_CONTRACT\_TYPE - There are very few Revolving loans than Cash loans

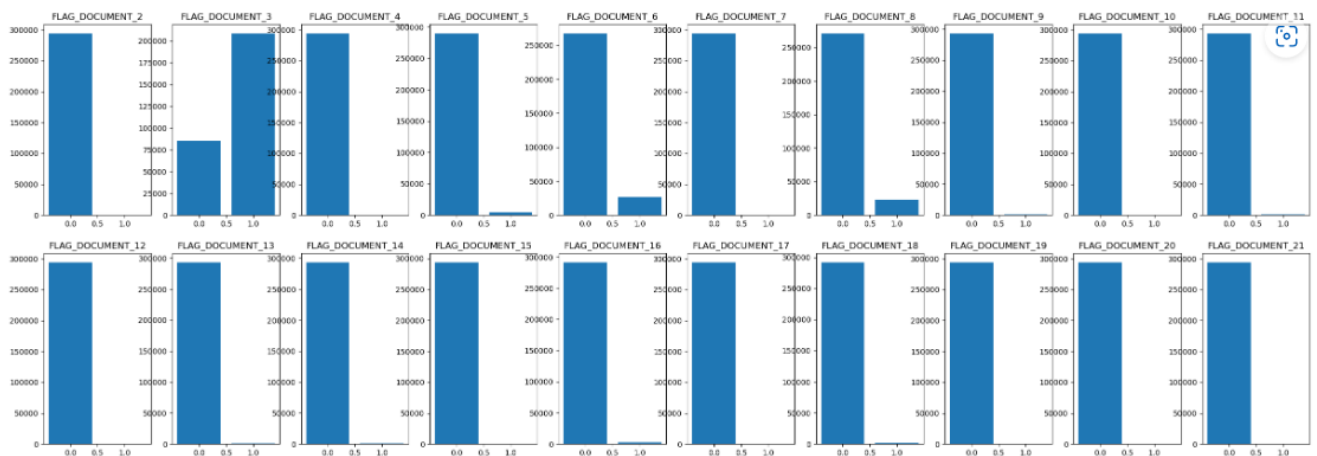
NAME\_EDUCATION\_TYPE - Most of the loans applied by Secondary/Secondary special educated people

NAME\_FAMILY\_STATUS - Most of the loans applied by Married people.

NAME\_HOUSING\_TYPE - Most of the application came from Home/apartment owner



The we understood the Flag data columns Data Imbalance



There is the data only in Flag Document 3

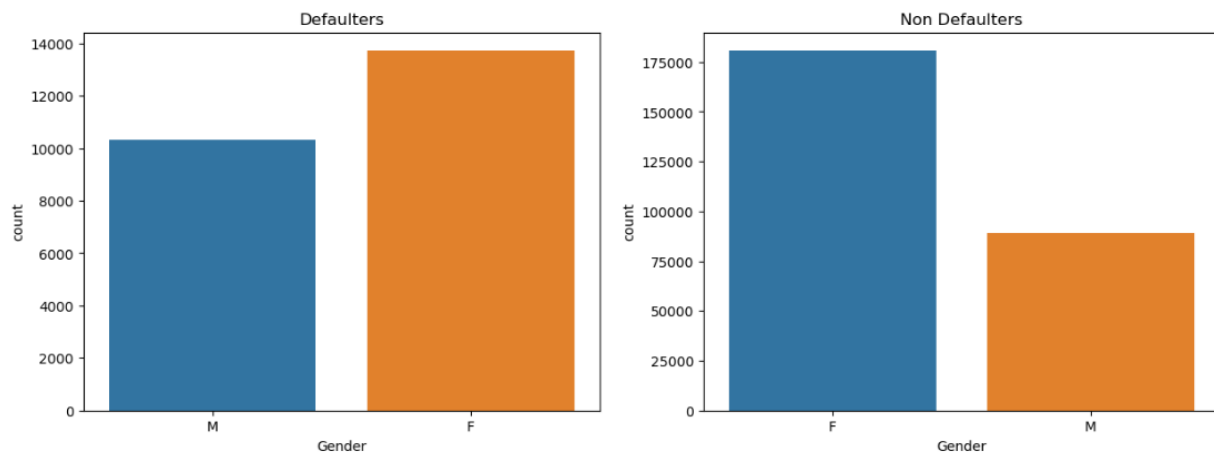
**Binning the columns:** Binning is the process of grouping a set of numerical values into a smaller number of "bins" or "buckets".

## Binning of AGE, AMT\_INCOME\_TOTAL, AMT\_CREDIT and EXT\_SOURCE\_SCORE columns

AGE	Years_Employed	Age_Group	Income_Group	Credit_Group	EXT_SOURCE_SCORE	Ext_source_group
26	2	Young	High	Low	0.20	Low
46	3	Mid Age	High	High	0.52	Medium
52	1	Mid Age	Low	Low	0.54	Medium
52	8	Mid Age	Medium	Low	0.60	High
55	8	Mid Age	Medium	Medium	0.38	Low

### Univariate analysis:

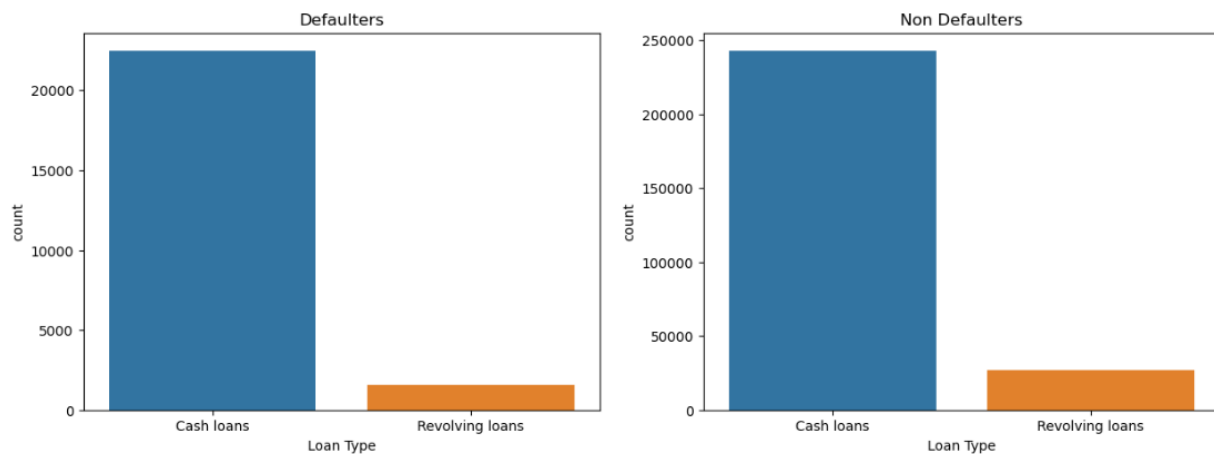
#### Count of defaulters and non-defaulters on the basis of gender



### Analysis

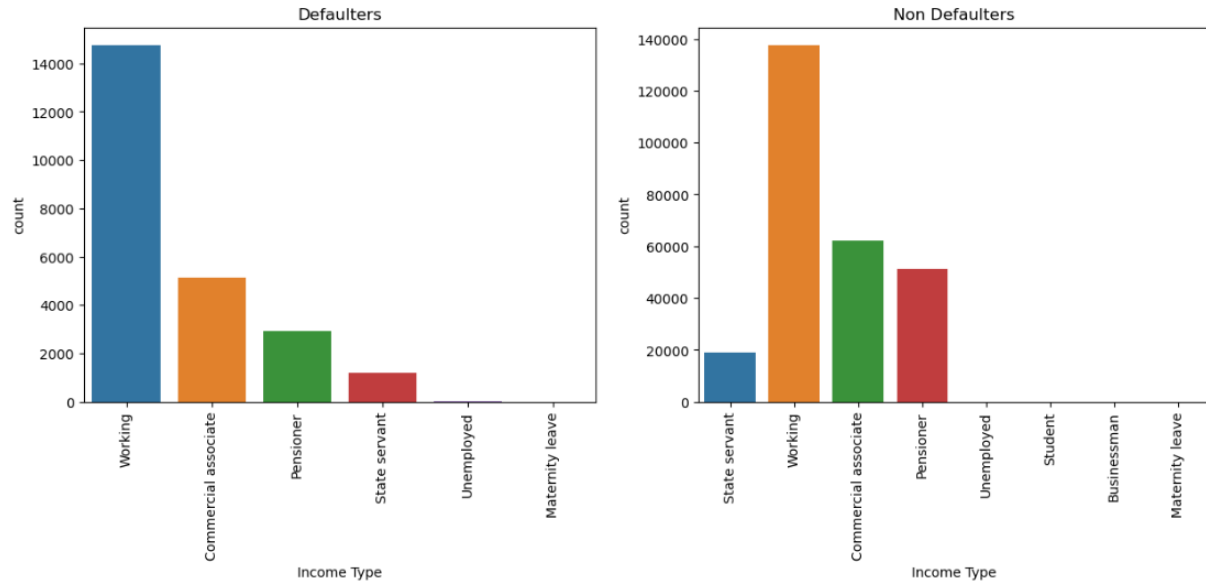
1. Defaulters - We can see that females are slightly more in number of defaulters than male.
2. Non-defaulters - The same pattern continues for non-defaulters as well. The females are more in number here than male.

#### Count of defaulters and non-defaulters on the basis of Loan Type



As we can see the Revolving loans are very less in number both in Defaulters and non-Defaulters

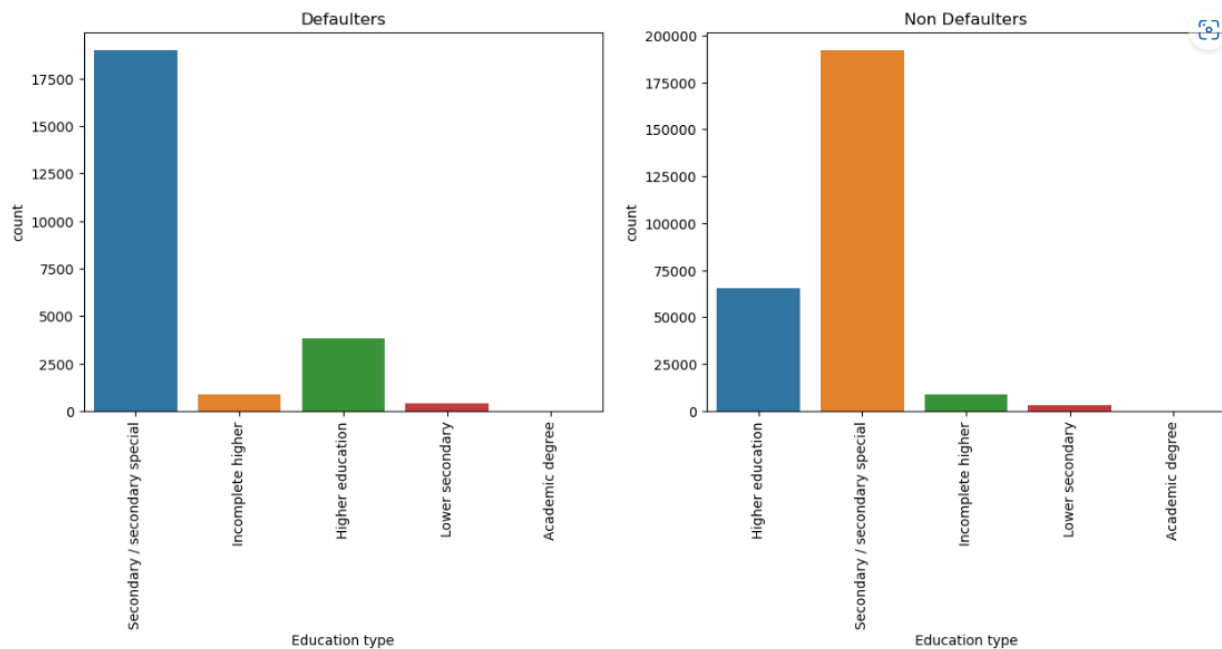
### Count of defaulters and non-defaulters on the basis Income Type



1. As we can see from the above plot Working type of income is highest in number both in Defaulter and non-defaulters

2. And also State Servant is comparatively low among all the types

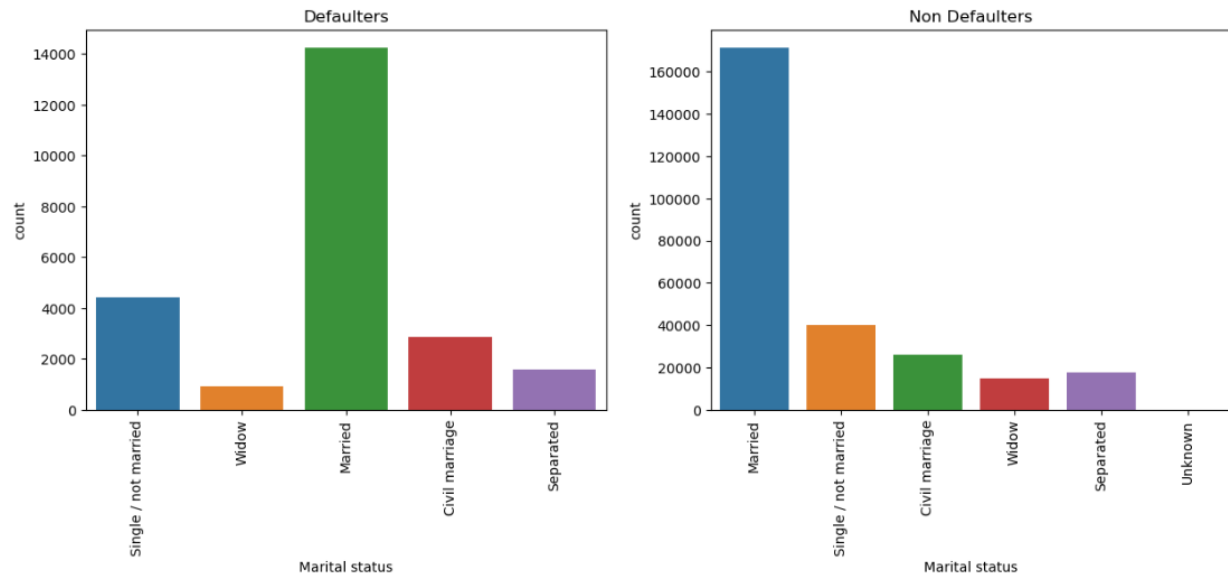
### Count of defaulters and non-defaulters on the basis Education type



### Analysis

The Secondary Special is comparatively high both in defaulters and non-defaulters

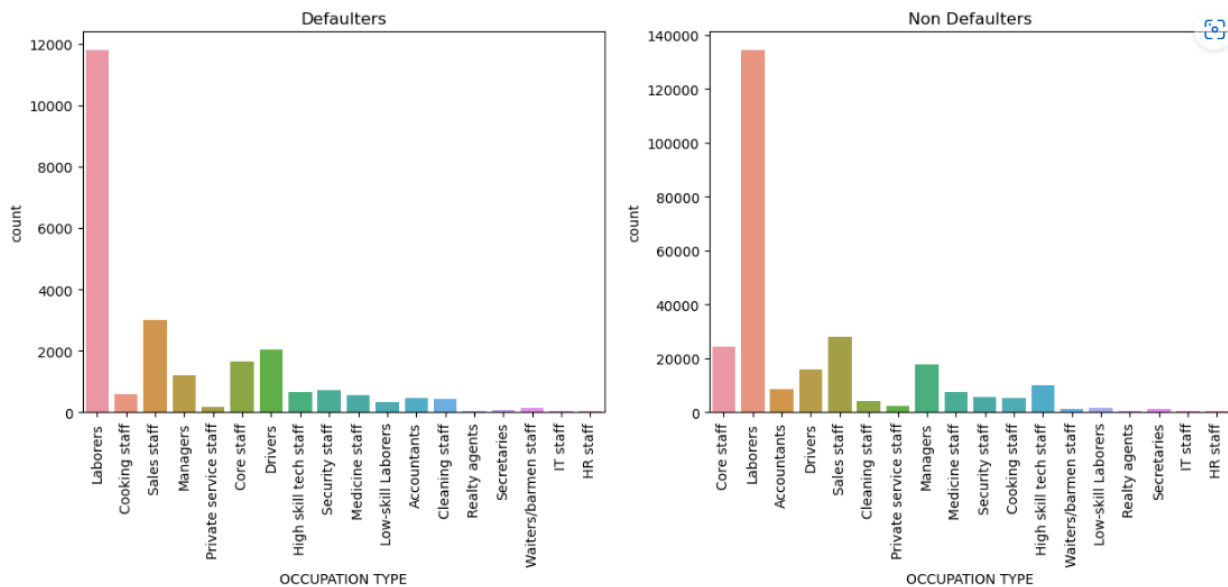
### Count of defaulters and non-defaulters on the basis Marital status



### Analysis

We can clearly understand that the marital status as Married are in more number both in defaulters and non-defaulters.

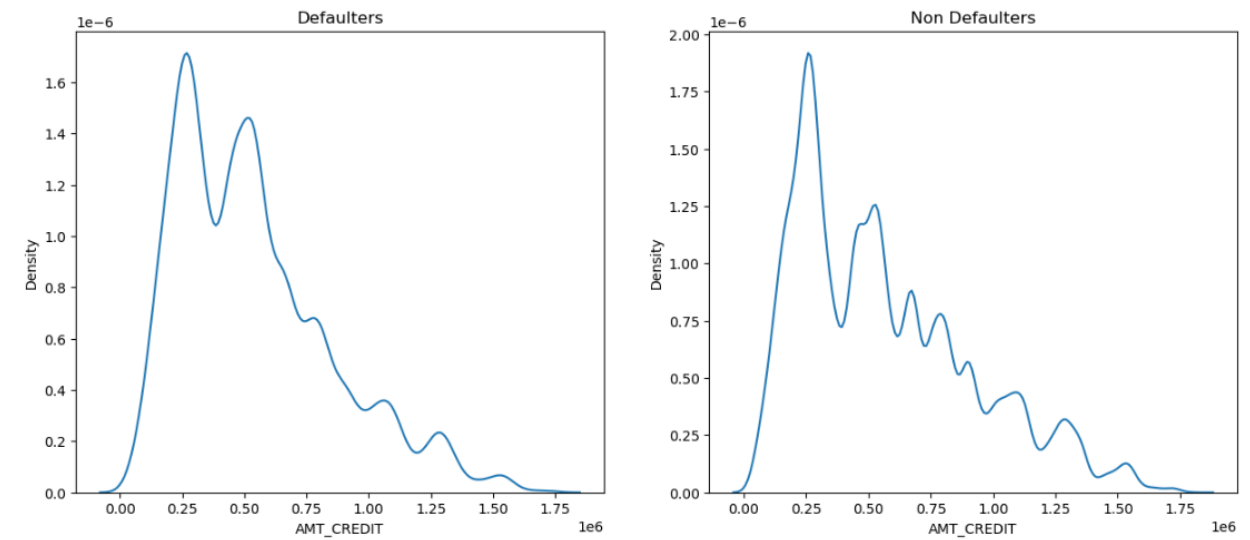
### Count of defaulters and non-defaulters on the basis OCCUPATION\_TYPE



### Analysis

Laborers is having more number comparing with other category both in defaulters and non-defaulters

### Defaulters and non-defaulters on the basis of credit amount of the loan

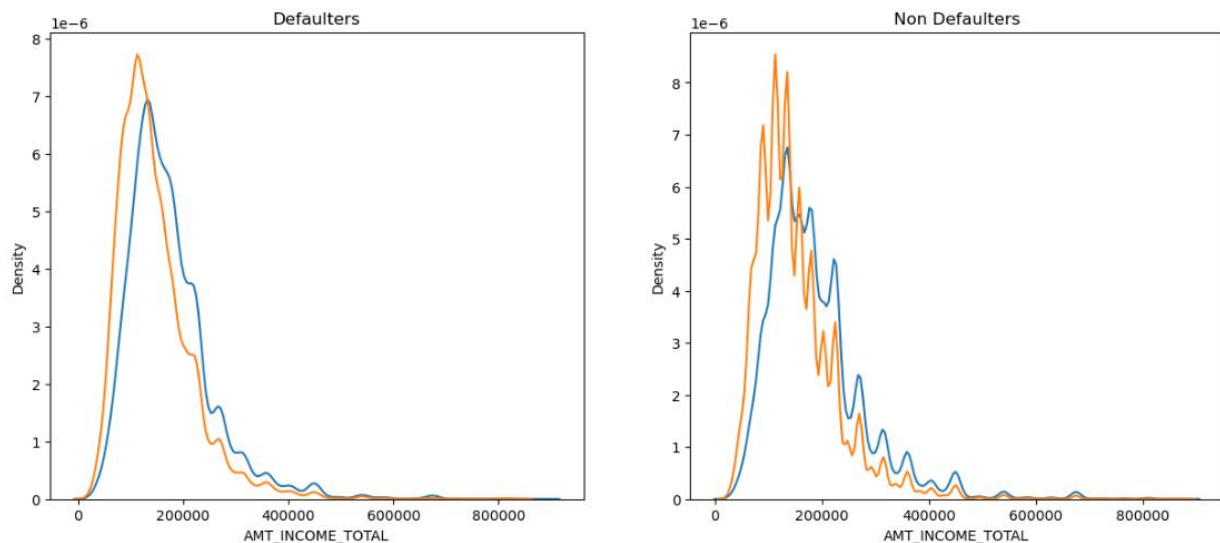


### Analysis

**Defaulters** - We can notice that the lesser the credit amount of the loan, the more chances of being defaulter. The spike is till 500000.

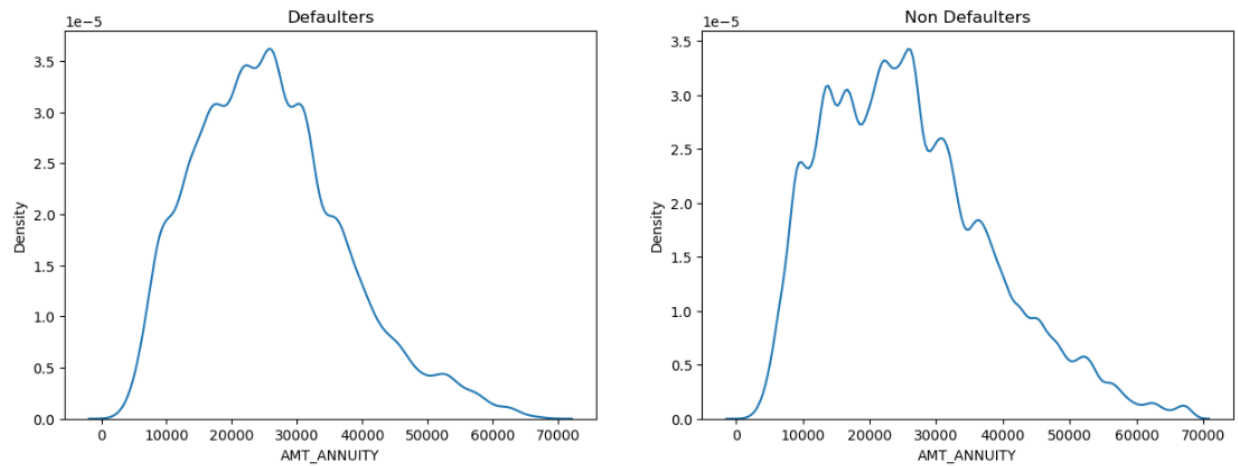
**Non defaulters** - If the credit amount is less, there is lesser chance of being defaulted. And gradually the chance is being decreased with the loan credit amount.

### *Defaulters and non-defaulters on the basis of gender and their total income*



1. **Defaulters** - We can notice by looking at the pattern that for being a defaulter both the genders (male and female) are almost equal in all income levels. The spike of being defaulters is from 50000 to 200000.
2. **Non defaulters** - Here we see an interesting pattern. Females are more non defaulter on the lower income level but lesser non defaluter in higher income level. The spike is more for both the genders from 75000 to 150000.

## Defaulters and non-defaulters on the basis of Loan annuity

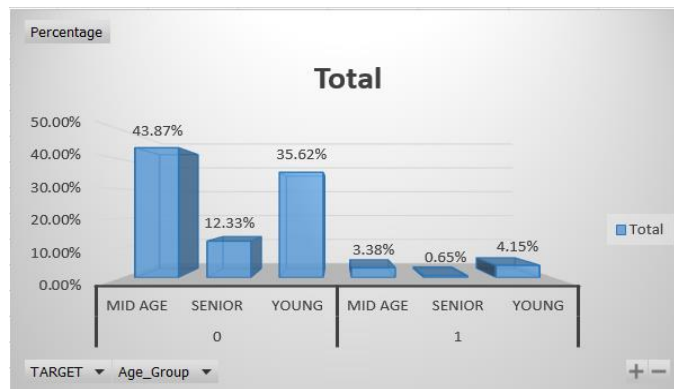


We can notice from the above distribution plot that in both the cases the loan annuity is concentrated more from 10000 to 40000.

## Bivariate Analysis:

### Age group & Target

Age_Group	Percentage
0	91.82%
Mid Age	43.87%
Senior	12.33%
Young	35.62%
1	8.18%
Mid Age	3.38%
Senior	0.65%
Young	4.15%
Grand Total	100.00%



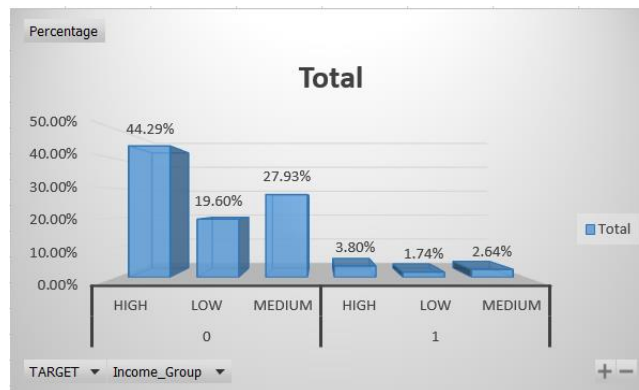
For defaulters, young age is given more preference compared to other groups.

For Non-defaulters, all age groups are given similar importance.

## Income group & Target:



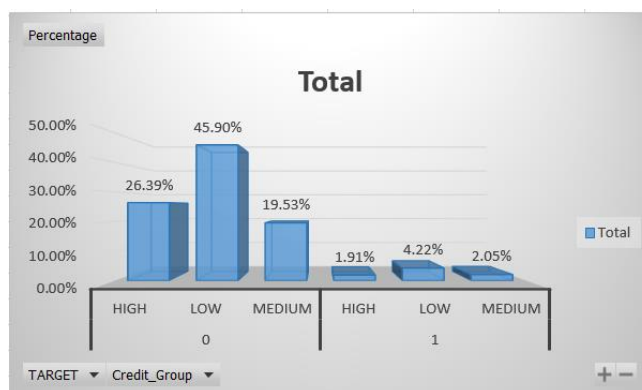
Income_Group	Percentage
0	91.82%
High	44.29%
Low	19.60%
Medium	27.93%
1	8.18%
High	3.80%
Low	1.74%
Medium	2.64%
<b>Grand Total</b>	<b>100.00%</b>



Low Income Group is given the least importance in both defaulters and non-defaulters.

#### Credit group & Target:

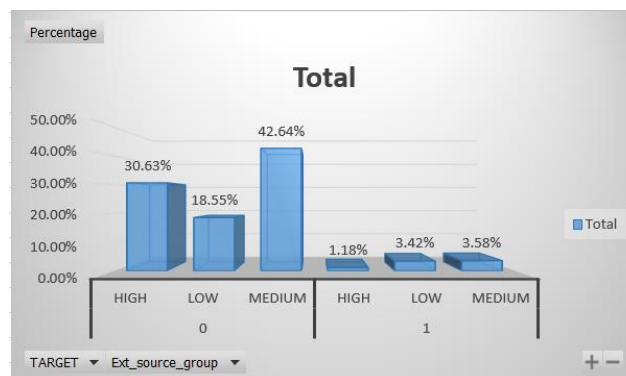
Credit_Groups	Percentage
0	91.82%
High	26.39%
Low	45.90%
Medium	19.53%
1	8.18%
High	1.91%
Low	4.22%
Medium	2.05%
<b>Grand Total</b>	<b>100.00%</b>



Credit groups having low value are more in number in both the defaulters and non-defaulters.

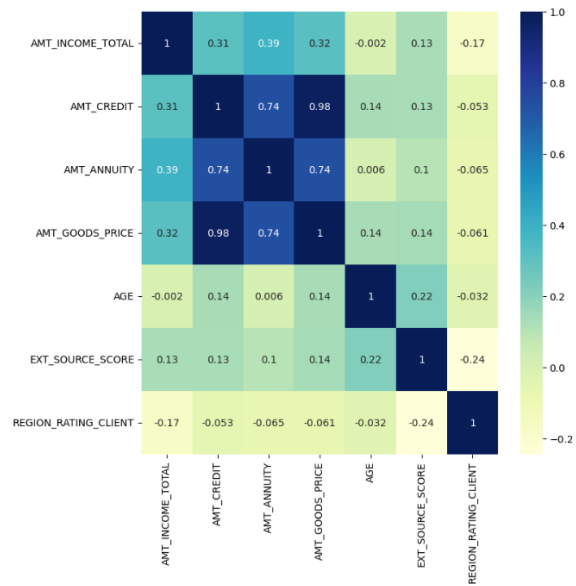
#### External Source Score group & Target:

Ext_score	Percentage
0	91.82%
High	30.63%
Low	18.55%
Medium	42.64%
1	8.18%
High	1.18%
Low	3.42%
Medium	3.58%
<b>Grand Total</b>	<b>100.00%</b>



Medium value dominates the range in both the categories.

#### Co-relation Analysis between top numerical columns:

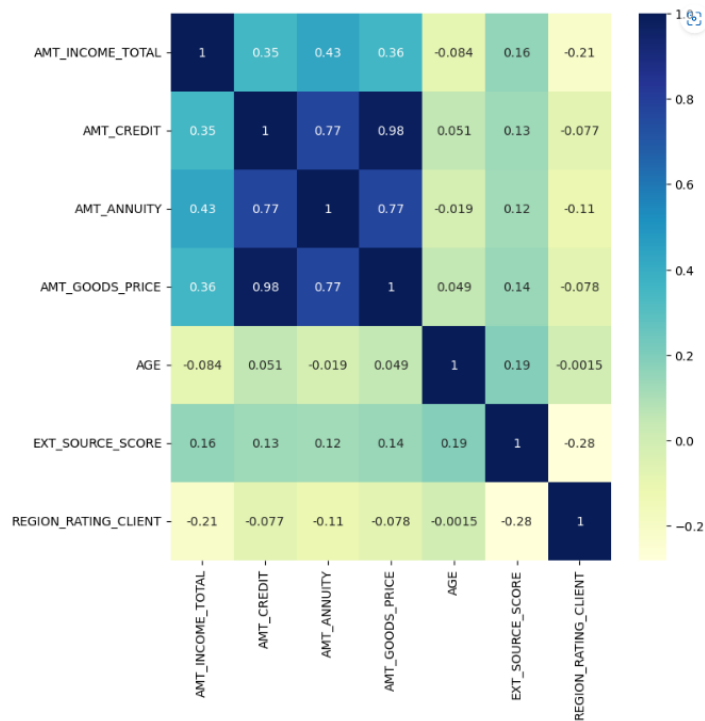


Highly co-relate columns for defaulters

AMT\_CREDIT and AMT\_ANNUITY (0.74)

AMT\_CREDIT and AMT\_GOODS\_PRICE (0.98)

AMT\_ANNUITY and AMT\_GOODS\_PRICE (0.74)



Highly correlate columns for non-defaulters

AMT\_CREDIT and AMT\_ANNUITY (0.76)

AMT\_CREDIT and AMT\_GOODS\_PRICE (0.98)

AMT\_ANNUITY and AMT\_GOODS\_PRICE (0.76)

Conclusion - We can see that for both defaulters and non-defaulters the same pairs of columns are highly correlated.

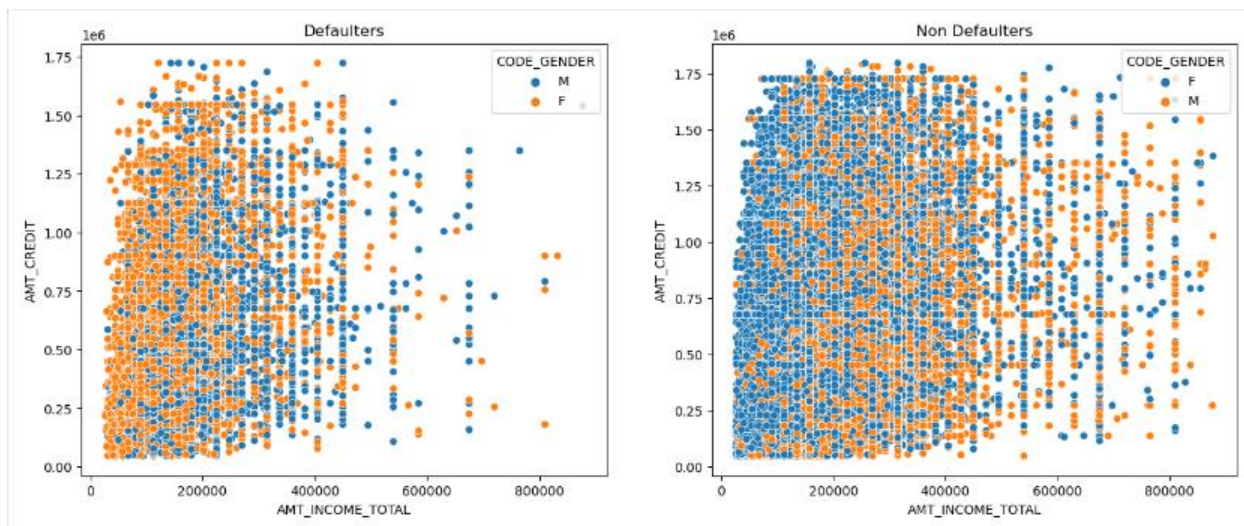
### ***Bivariate analysis on continuous variable***

#### ***Credit amount of the loan on the basis of client income for both male and female***

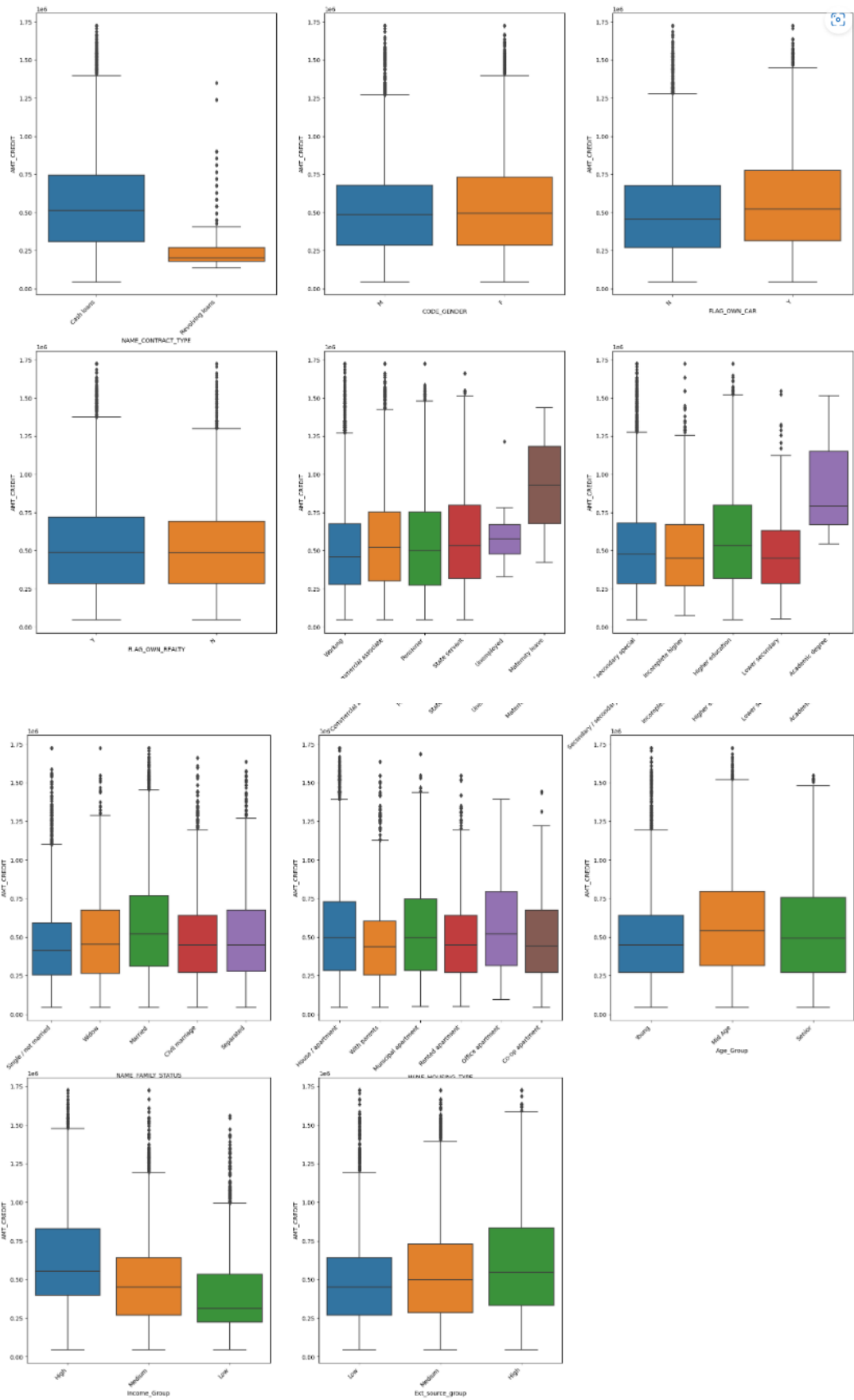
Analysis

Defaulters - We can slightly figure out that the values are more concentrated on the lower income and lower credit of the loan. That means as the income is increased, the amount of loan is also increased. This is true for both the genders.

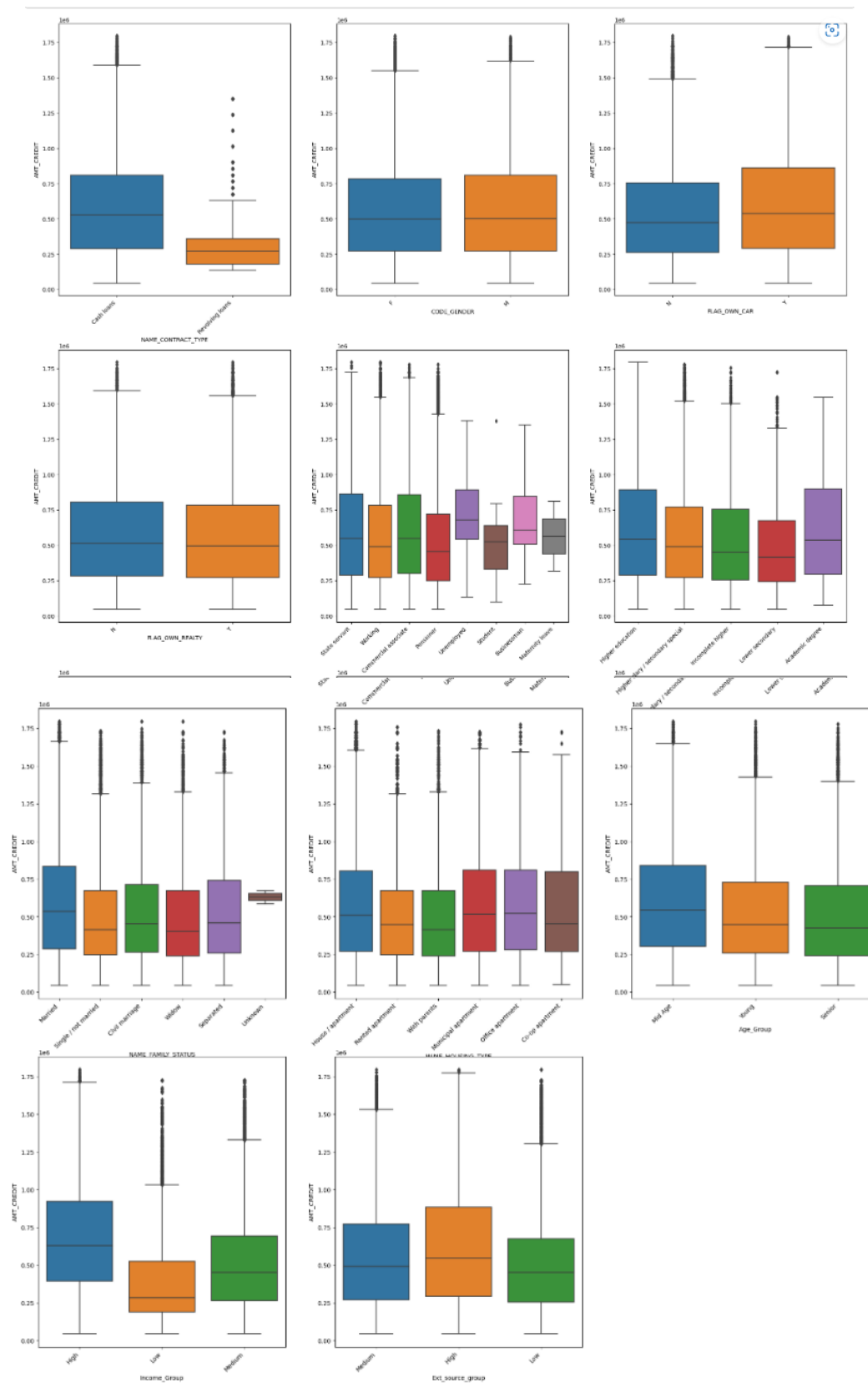
Non defaulters - We can hardly figure out any pattern out of this.



**Credit amount of the loan of various categories:**

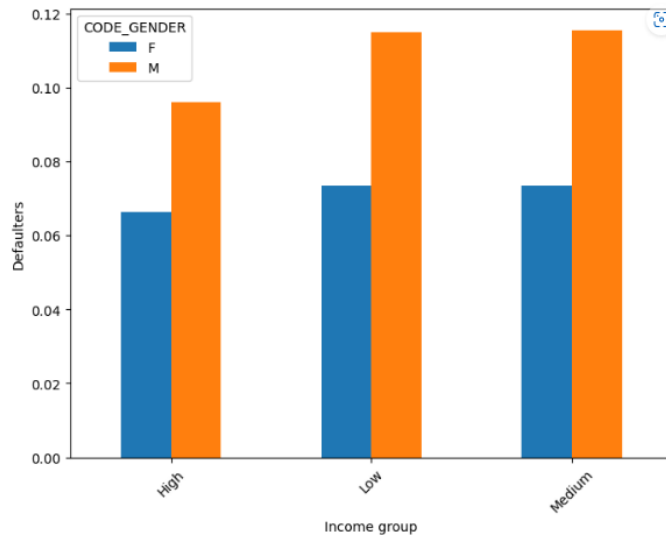


### Non-defaulters



## Analysis of two segmented variables

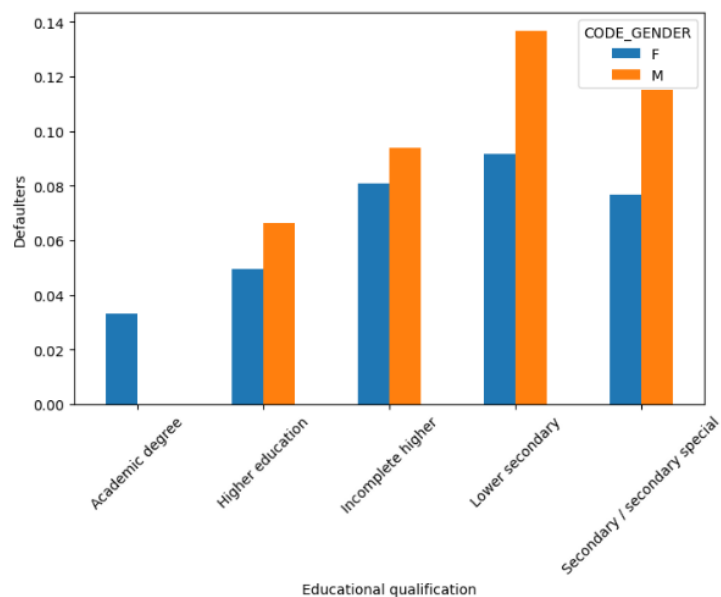
### Income group and gender



#### Analysis

We can see that Males are more likely defaulted than Females across all income groups

### Code Gender and Education Qualification



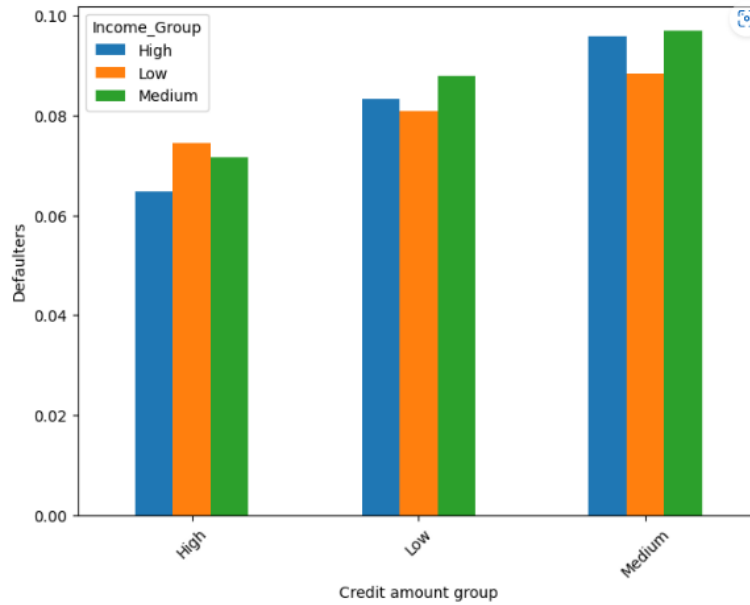
#### Analysis

Lower secondary educated clients are more defaulted followed by Secondary and Incomplete higher educated clients.

The Higher educated people are less defaulted.

Across all educated level Females are less defaulted than male.

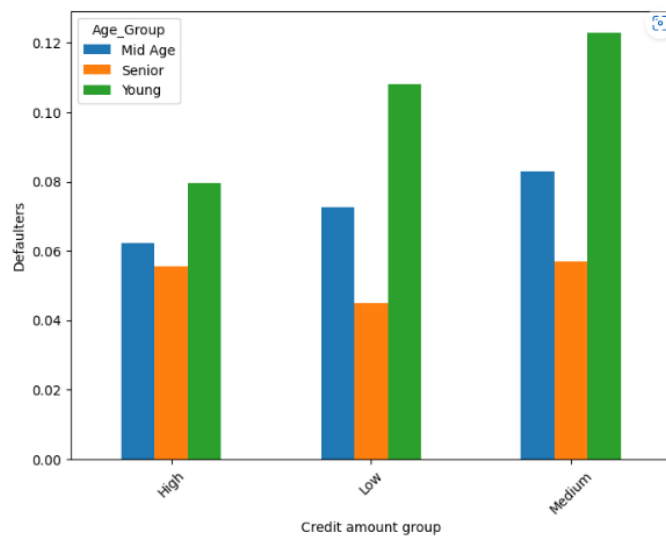
### Credit amount group and Income group



#### Analysis

Medium credit amount group are highly defaulted in all income groups.  
High credit amount groups are less likely to default in all income groups.

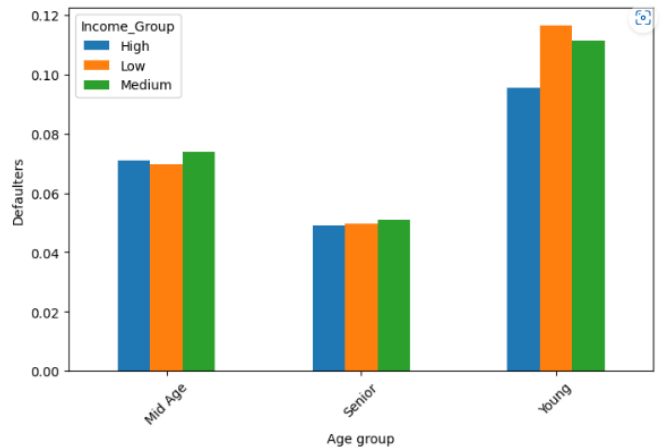
#### ***Credit amount group and Age group***



#### Analysis

Young clients with medium and low credit amount group are highly defaulted.  
Senior citizens across all credit amount groups are less likely defaulted.

#### ***Age group and Income group***



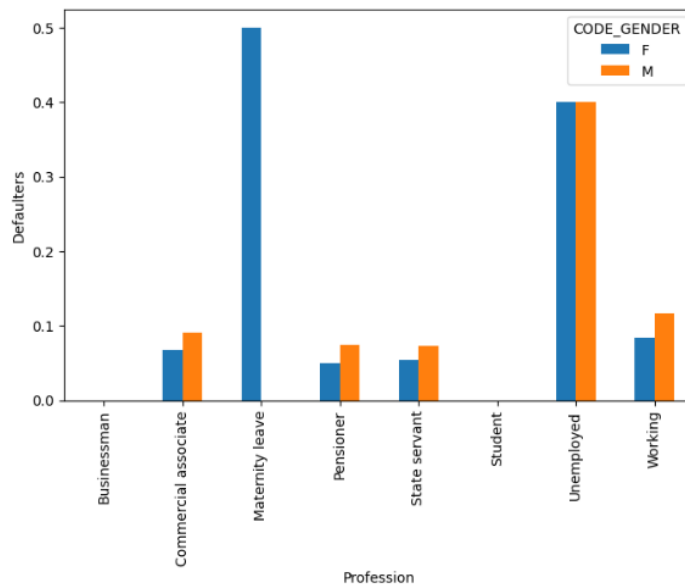
### Analysis

Young clients are more defaulted than Mid age and senior.

Young low income people are more defaulted.

For Mid age and senior people, the default rate is almost same in all income group.

### Profession and Gender



### Analysis

No surprise the unemployed clients are more defaulted.

Clients with maternity leave are expected to be defaulted more.

The default rate is lesser in all other professions.

Males are more defaulted with their respective professions compared to females.

### Analysis on Previous application:

Previous application dataset size (1670214, 37)

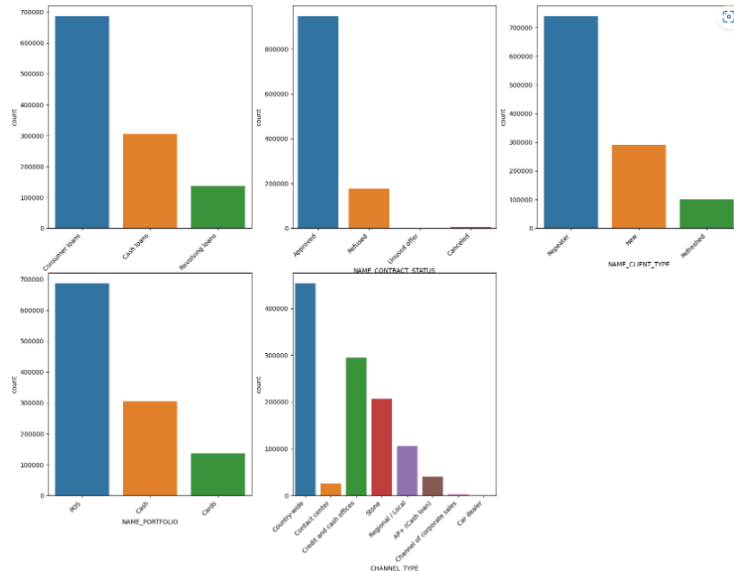
Similar to the application\_data we have also analyzed the previous application data

The outlier Analysis were also done: In our data cleaning exercise we have lost around 32% rows.

As this is huge dataset, we still have 1129387 records for analysis.



## Checking data imbalance



We can see that there is data imbalance in below columns: -

NAME\_CONTRACT\_TYPE - There are very few Revolving Loans

NAME\_CONTRACT\_STATUS - There are very few Refused loans. Almost negligible Canceled loans.

NAME\_CLIENT\_TYPE - There are very few New applicants. Even fewer Refreshed applicants.

NAME\_PORTFOLIO - Very few applications for Cards and Cars

CHANNEL\_TYPE - Except Country-Wide, Credit and Cash offices and Stone all other channels are very few in number

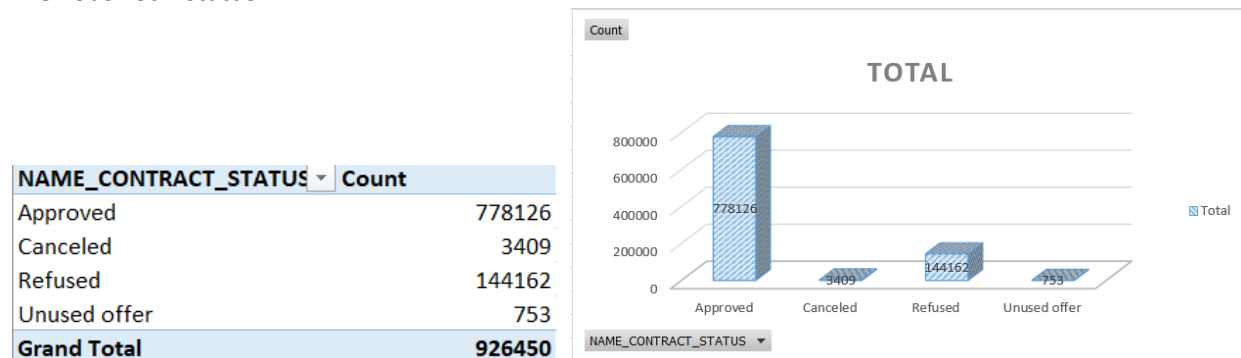
## Merging Current application and Previous application dataset

All the initial data exploration is done in the merged dataset

<https://drive.google.com/file/d/1YtliyTrEZNzTvBTLiKxQWYFq6mQqUBz/view?usp=sharing>

## Univariate analysis on unordered categorical variable:

Previous Loan status

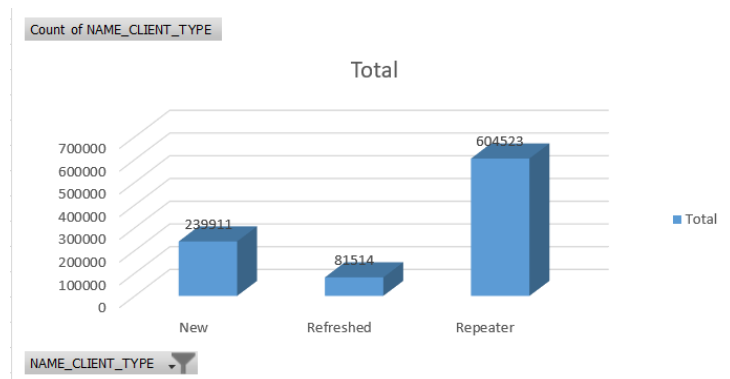


## Analysis

There are huge number of Approved loan than Refused. Hardly, there are any Canceled or Unused offer loan.

## Client type

NAME_CLIENT_TYPE	Count of NAME_CLIENT_TYPE
New	239911
Refreshed	81514
Repeater	604523
<b>Grand Total</b>	<b>925948</b>

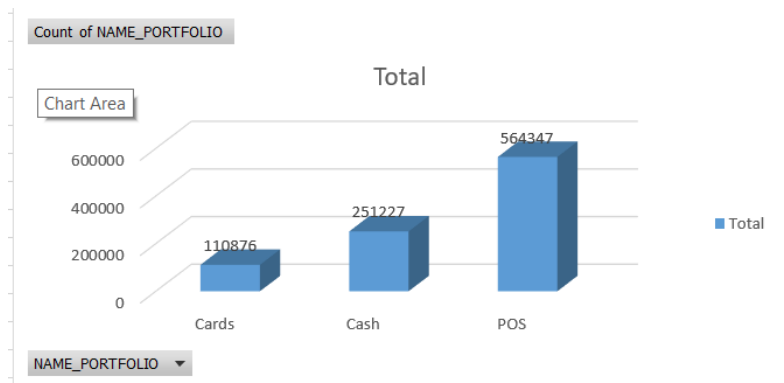


## Analysis

Mostly the applicants were Repeater

## Portfolio of the previous applications

NAME_PORTFOLIO	Count of NAME_PORTFOLIO
Cards	110876
Cash	251227
POS	564347
<b>Grand Total</b>	<b>926450</b>

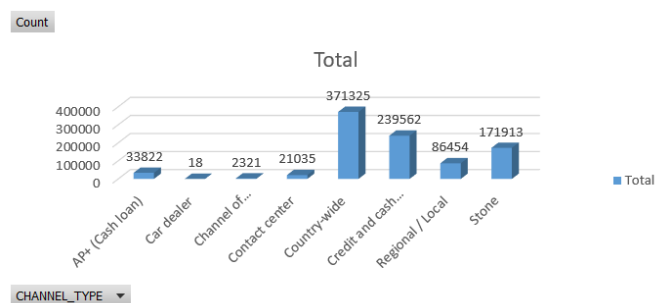


## Analysis

The highest number of the previous applications was for POS. Applications for Cash also has good number. Applications for Cards were very few.

## Application channel type

CHANNEL_TYPE	Count
AP+ (Cash loan)	33822
Car dealer	18
Channel of corporate sa	2321
Contact center	21035
Country-wide	371325
Credit and cash offices	239562
Regional / Local	86454
Stone	171913
<b>Grand Total</b>	<b>926450</b>

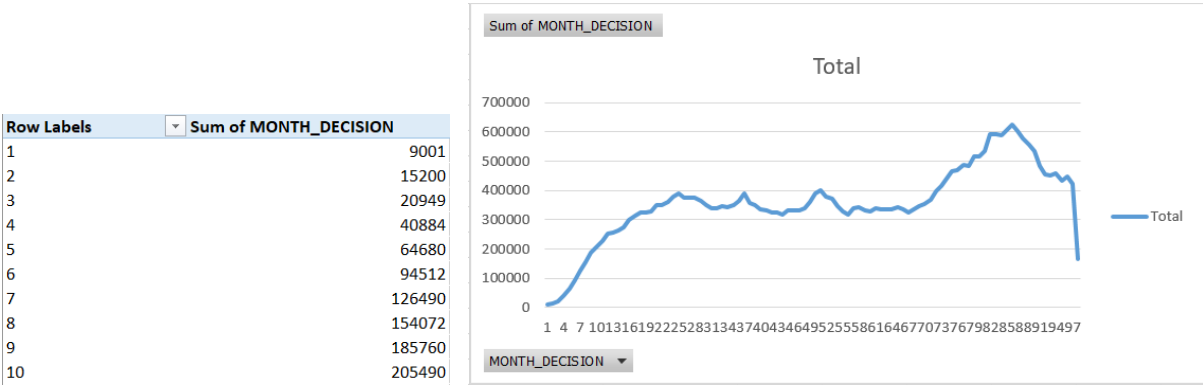


## Analysis

We see that Country-wide was heavily used for previous applications followed by Credit and Cash offices, Stone and Regional. Rest other channels are hardly used.

Univariate analysis for continious variables

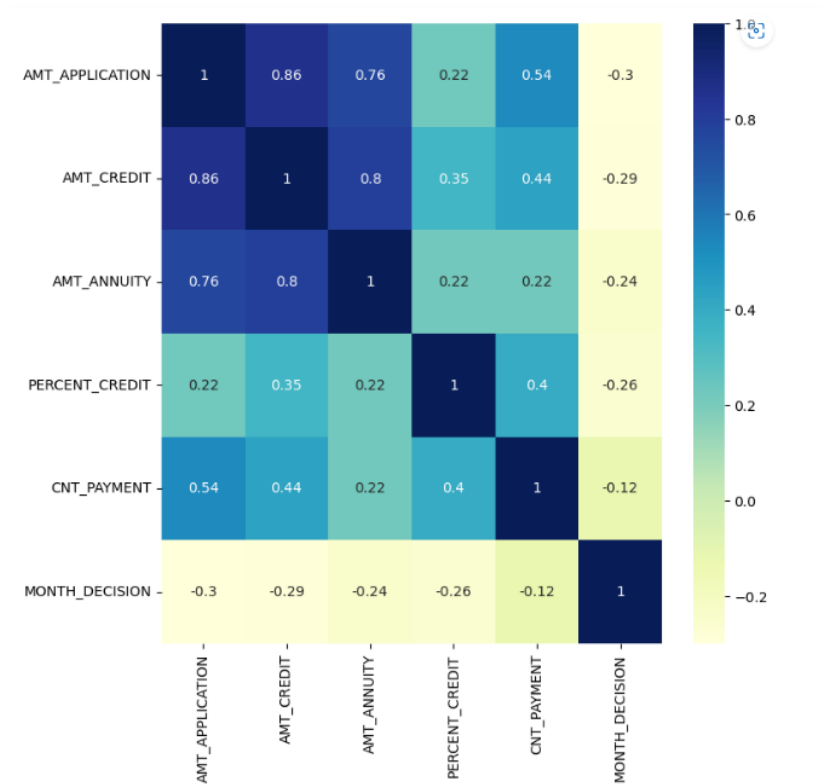
Months took for the pervious application decision relative to the current application



Analysis

We can see that most of the applications decision took approximately 30 months. The time taken spread up to 100 months

Correlation of relevant numerical columns

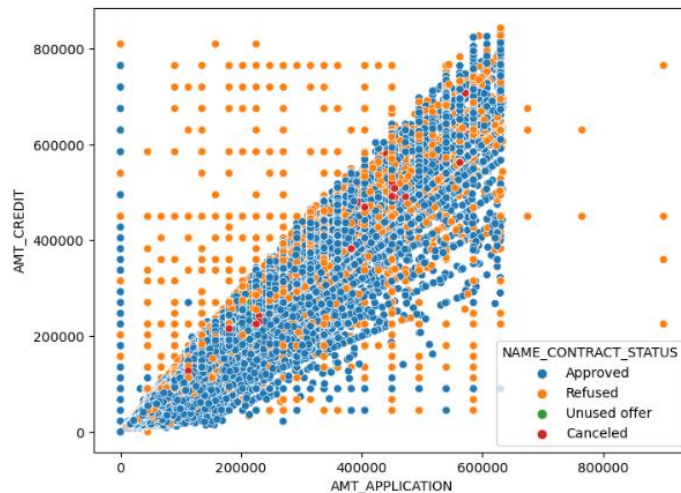


Highly correlate columns

AMT\_APPLICATION and AMT\_CREDIT  
AMT\_APPLICATION and AMT\_ANNUITY  
AMT\_CREDIT and AMT\_ANNUITY  
**Moderately correlated columns**  
AMT\_APPLICATION and CNT\_PAYMENT  
AMT\_CREDIT and CNT\_PAYMENT

### Bivariate analysis on continuous variable

#### Application amount and credited amount

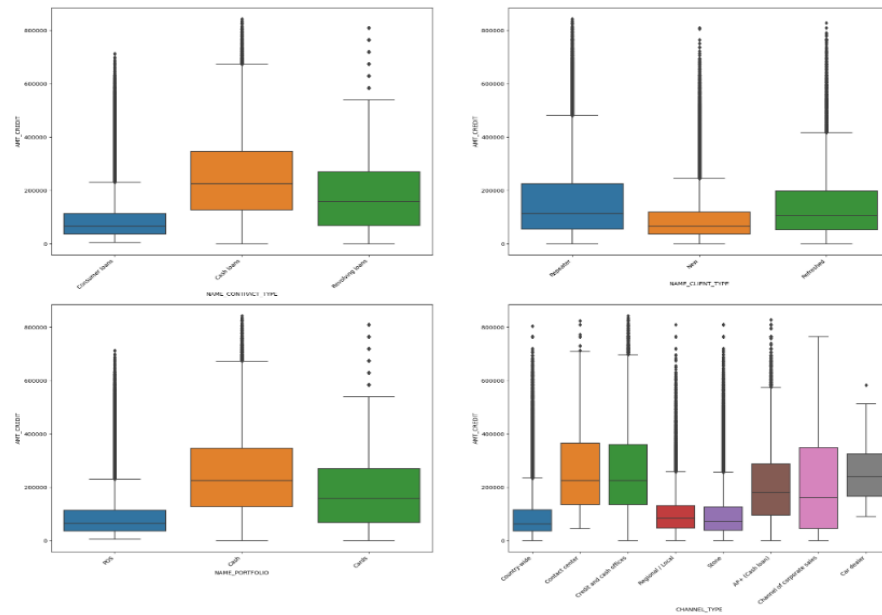


### Analysis

We can see a pattern here that the more the application amount of the loan, the lesser the months taken prior to current application. That means, most of the higher amount of the loan application decision made in the recent time compared to the lower loan amount application.

### Bivariate analysis on categorical variable

#### Credit amount of the loan of various categories



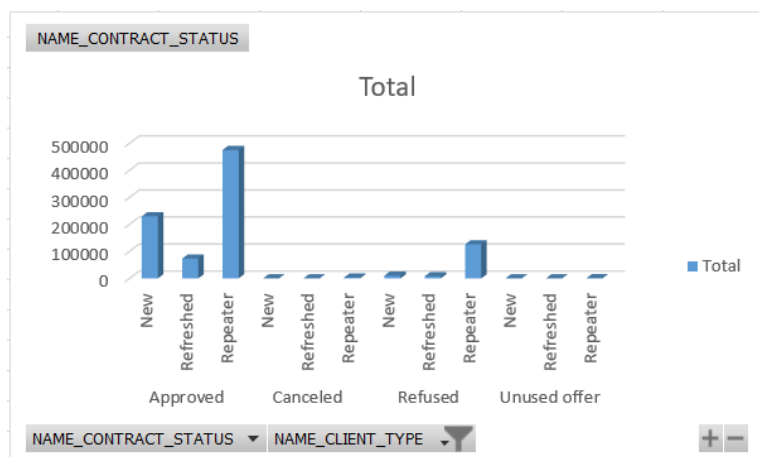
## Analysis

1. Cash loans are more credited in amount than Revolving and Consumer loans.
2. Repeater clients get more amount loan than new and refreshed clients.
3. The loan with portfolio Cars are more amount credited followed by Cash.
4. The credit amount of the loan is more from the application channel type as car dealer followed by Channel of corporate sales, Credit and cash offices and Contact center. The amount is very less for Regional, Stone and Country-wide channels.

## Analysis of two segmented variables

### Status and Client type:

Row Labels		NAME_CONTRACT_STATUS
Approved		777763
New		229541
Refreshed		73282
Repeater		474940
Canceled		3404
New		114
Refreshed		308
Repeater		2982
Refused		144029
New		10234
Refreshed		7842
Repeater		125953
Unused offer		752
New		22
Refreshed		82
Repeater		648
Grand Total		925948

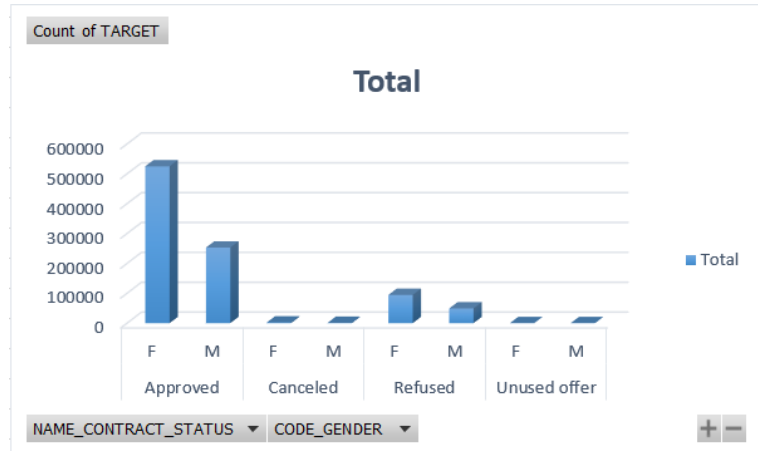


## Analysis

We see that the Repeater clients have more approved loans than New and Refreshed clients.

## Current loan defaulter status with respect to previous loan application status

Row Labels	Count of TARGET
<b>Approved</b>	<b>778126</b>
F	525008
M	253118
<b>Canceled</b>	<b>3409</b>
F	2399
M	1010
<b>Refused</b>	<b>144162</b>
F	94665
M	49497
<b>Unused offer</b>	<b>753</b>
F	449
M	304
<b>Grand Total</b>	<b>926450</b>

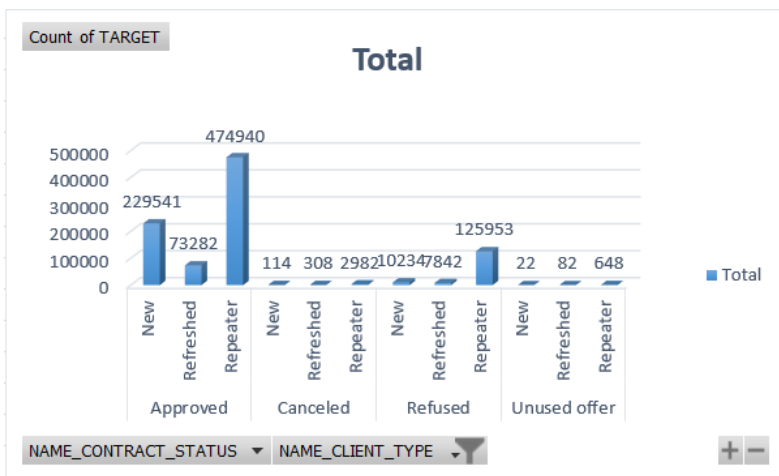


### Analysis

We see that previously Refused client is more defaulted than previously Approved clients. Also, in all the cases the Males are more defaulted than Females.

## Current loan defaulter status with respect to previous loan application status and client types

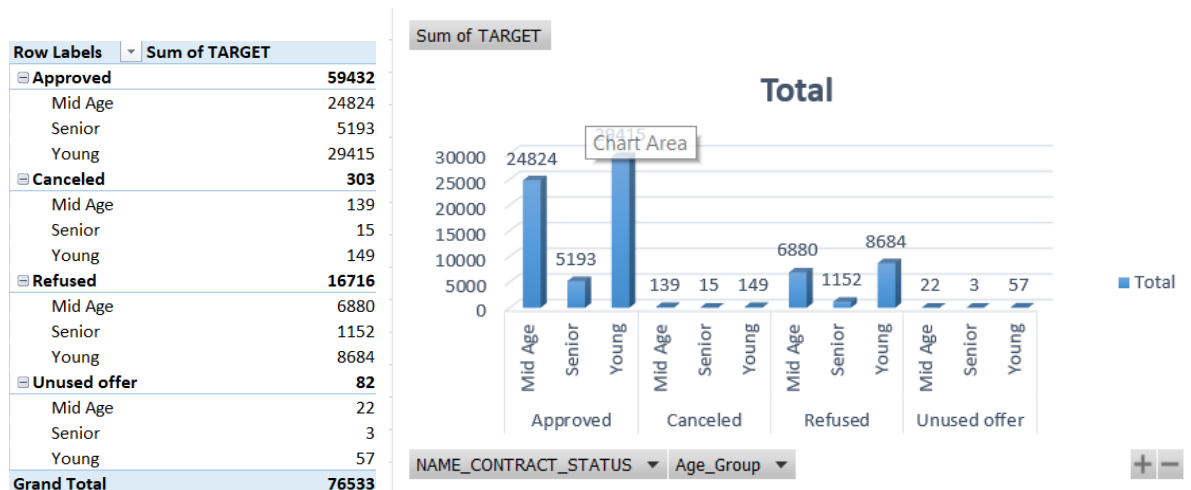
Row Labels	Count of TARGET
<b>Approved</b>	<b>777763</b>
New	229541
Refreshed	73282
Repeater	474940
<b>Canceled</b>	<b>3404</b>
New	114
Refreshed	308
Repeater	2982
<b>Refused</b>	<b>144029</b>
New	10234
Refreshed	7842
Repeater	125953
<b>Unused offer</b>	<b>752</b>
New	22
Refreshed	82
Repeater	648
<b>Grand Total</b>	<b>925948</b>



### Analysis

1. We can see that the Defaulters are more for previously Unused offers loan status clients, who were New.
2. For previously Approved status the New clients were more defaulted followed by Repeater.
3. For previously Refused applicants the Defaulters are more Refreshed clients.
4. For previously Canceled applicants the Defaulters are more New clients.

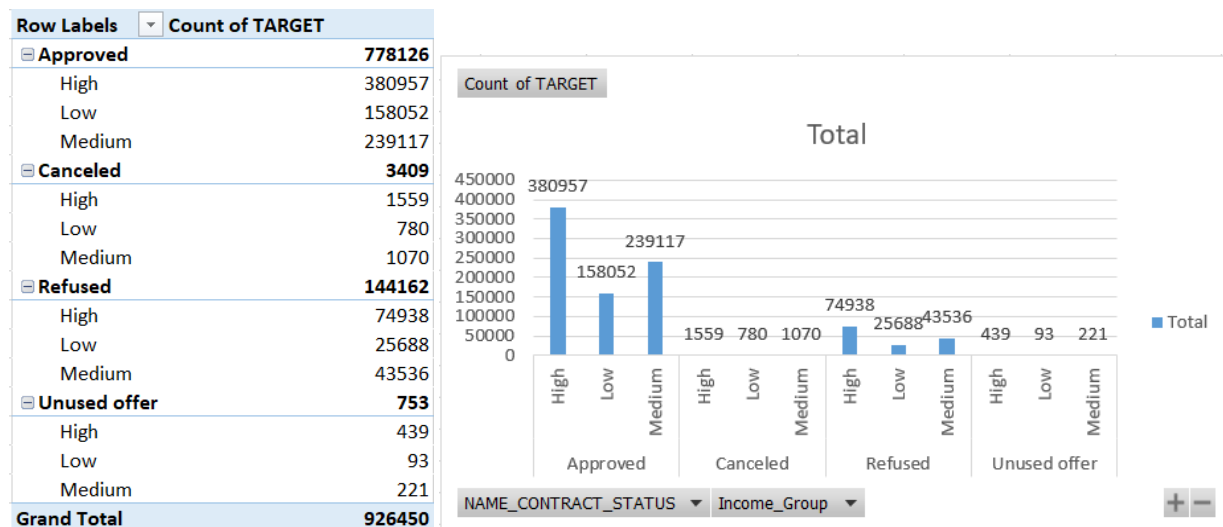
## Current loan defaulter status with respect to previous loan application status and age group



## Analysis

1. For all the previous status Young applicants are more defaulted.
2. For all the previous status Senior applicants are less defaulted compared to others.

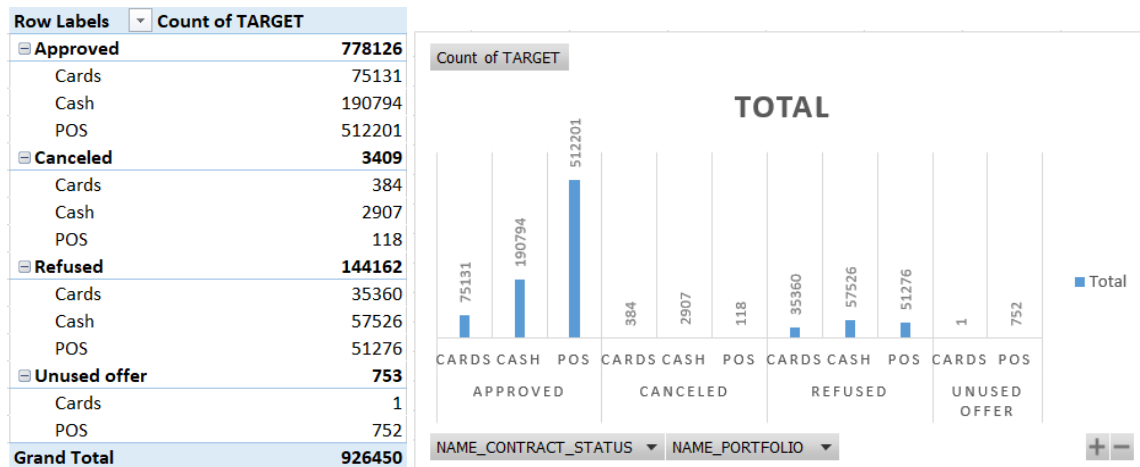
## Current loan defaulter status with respect to previous loan application status and income group



## Analysis

1. For previously Unused offer the Medium income group was more defaulted and Low income group is the least.
2. For other application status more or less all the income groups are equally defaulted.

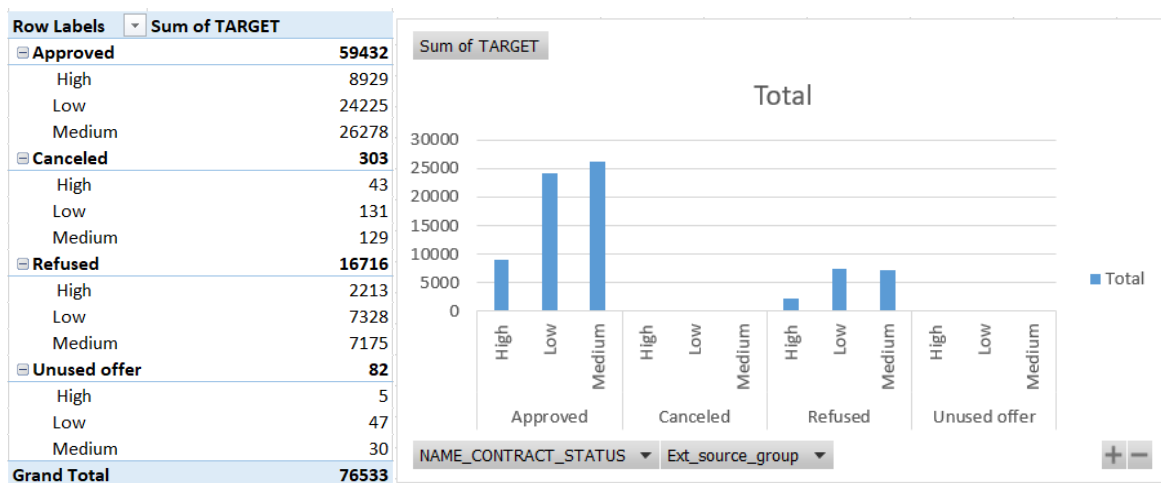
## Current loan defaulter status with respect to previous loan application status and portfolio of the loan



## Analysis

1. Most of the clients were defaulted, who previously applied loan for Cards.
2. For approved loan status the clients applied for Cars are less defaulted.
3. For Refused loan status the clients applied for POS are less defaulted.

## Current loan defaulter status with respect to previous loan application status and external source score category



1. Applicants with low external source score are highly defaulted.
2. Higher scorer applicants are very unlikely to default irrespective of their previous loan status.

**Result:** Working with such industry kind data set is very good experience. I have learned to do Exploratory Data Analysis on the real time Data set. Handling such a huge dataset was a task but I gradually learned to handle it better and use the necessary tools.

**Drive Link:**



**Jupyter**

**Notebook**

**file**

<https://drive.google.com/file/d/1nozOKMVvqjVV3zOCHxGLcX1mEcoCGlcU/view?usp=sharing>

**Univariate Analysis on Application Dataset**

<https://docs.google.com/spreadsheets/d/1uv9Tf9Y758DBw6OMpA5L4vcXGLXIURZB/edit?usp=sharing&oid=102292071000492049204&rtpof=true&sd=true>

**Univariate and Bivariate Analysis on merged dataset**

<https://docs.google.com/spreadsheets/d/1OSmqkiOMppKq7oVSK41hc5iw1CP6ON2O/edit?usp=sharing&oid=102292071000492049204&rtpof=true&sd=true>

**Whole File**

<https://docs.google.com/spreadsheets/d/1OSmqkiOMppKq7oVSK41hc5iw1CP6ON2O/edit?usp=sharing&oid=102292071000492049204&rtpof=true&sd=true>