



# CREDIT EDA CASE STUDY

By-

Anurag Sharma &  
Shubham Deshpande

# PROBLEM STATEMENT

- To find the possible patterns in the data so that a probability can be set up to find if a particular applicant will default or not and also a genuine applications is not rejected.



# HYPOTHESIS

- $H_0$  : The applicant will not default and will repay the loan on time.
- $H_1$  : The applicant will default and will not repay the loan on time.



# TYPES OF ANALYSIS DONE

- Handling of null values.
- Standardization of columns
- Univariate Analysis
- Bivariate Analysis
- Multivariate Analysis



# HANDLING NULL VALUES FOR CATEGORICAL AND NON-CONTINUOUS NUMERICAL VARIABLES.

The below mentioned columns are categorical/non-continuous numerical columns, we can use the mode value (most frequent value) for each ach column to impute the null values.

- 1) OCCUPATION\_TYPE
- 2) AMT\_REQ\_CREDIT\_BUREAU\_YEAR
- 3) AMT\_REQ\_CREDIT\_BUREAU\_MON
- 4) AMT\_REQ\_CREDIT\_BUREAU\_WEEK
- 5) AMT\_REQ\_CREDIT\_BUREAU\_DAY
- 6) AMT\_REQ\_CREDIT\_BUREAU\_HOUR
- 7) AMT\_REQ\_CREDIT\_BUREAU\_QRT
- 8) NAME\_TYPE\_SUITE

Below is a screenshot from Jupyter notebook denoting the values to be used for imputing null values.

```
--> For column OCCUPATION_TYPE, the mode value - Laborers can used to impute null values
--> For column AMT_REQ_CREDIT_BUREAU_YEAR, the mode value - 0.0 can used to impute null values
--> For column AMT_REQ_CREDIT_BUREAU_MON, the mode value - 0.0 can used to impute null values
--> For column AMT_REQ_CREDIT_BUREAU_WEEK, the mode value - 0.0 can used to impute null values
--> For column AMT_REQ_CREDIT_BUREAU_DAY, the mode value - 0.0 can used to impute null values
--> For column AMT_REQ_CREDIT_BUREAU_HOUR, the mode value - 0.0 can used to impute null values
--> For column AMT_REQ_CREDIT_BUREAU_QRT, the mode value - 0.0 can used to impute null values
--> For column NAME_TYPE_SUITE, the mode value - Unaccompanied can used to impute null values
```

# HANDLING NULL VALUES FOR NUMERICAL VARIABLES.

Below is the logic used for imputing null values:

- 1) If there are outliers, we'll use median value to impute the missing values.
- 2) if there are no outliers, we'll use the mean value to impute the missing values.

The below columns had outliers, 'median' values used to replace null columns:

- 1) AMT\_ANNUIITY
- 2) AMT\_GOODS\_PRICE
- 3) CNT\_FAM\_MEMBERS

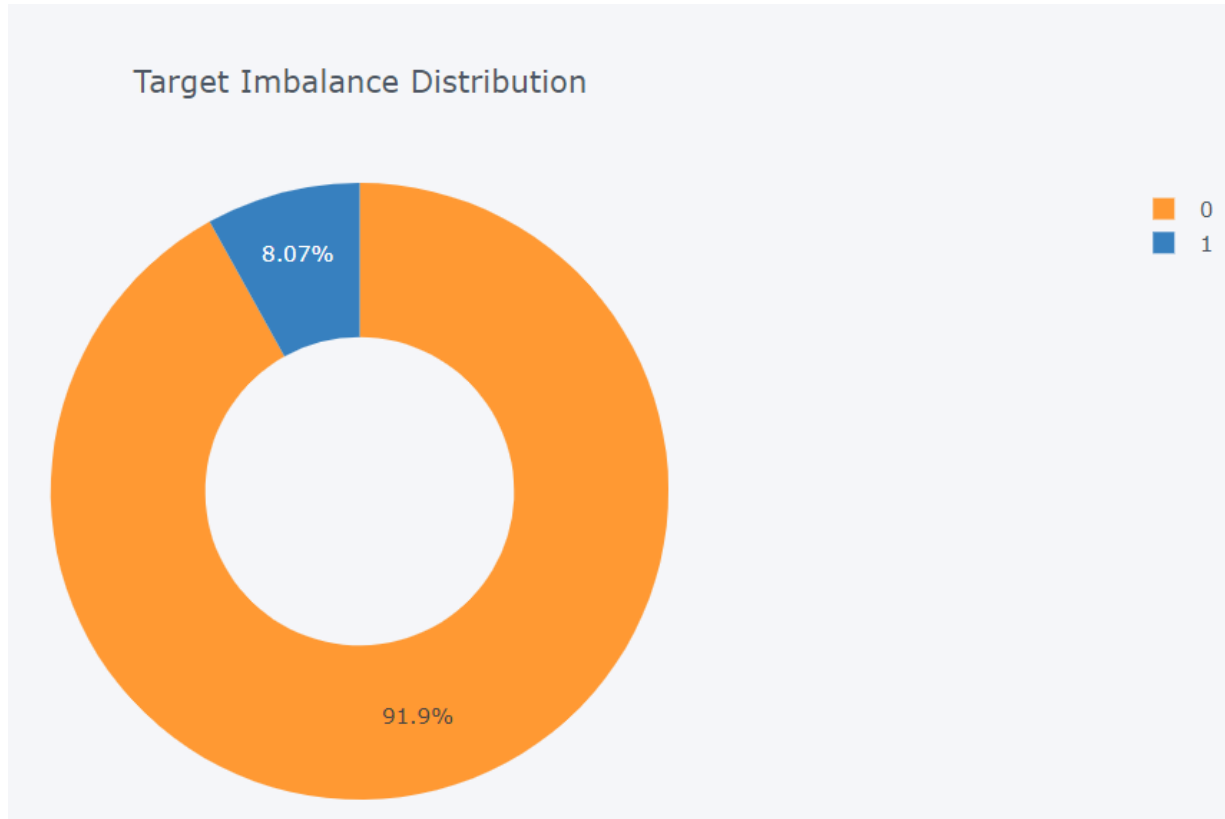
The below columns don't have any outliers, 'mean' values used to impute the missing values:

- 1) EXT\_SOURCE\_2
- 2) EXT\_SOURCE\_3

Below is the screenshot depicting the mean/median values that should be used for imputing null values.

```
--> For column AMT_ANNUIITY, the median value - 24903.0 can used to impute null values
--> For column AMT_GOODS_PRICE, the median value - 450000.0 can used to impute null values
--> For column CNT_FAM_MEMBERS, the median value - 2.0 can used to impute null values
--> For column EXT_SOURCE_2, the mean value - 0.5143926741308463 can used to impute null values
--> For column EXT_SOURCE_3, the mean value - 0.5108529061800121 can used to impute null values
```

# TARGET IMBALANCE



- From the pie chart above, we can clearly observe that imbalance is very high between the target variables i.e. applicants with no payment difficulties (~8%) are way more than applicants with payment difficulties (~92%).

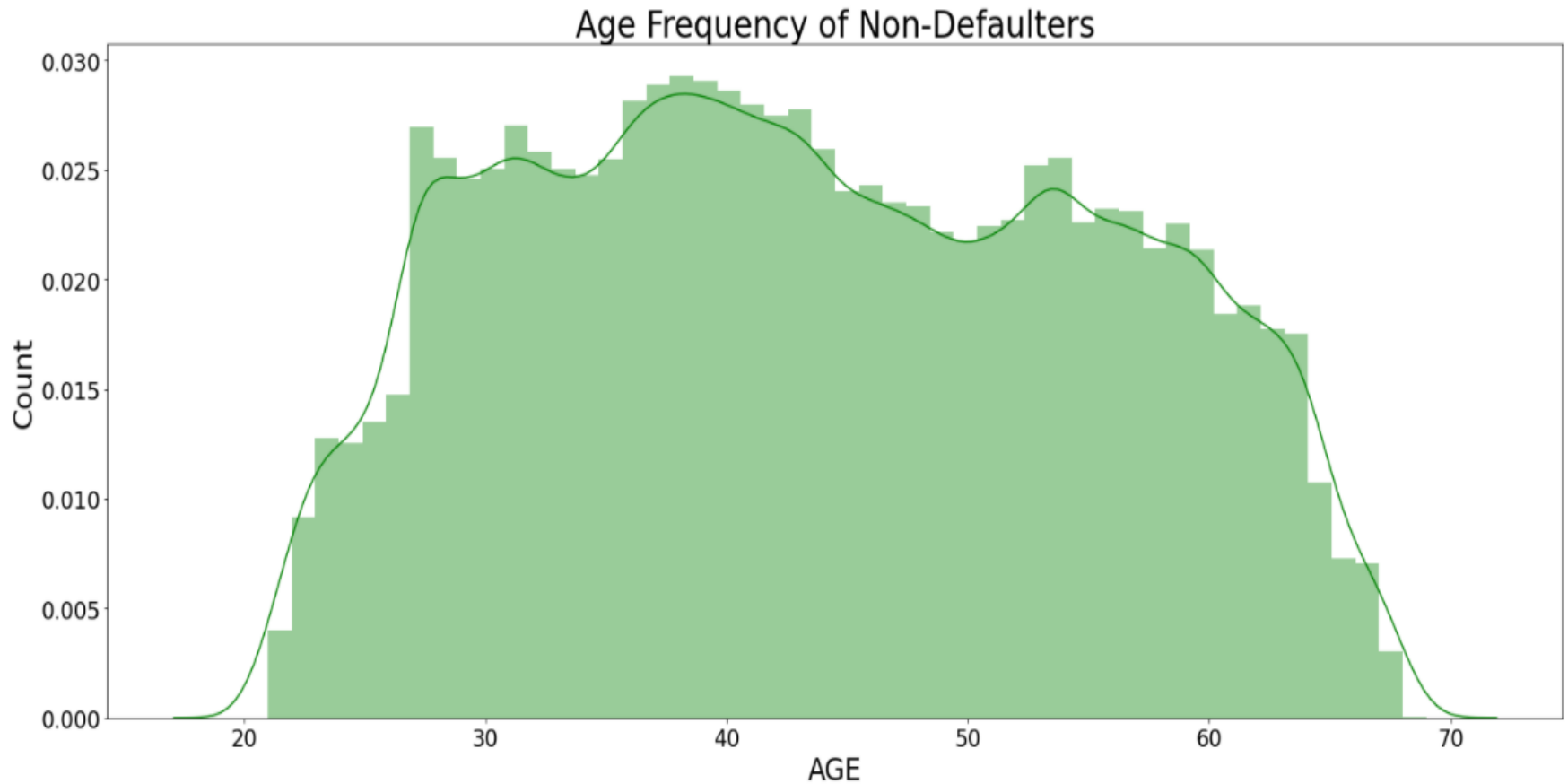


# UNIVARIATE ANALYSIS – NUMERIC DATA ON TARGET 0





# AGE

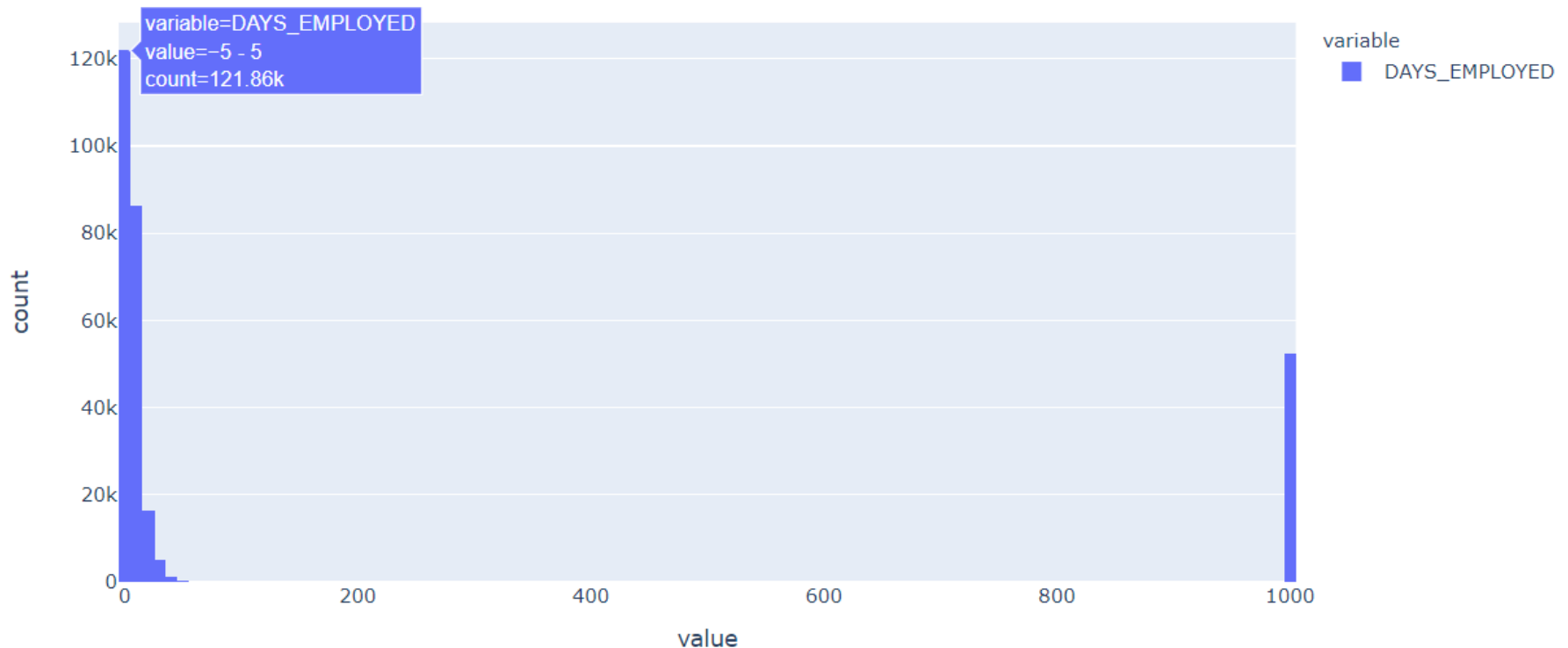


- Most of the applicants are middle aged. To be more specific, in middle aged they lie in the range of 35-45.



# NO OF DAYS EMPLOYED

No of Employed Days of Non-Defaulters



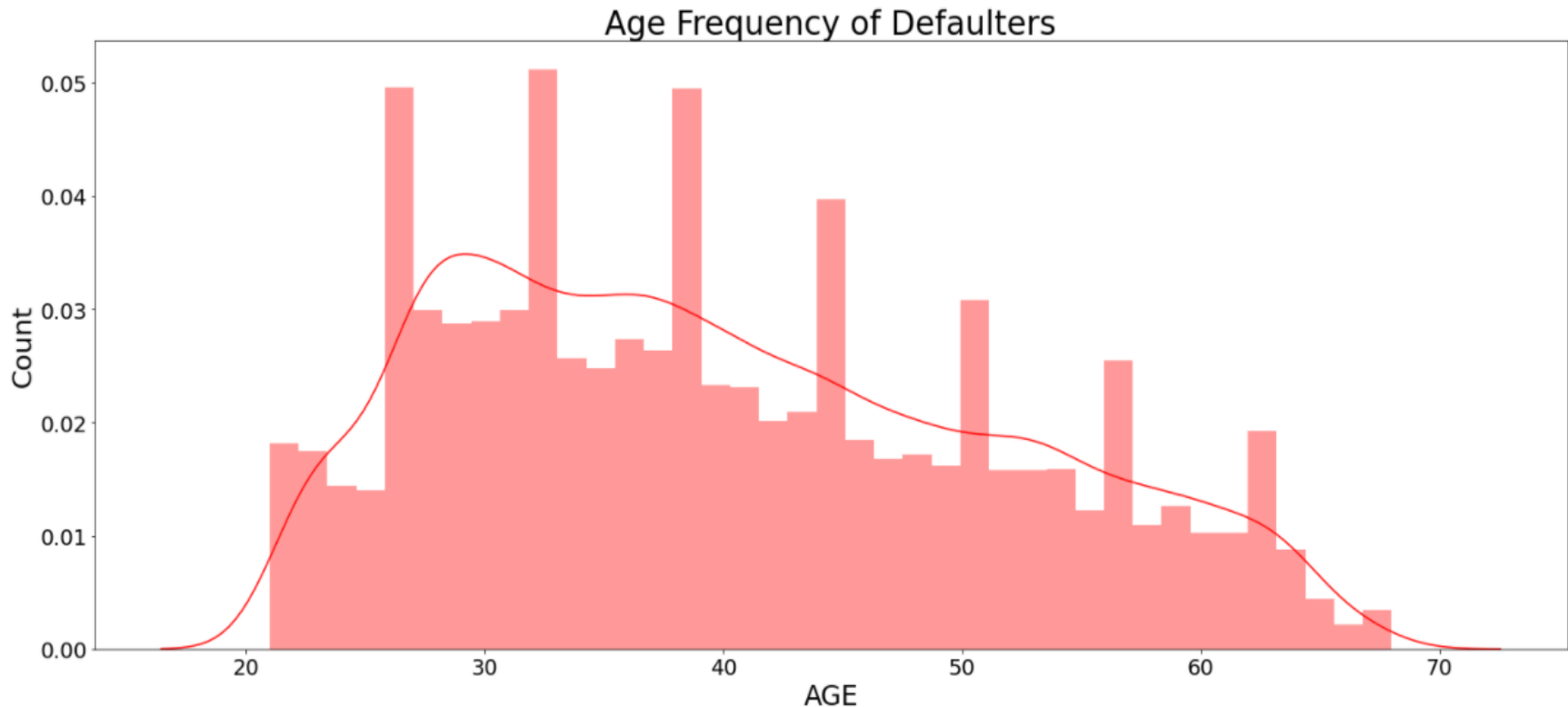
- People who are employed for a period of 5 years have highest probability that they will not default.



# UNIVARIATE ANALYSIS – NUMERIC DATA ON TARGET 1



# AGE

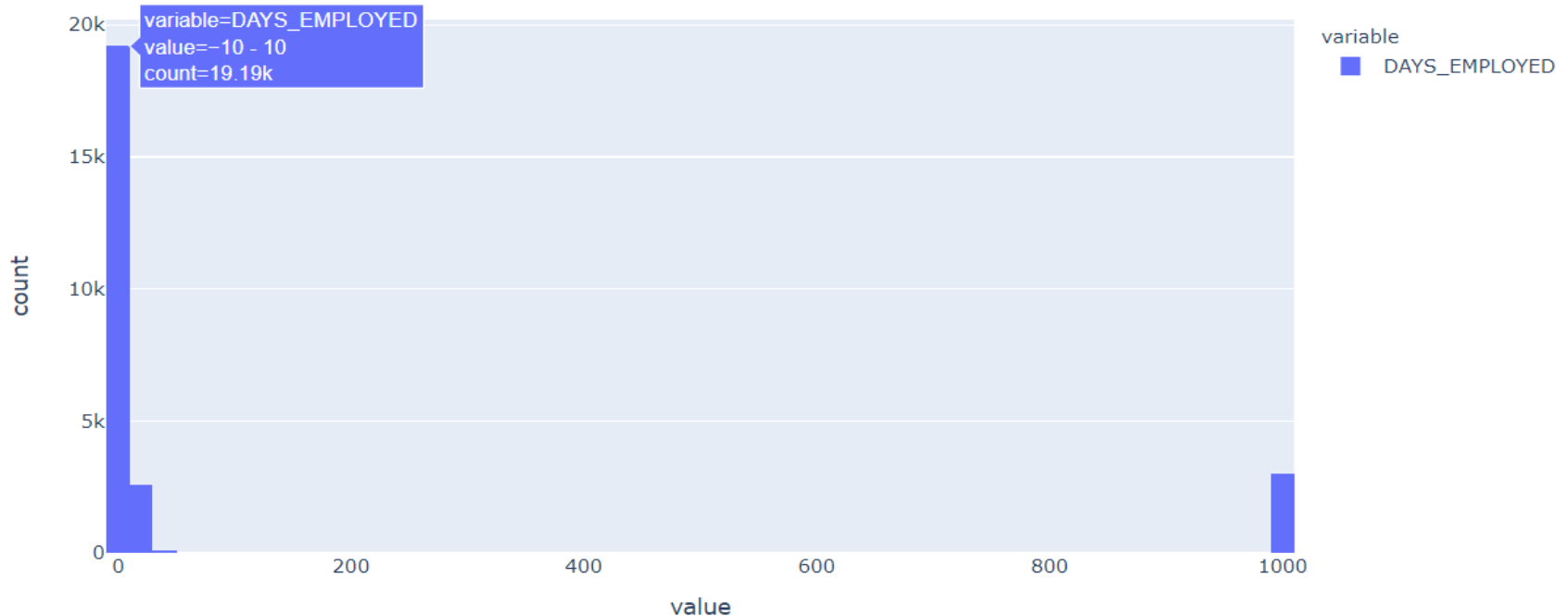


- It can be seen here that the highest number of defaulters lie in the age group of around 30.
- As the age increases, the possibility of the person defaulting decreases.



# NO OF DAYS EMPLOYED

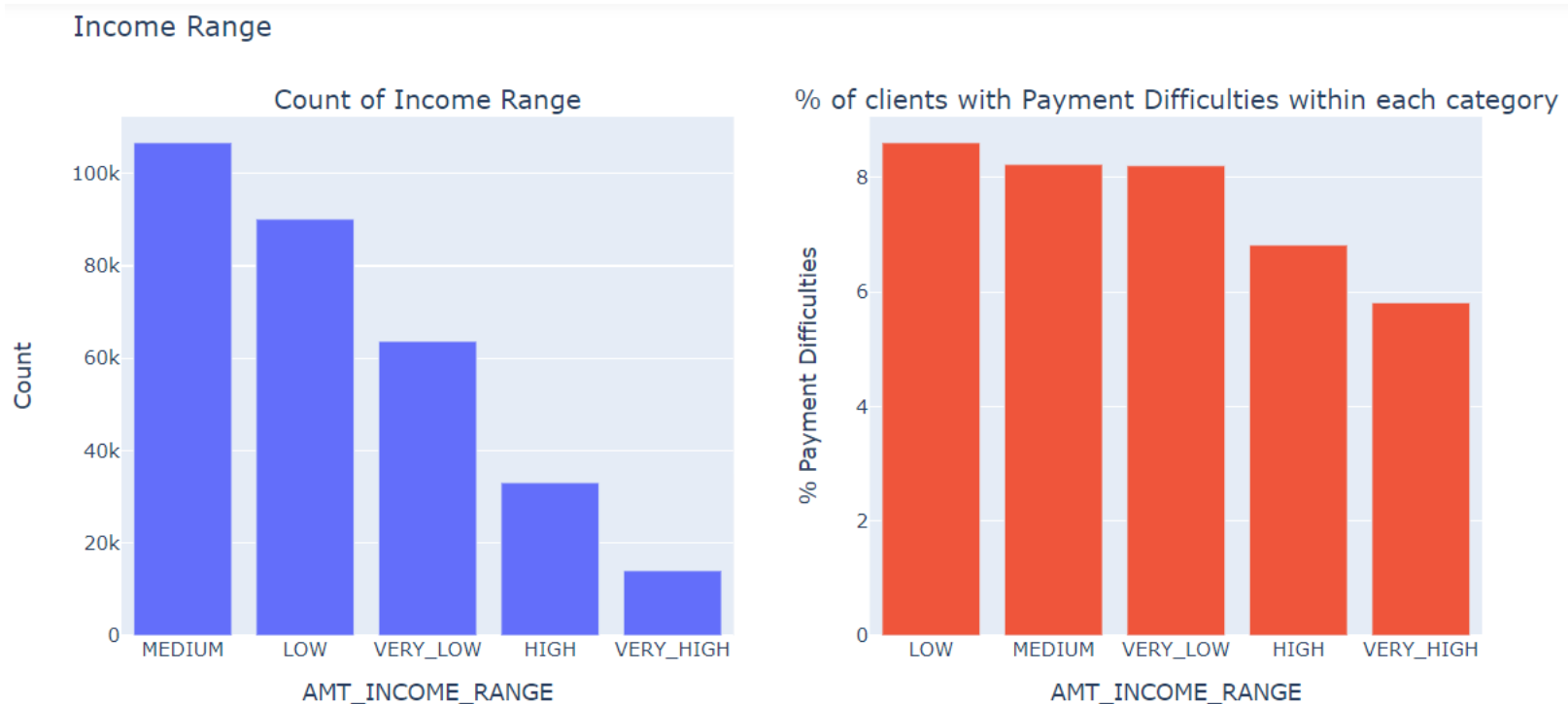
No of Employed Days of Defaulters



- People who are employed for a period of 10 years have highest probability that they will default.



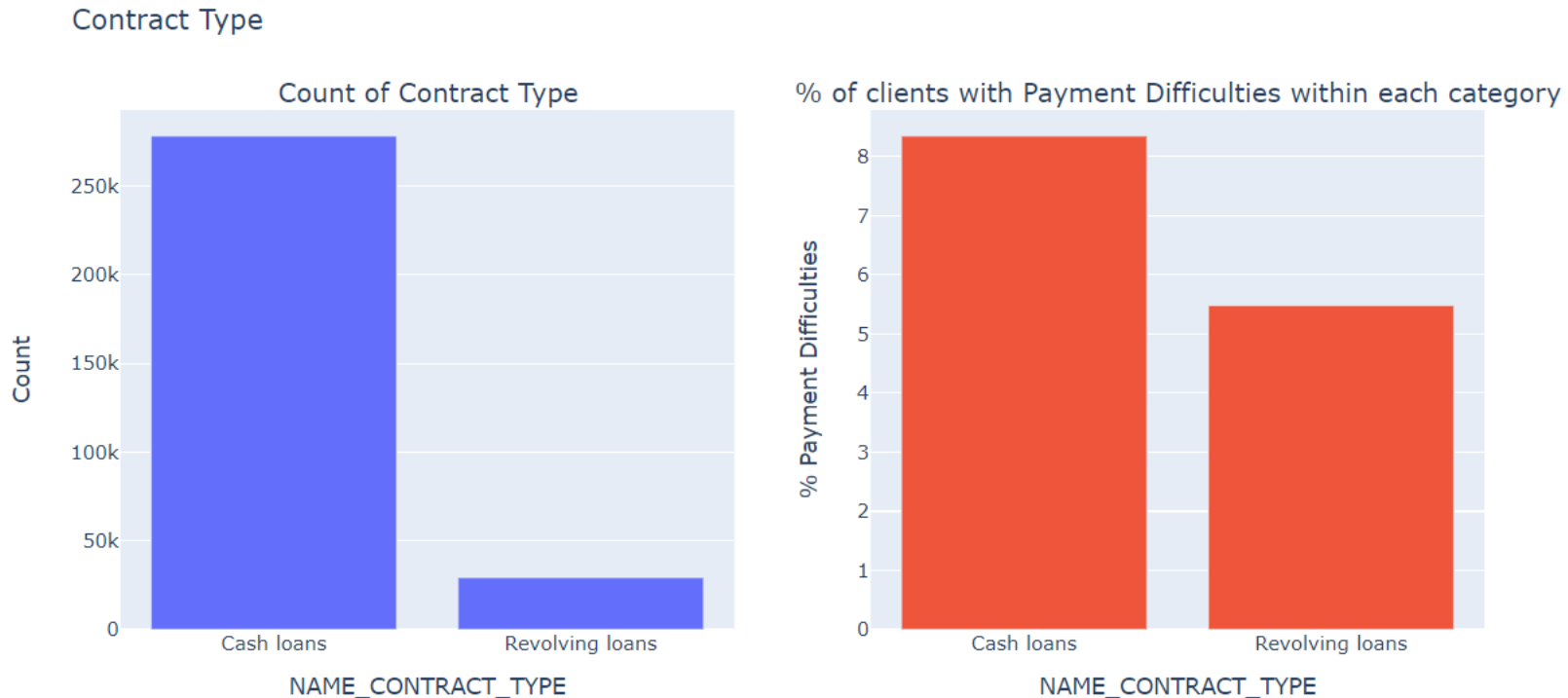
# DEFAULTERS % : INCOME RANGE



- Highest defaulters lie in the low income range followed by medium income range.



# DEFAULTERS % : CONTRACT TYPE

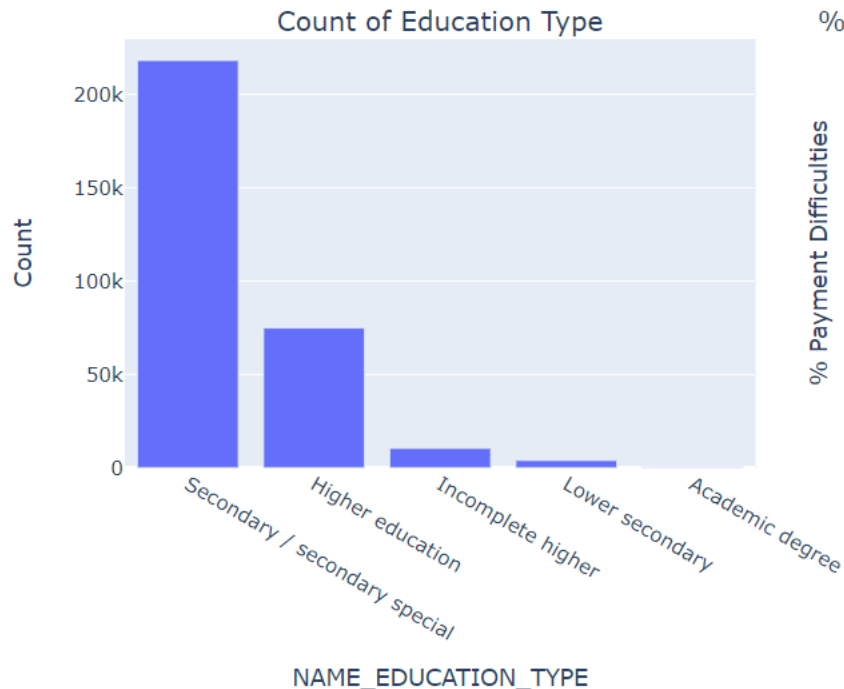


- Even though there are less applications of revolving loans yet the number of defaulters is higher in revolving loans when proportion is compared.

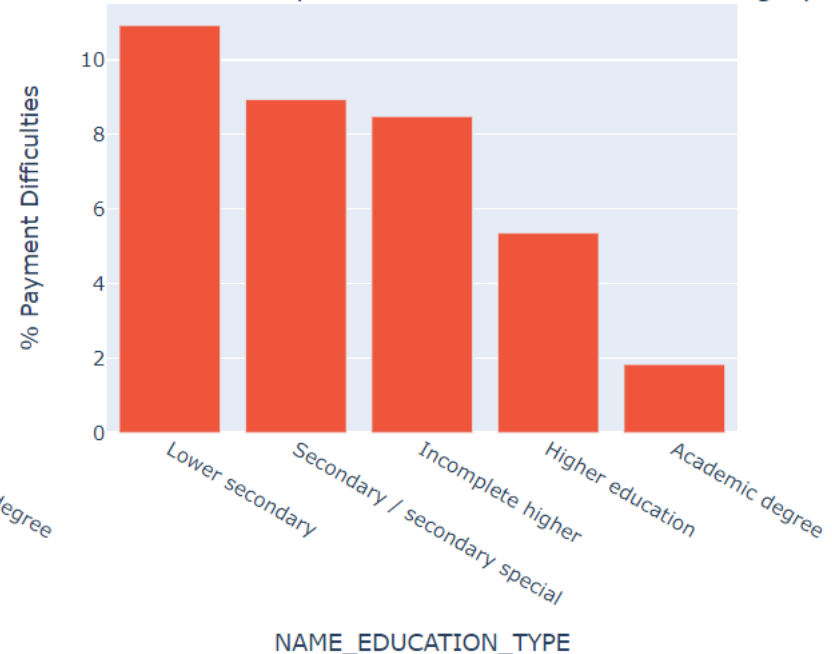


# DEFAULTERS % : EDUCATION TYPE

Education Type



% of clients with Payment Difficulties within each category



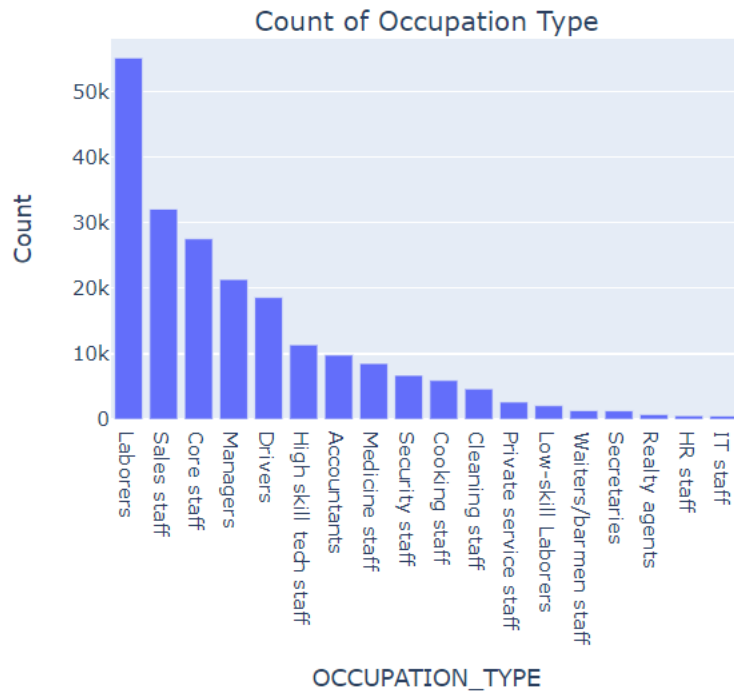
- Lower secondary, secondary and Incomplete higher have most of the defaulters in proportion to number of applications.



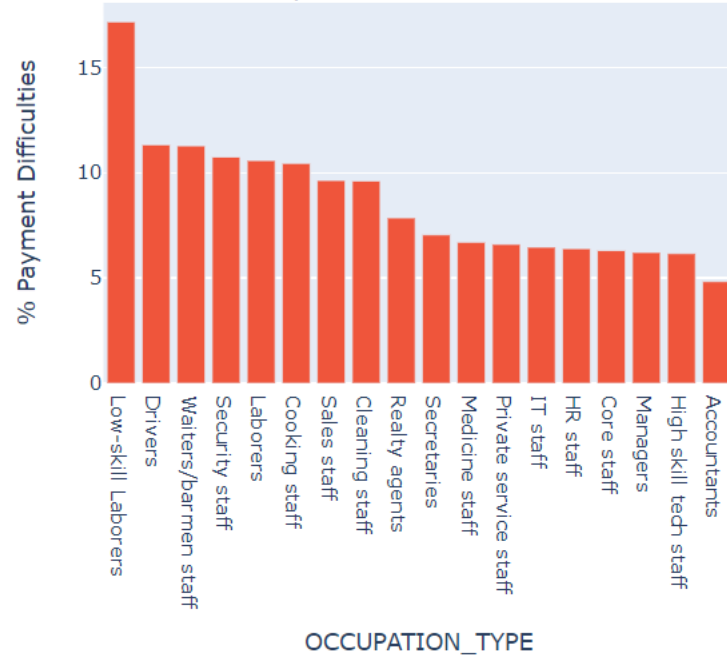


# DEFAULTERS % : OCCUPATION TYPE

Occupation Type



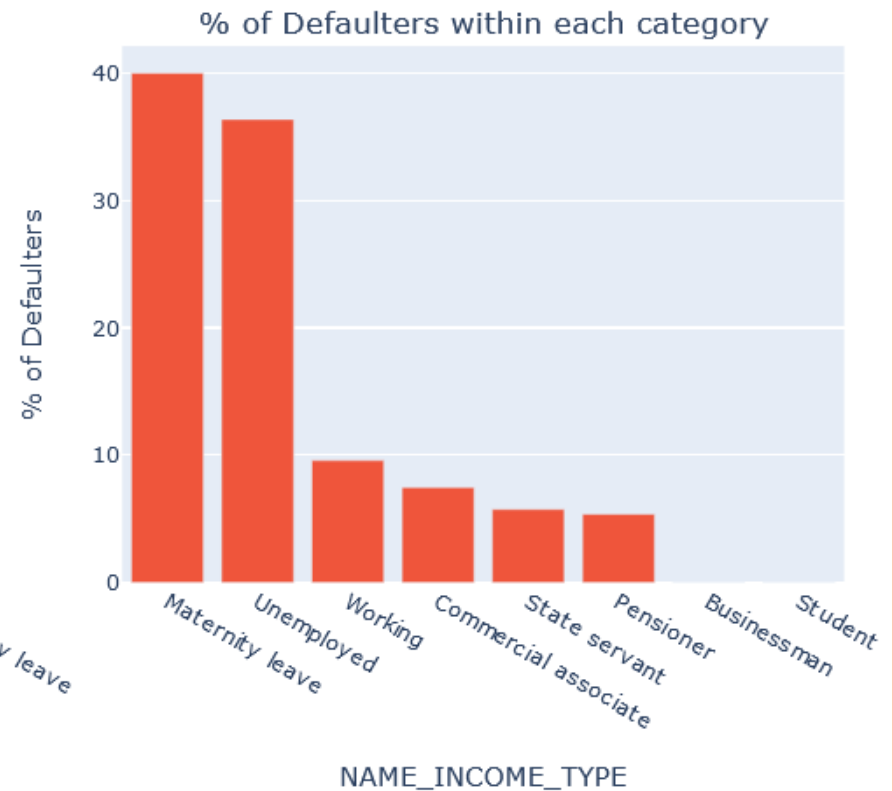
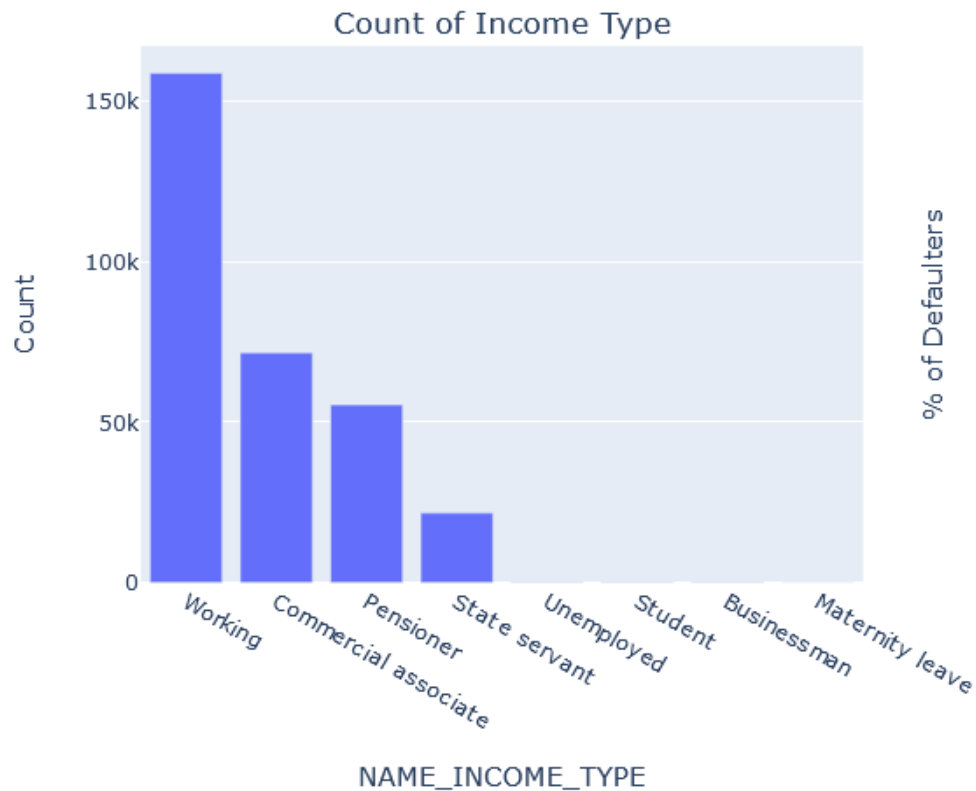
% of clients with Payment Difficulties within each category



- Low skill laborers are less when it comes to number of loan applications but they are the top most in terms of defaulters.



# DEFAULTERS % : INCOME TYPE



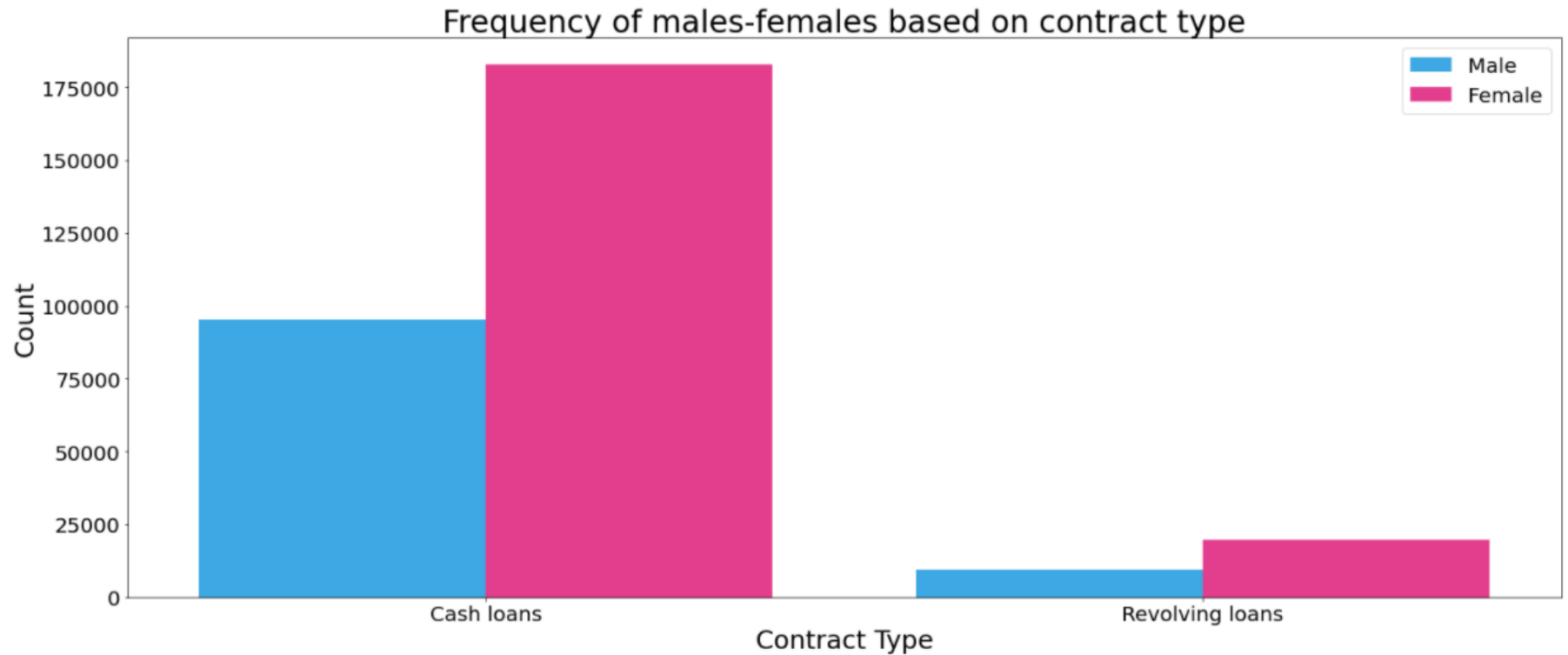
- From the plots above we can say that clients with 'Maternity leave' Income type have maximum % of payment difficulties followed by 'UNEMPLOYED' although the count of both these categories is far less than others.



# UNIVARIATE ANALYSIS – CATEGORICAL DATA



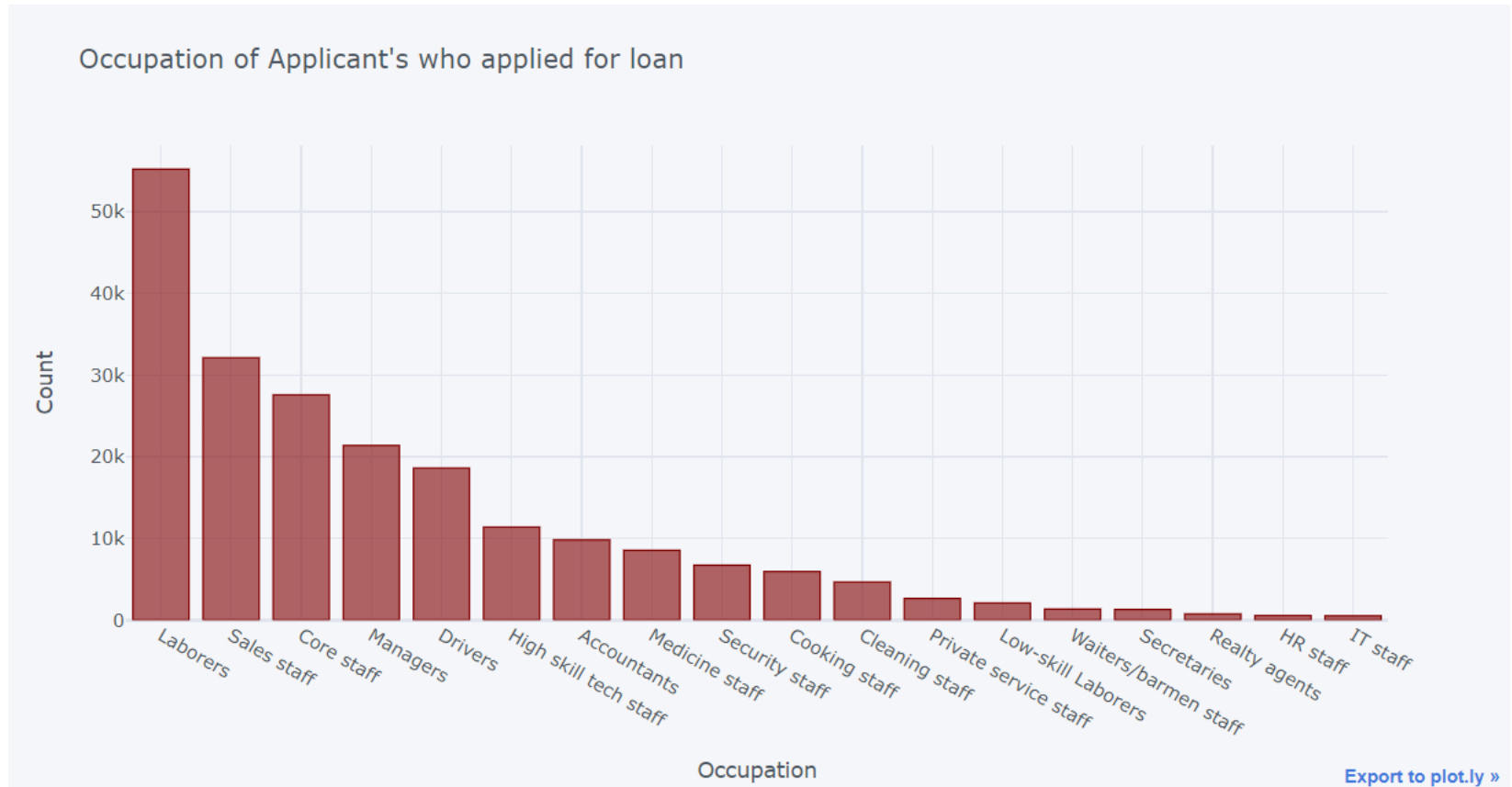
# FREQUENCY OF CONTRACT TYPE



- Cash loans is the highest amongst applications.
- Females are highest in both the contracts.



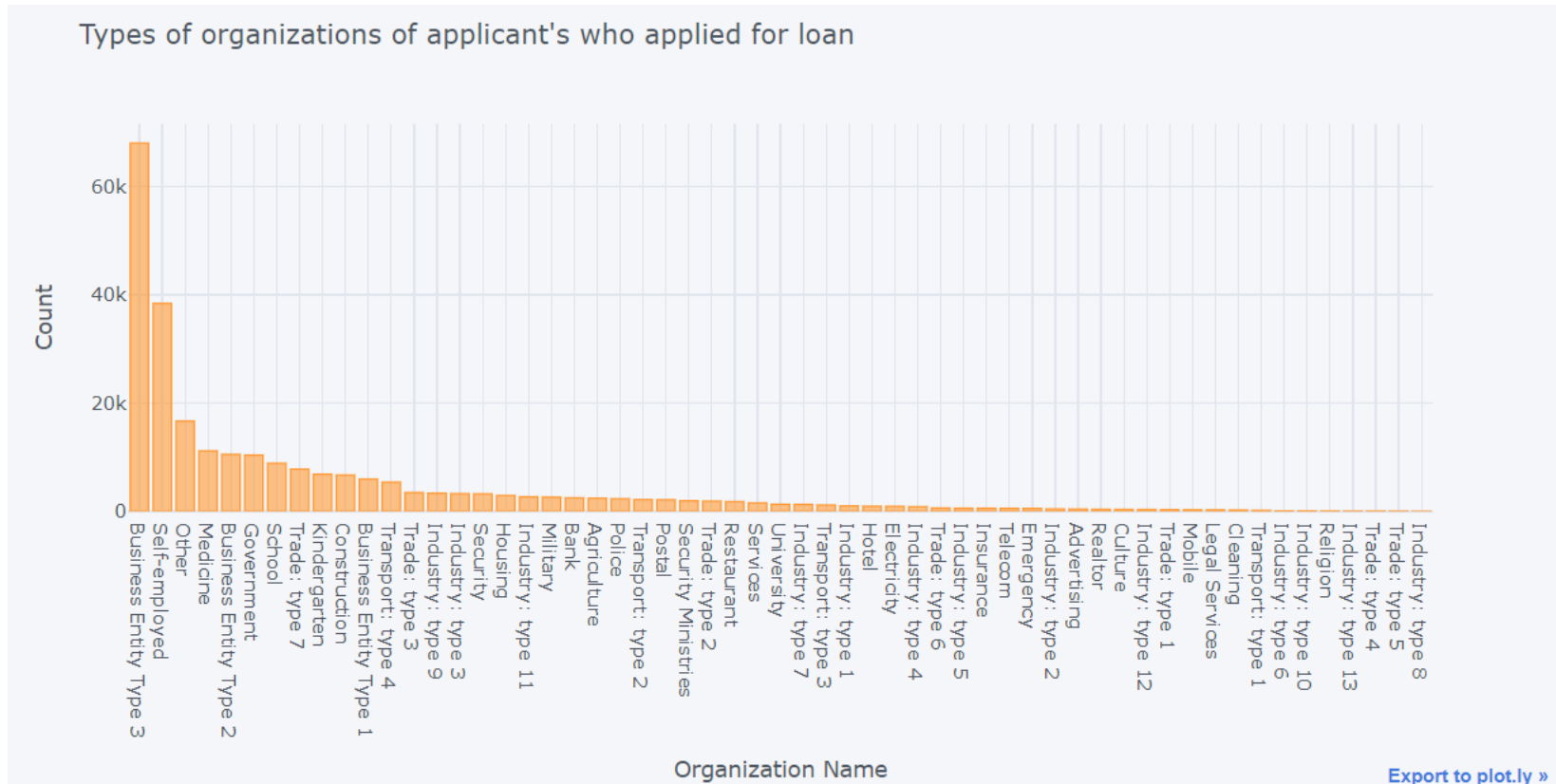
# OCCUPATION OF THE APPLICANTS



- Laborers is the highest occupation type among applicants while IT staff is the lowest

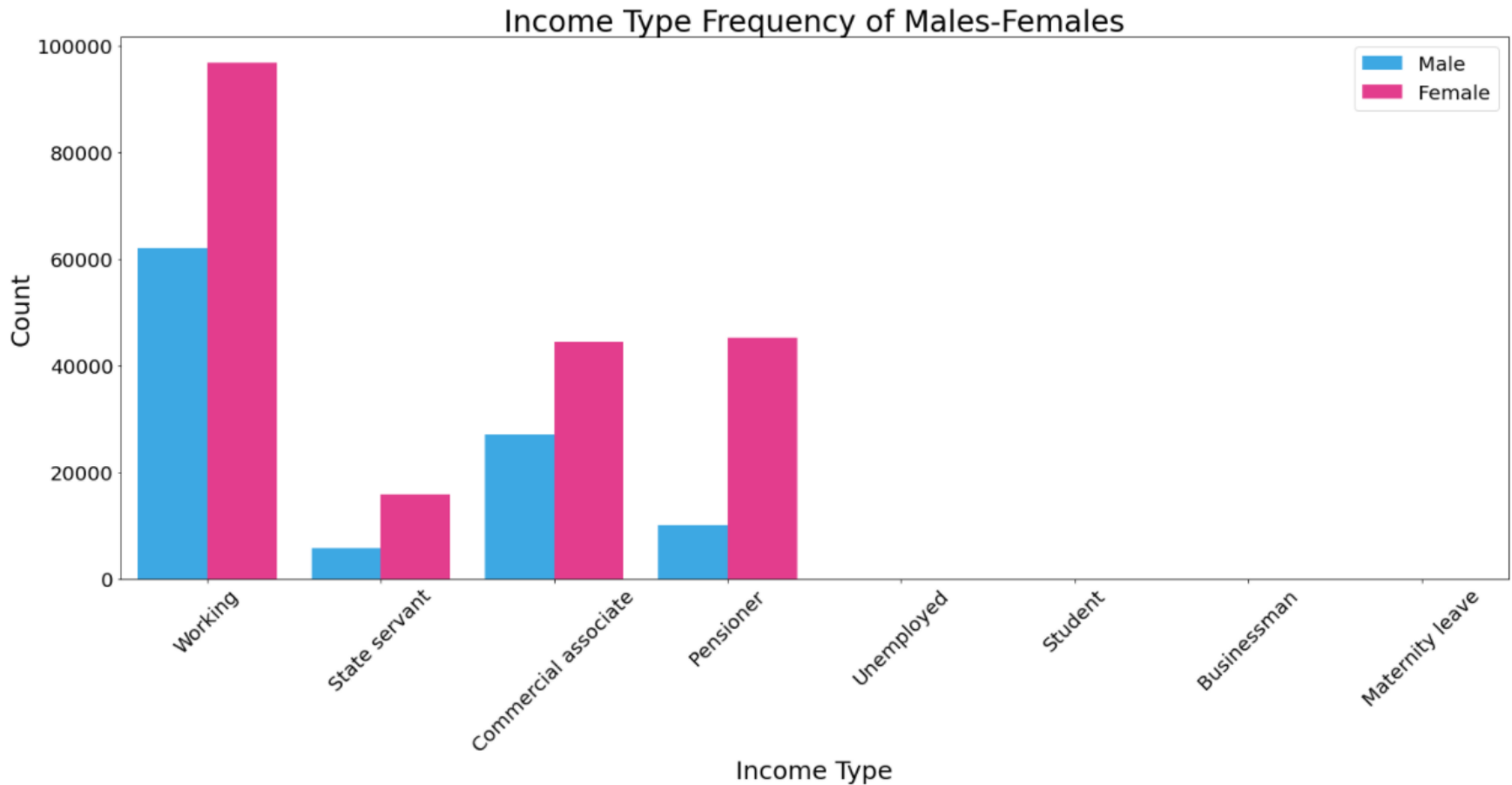


# ORGANIZATIONS OF APPLICANTS



- Business Entity Type 3 and self employed are the highest for applying loans
- There are less clients from 'Industry: Type 8', 'Trade : Type 5'

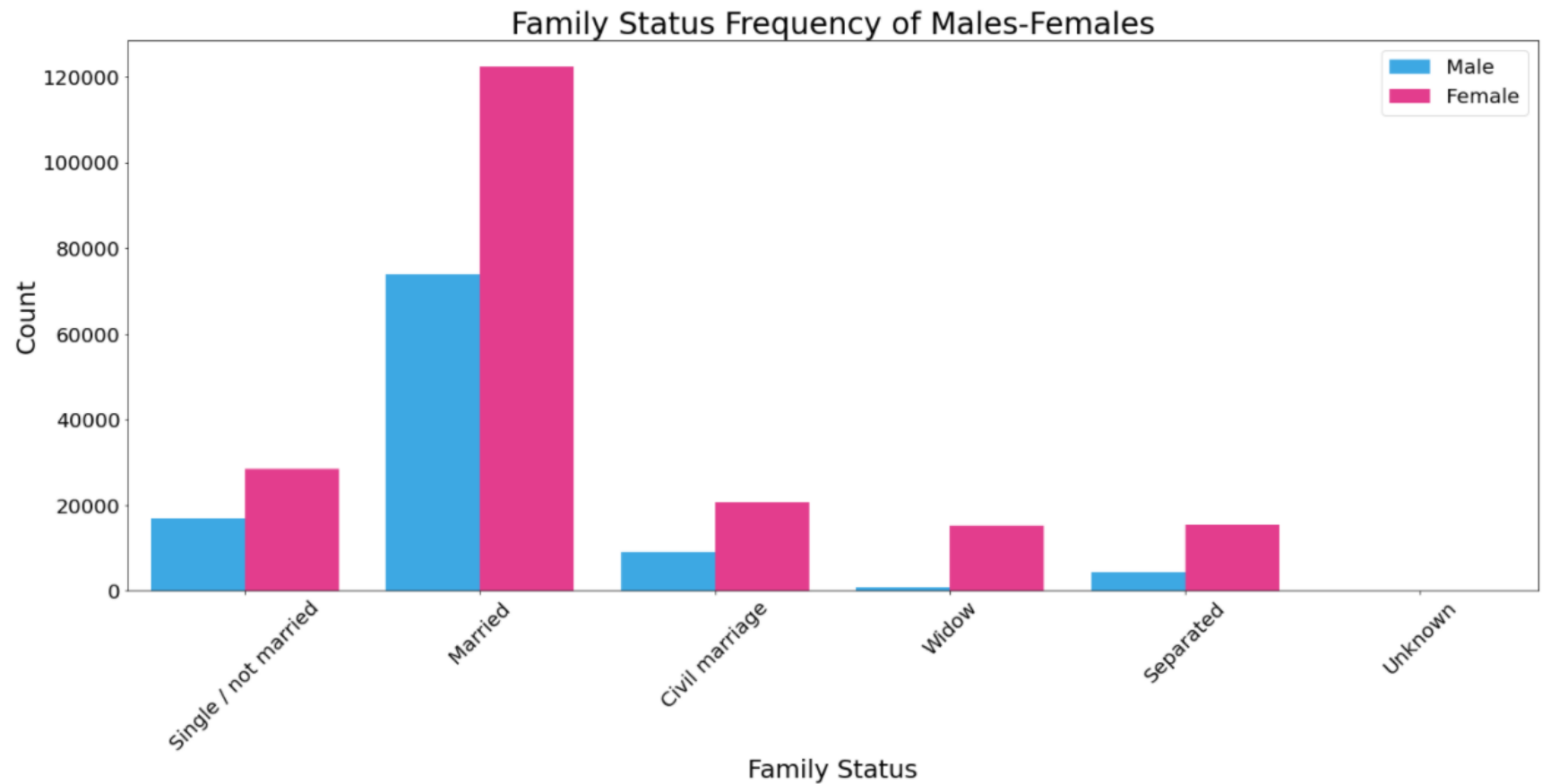
# INCOME SOURCE



- Working people have the highest number of applications.



# FAMILY STATUS

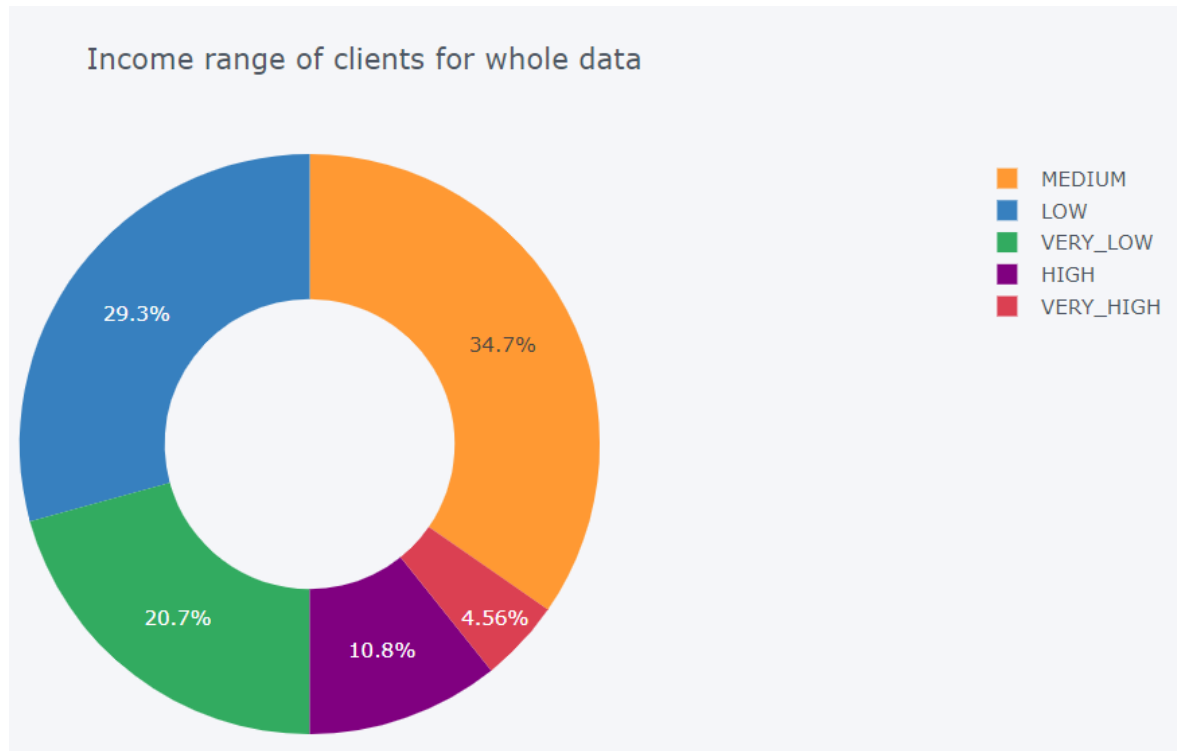


- Married people have the highest number of applications
- Widow males apply for loans very less than that of females.





# INCOME RANGE



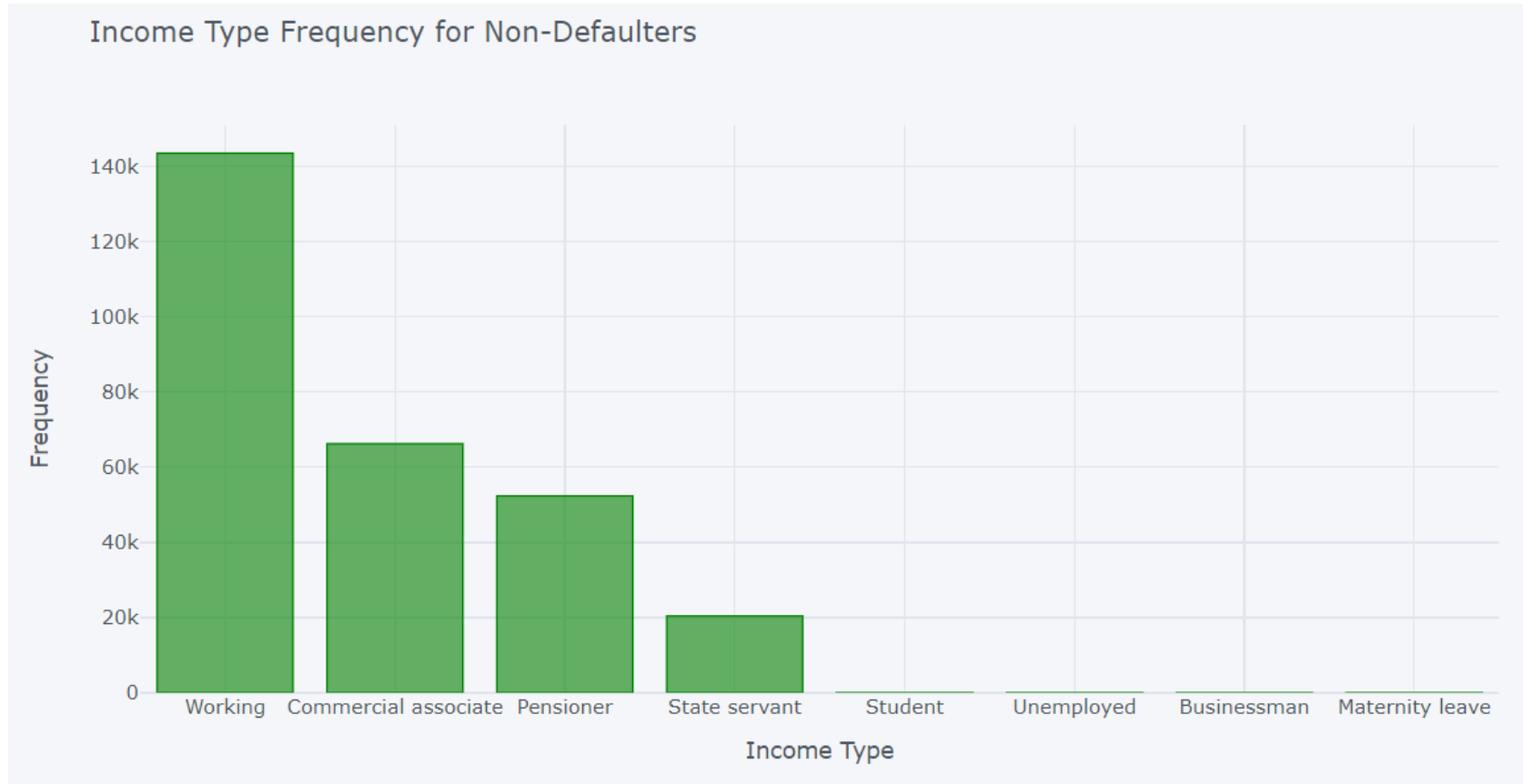
- As the income of people increases, they do not take much loans
- Medium income range people have the highest number of applications.



# UNIVARIATE ANALYSIS – CATEGORICAL DATA ON TARGET 0



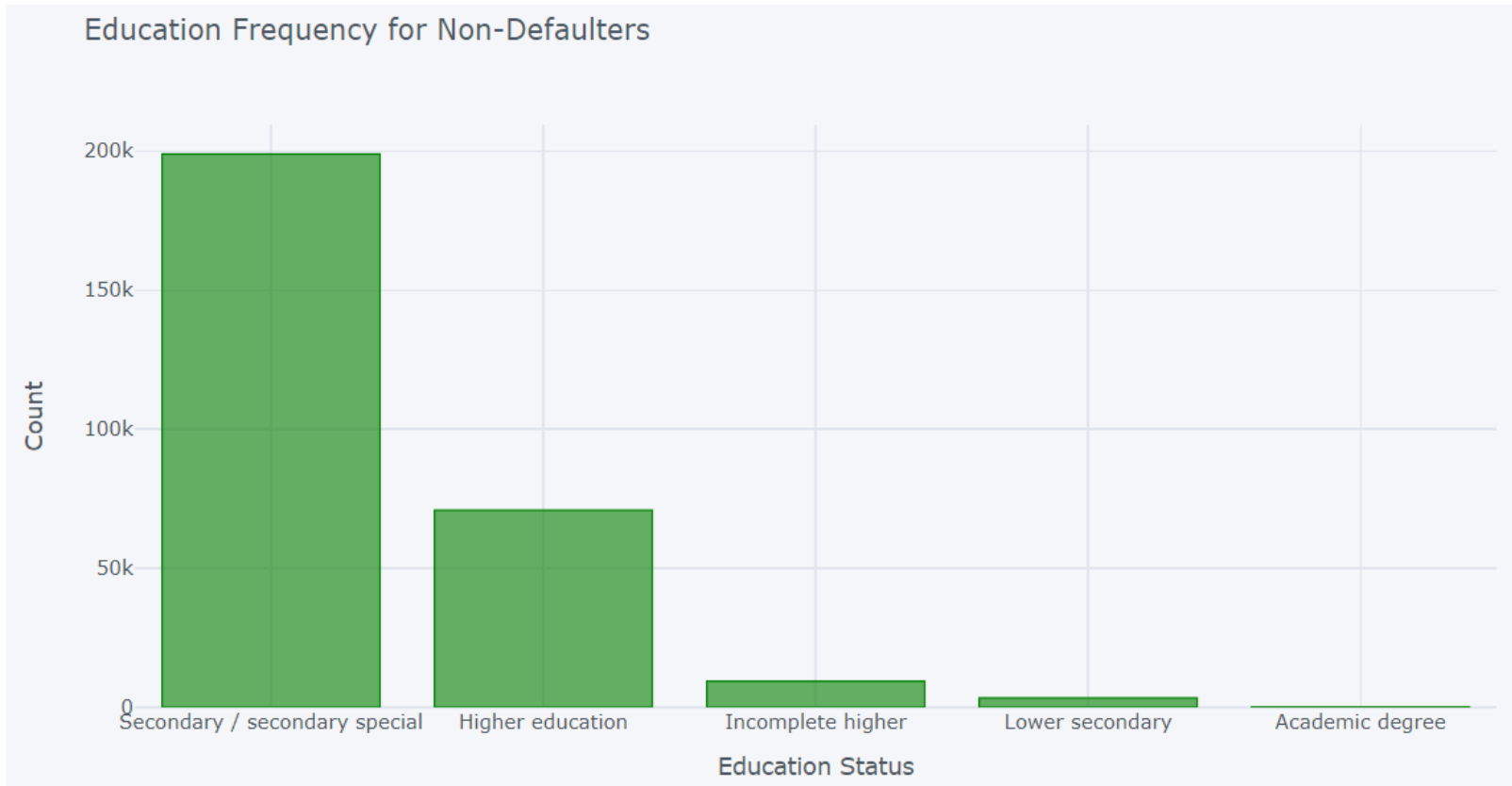
# INCOME SOURCE



- Working people have the highest application.



# EDUCATION



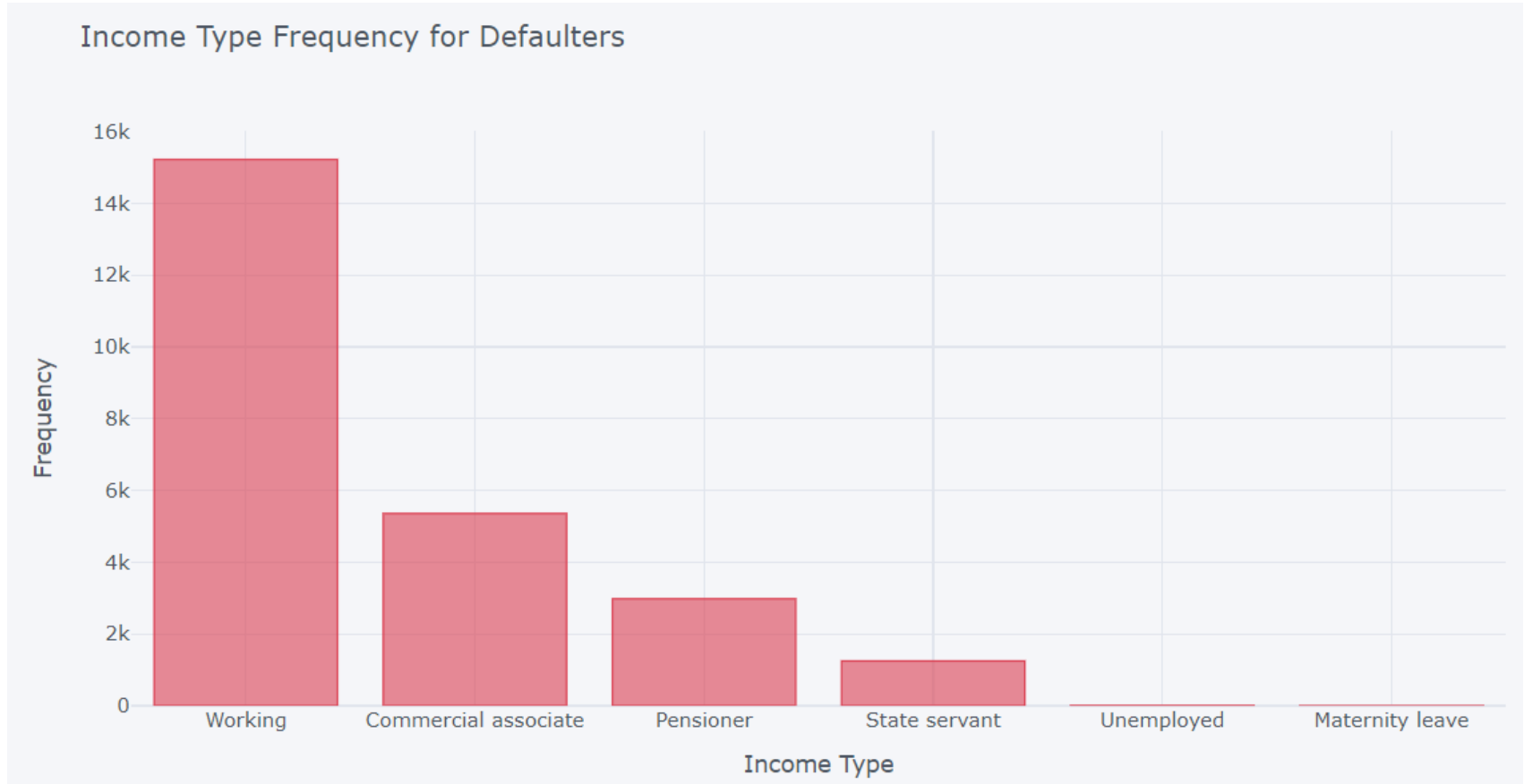
- People with Secondary Education have the highest applications of loans.



# UNIVARIATE ANALYSIS – CATEGORICAL DATA ON TARGET 1



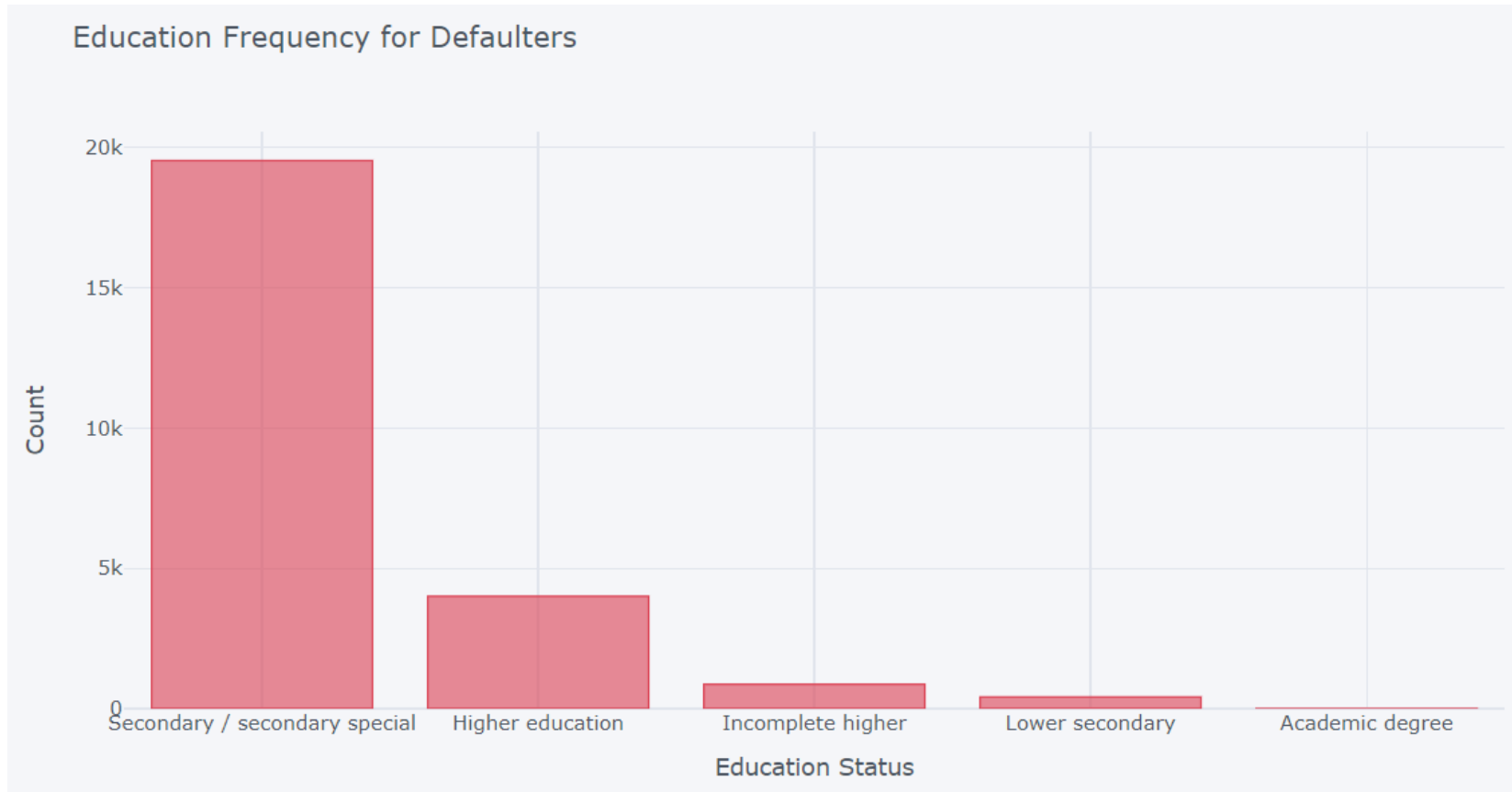
# INCOME SOURCE



- Students do not default according to the data we have



# EDUCATION



- People with Higher education defaults less when compared to the non-default rate.



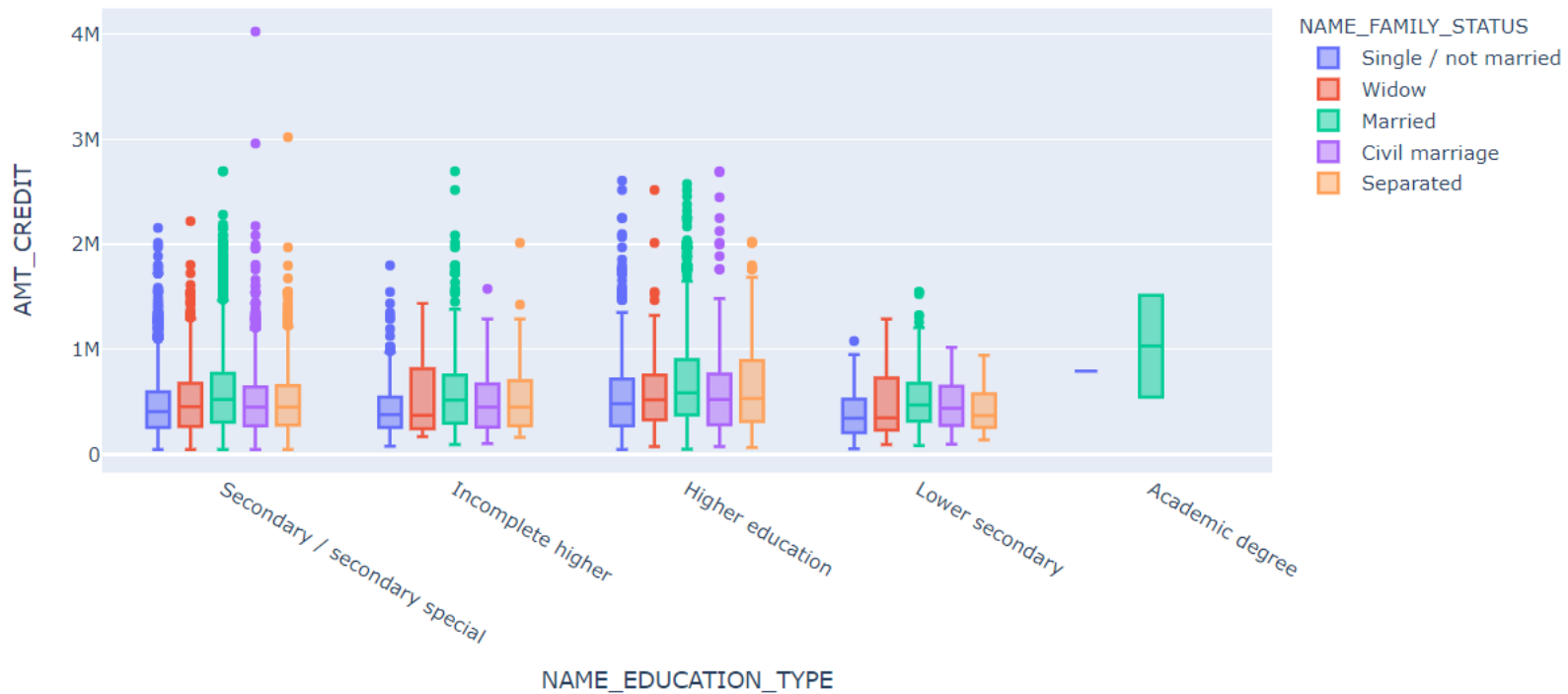
# BIVARIATE ANALYSIS – CATEGORICAL VS NUMERICAL





# CREDIT AMOUNT VS EDUCATION

Credit amount vs Education for clients with Payment Difficulties

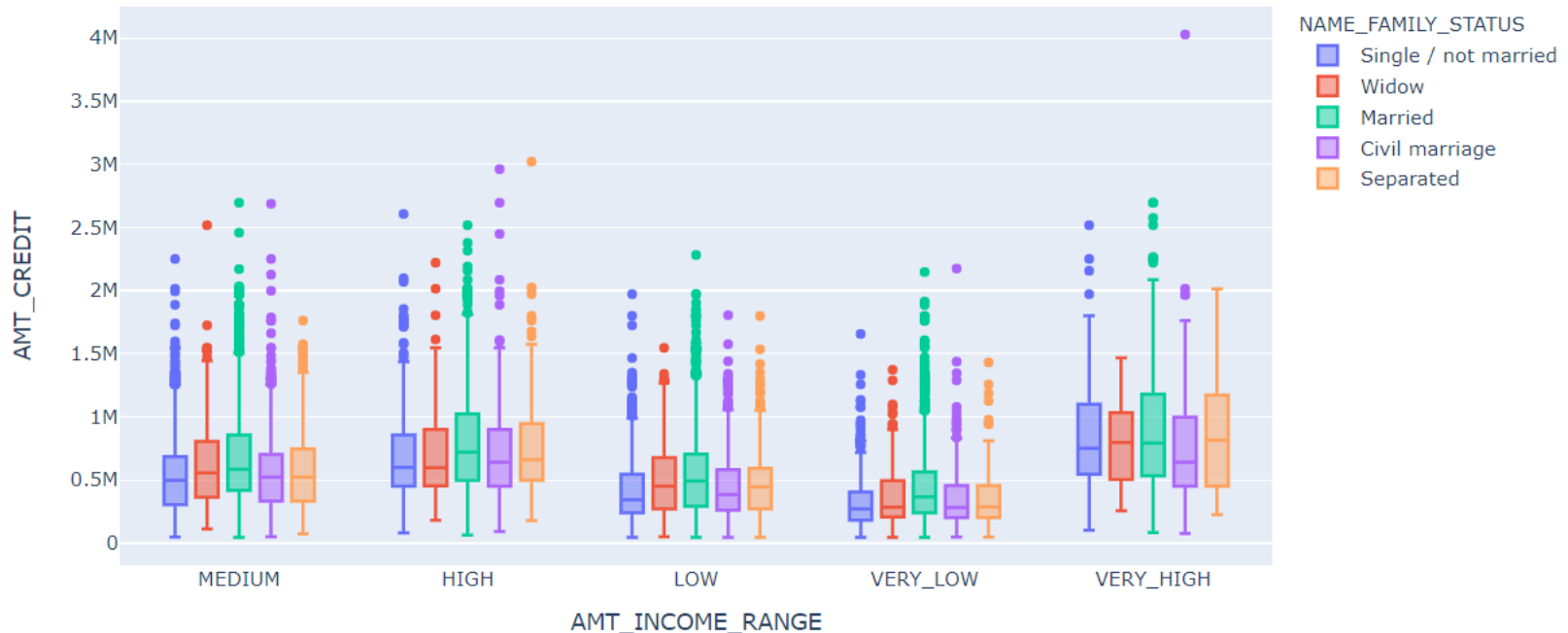


- It can be said that people who have academic degree and have status other than married are safe to provide loans.



# INCOME RANGE VS CREDIT AMOUNT

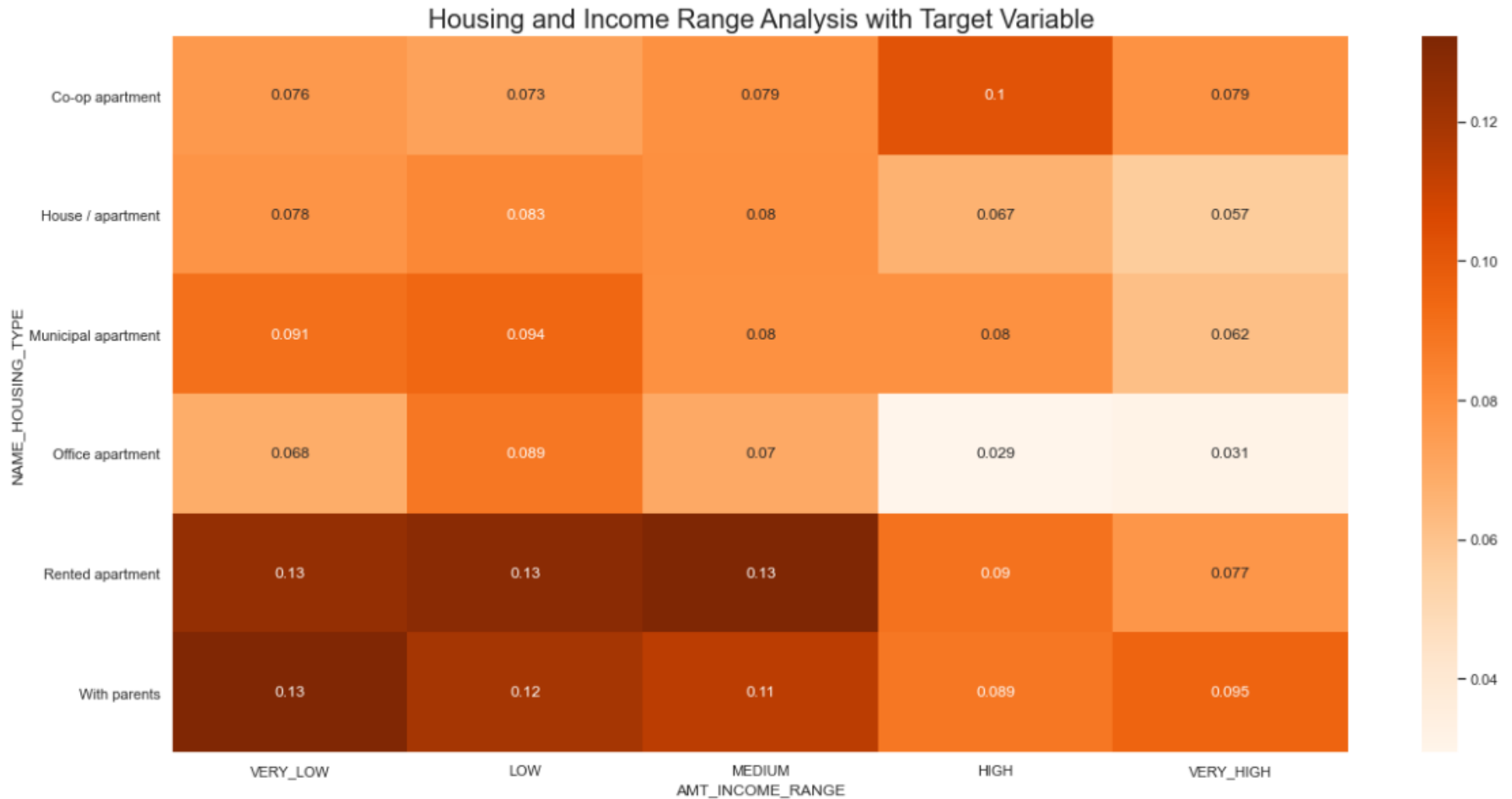
Income range vs Credit amount for clients with Payment Difficulties



- It can be seen here that the people who have taken credit amount for more than 3M have less probability of defaulting for Medium, High and very high income range.



# HOUSING – INCOME RANGE WITH PAYMENT DIFFICULTIES

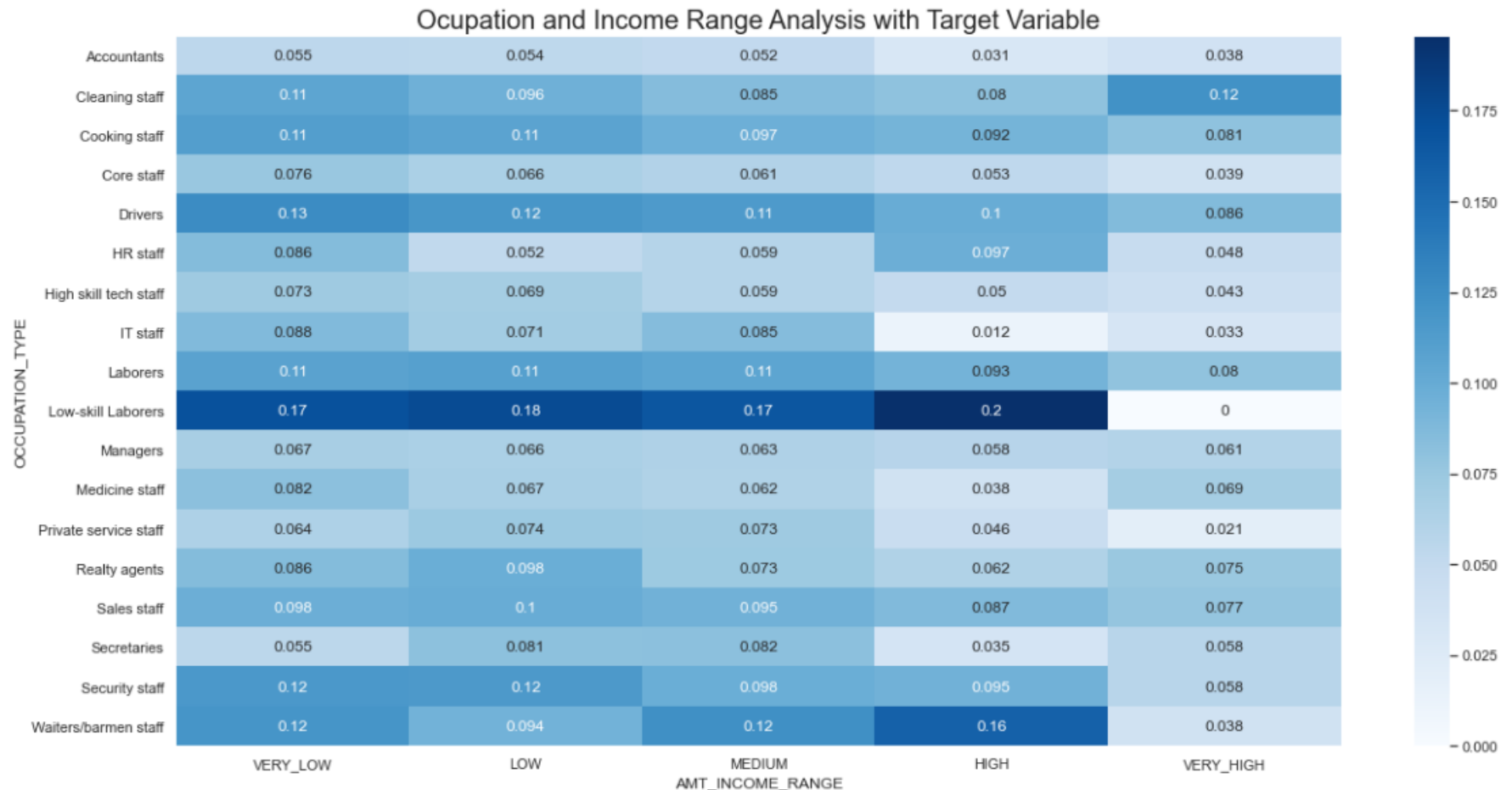


# HOUSING – INCOME RANGE WITH PAYMENT DIFFICULTIES : INSIGHTS

- People staying with parents or staying in rented apartments with income of medium or less have high chances of defaulting.
- People staying in office apartments and have very high income have very less chances of defaulting.



# OCCUPATION – INCOME RANGE WITH PAYMENT DIFFICULTIES



# OCCUPATION – INCOME RANGE WITH PAYMENT DIFFICULTIES : INSIGHTS

- Low skill laborers have high chance of defaulting.
- IT staff with High salary income has the lowest chances of defaulting.
- Accountants overall have less chance of defaulting.

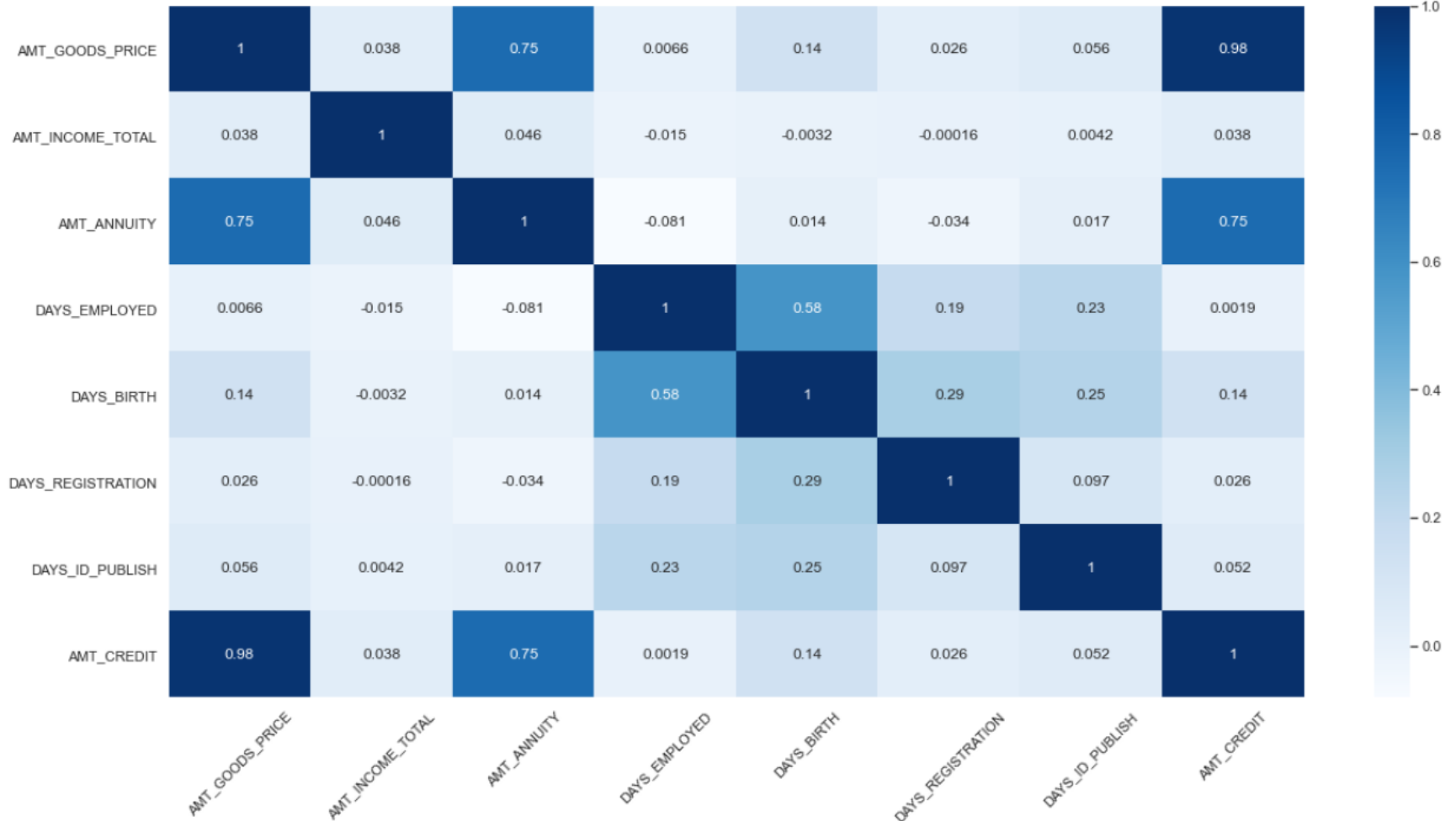


# BIVARIATE ANALYSIS – NUMERICAL VS NUMERICAL



# LOAN – WITH PAYMENT DIFFICULTIES

Loan - With Difficulties





# LOAN – WITH PAYMENT DIFFICULTIES

## INSIGHTS

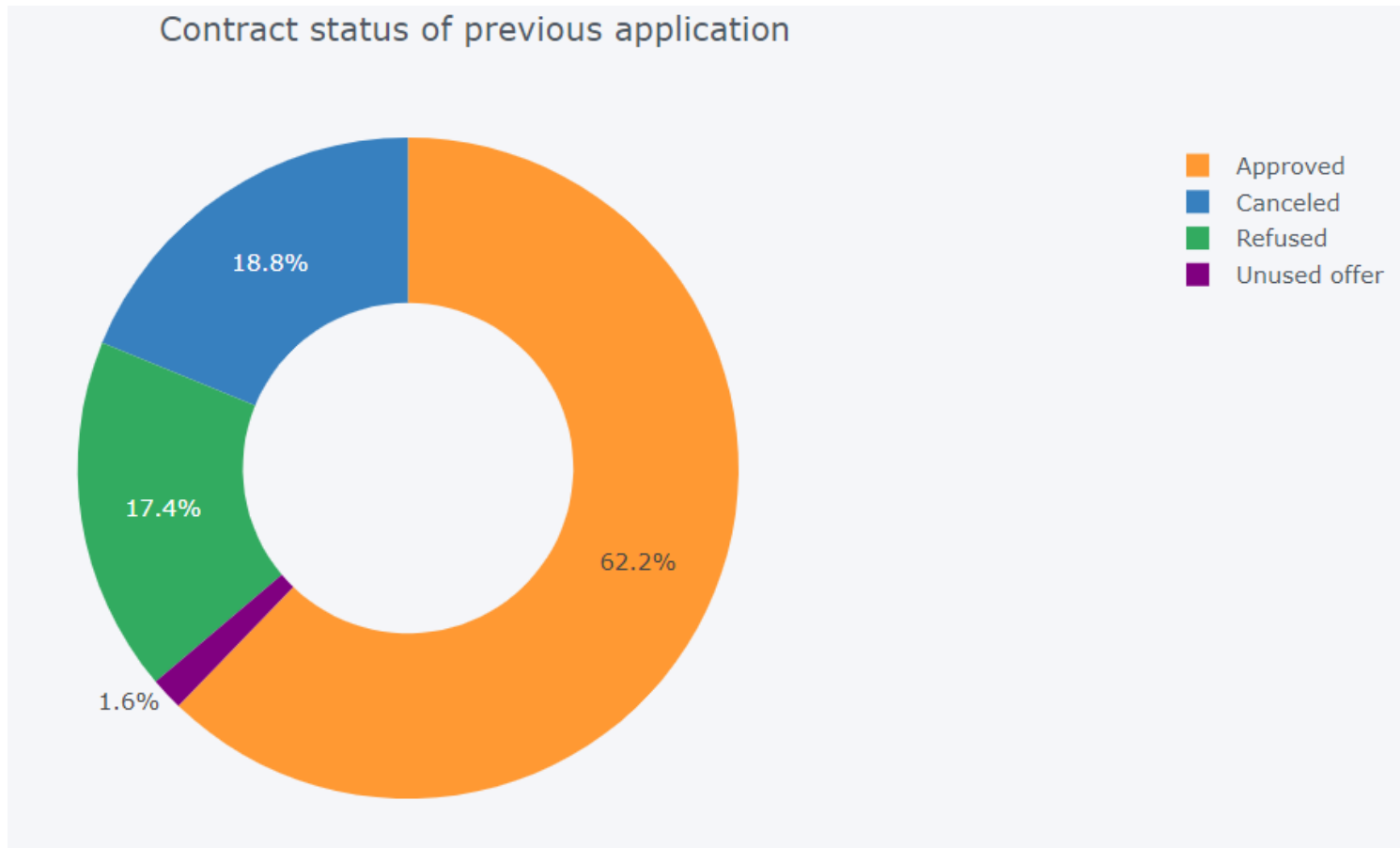
- We have high correlation between credit price and goods price, as it can be because the loan is credited only as much as the price of the goods.
- We also have high correlation between credit price and annuity as this can be because higher the credit is given, higher will be the installments.



# PREVIOUS APPLICATION DATASET



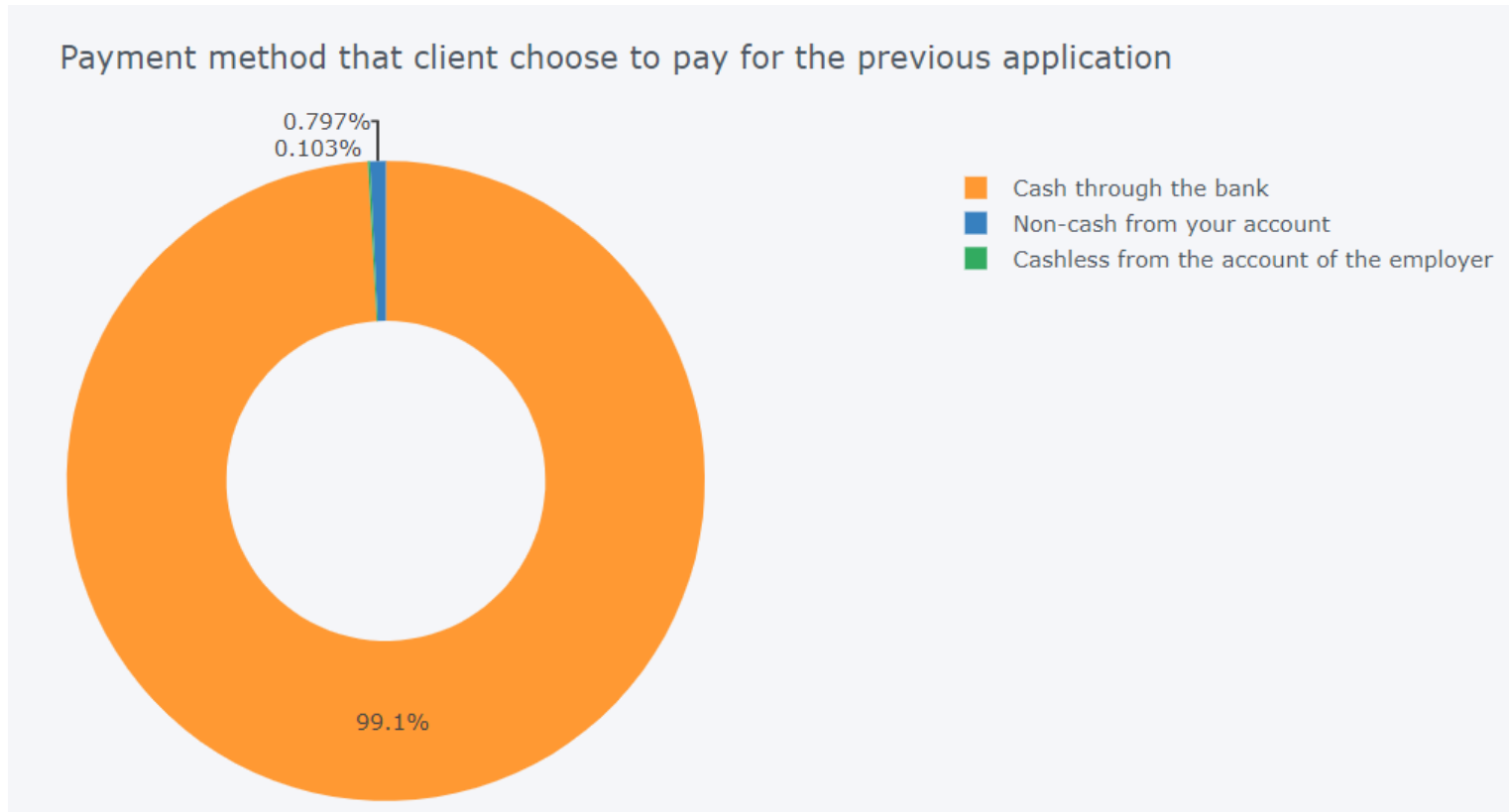
# CONTRACT STATUS



- Most of the loans do get approved.



# PAYMENT METHOD

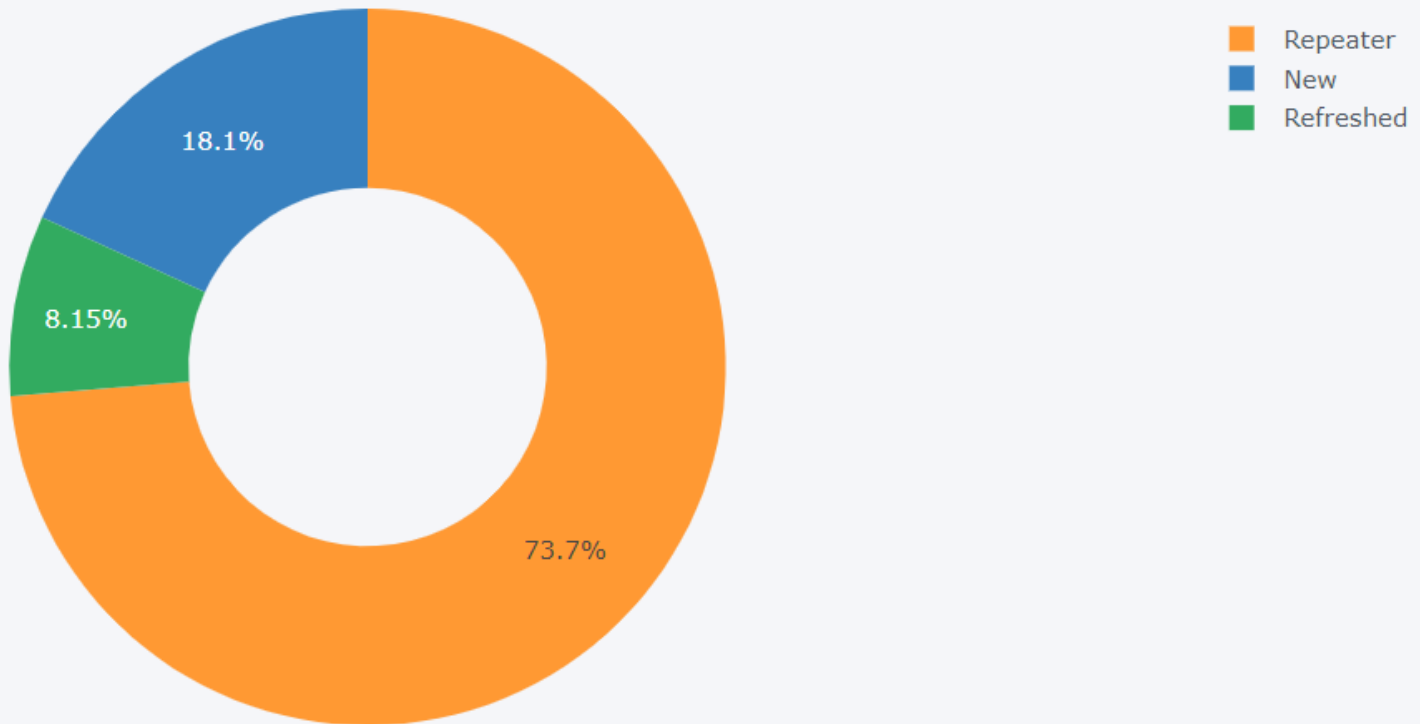


- 99% of the people uses cash through the bank for payments.



# CLIENT TYPE

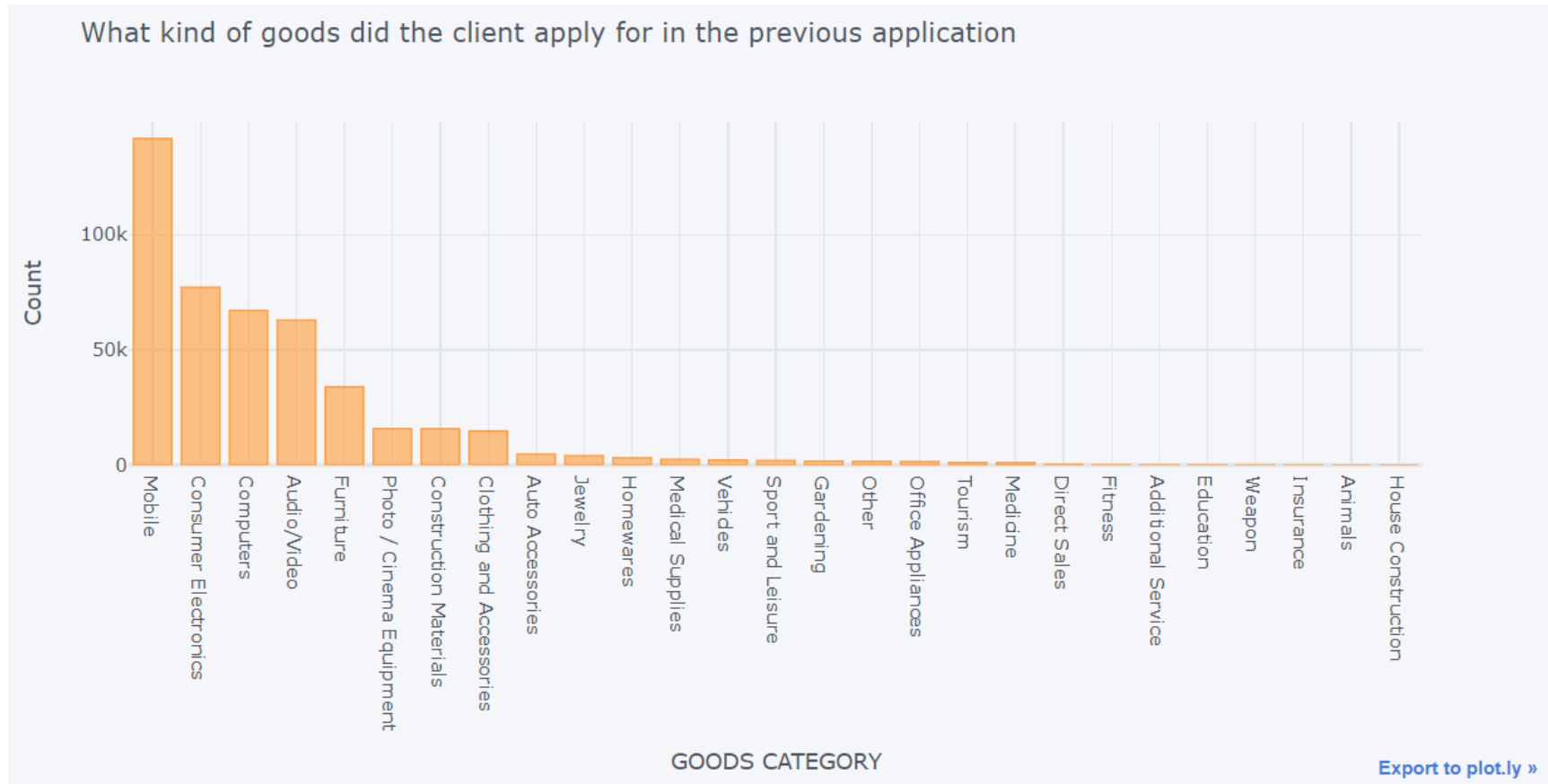
Was the client old or new client when applying for the previous application



- Maximum people are the ones who have taken loan previously.



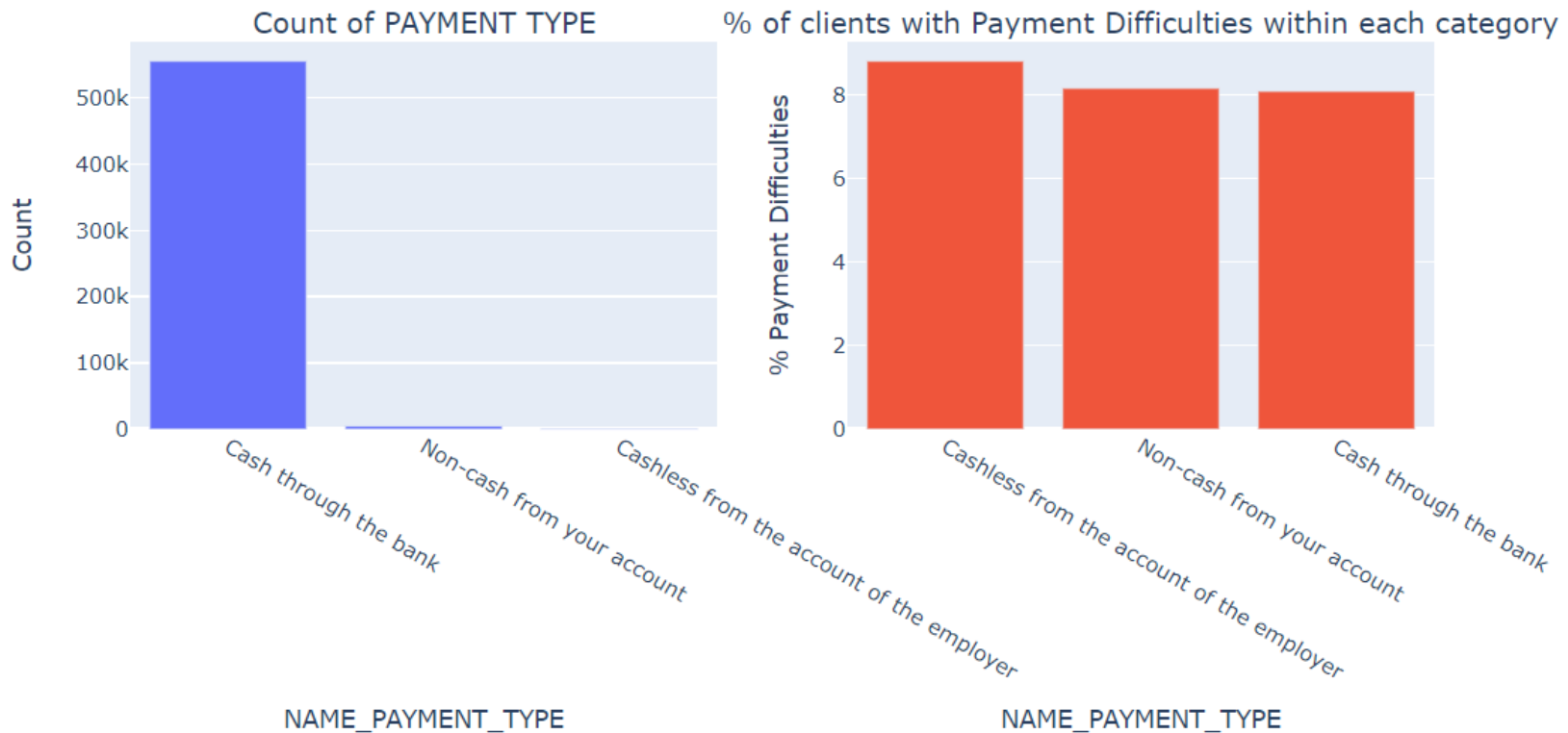
# TYPE OF GOODS



- Maximum loans are taken to buy mobile phones while least loans are taken for house construction.

# DEFAULTERS IN PAYMENT TYPE

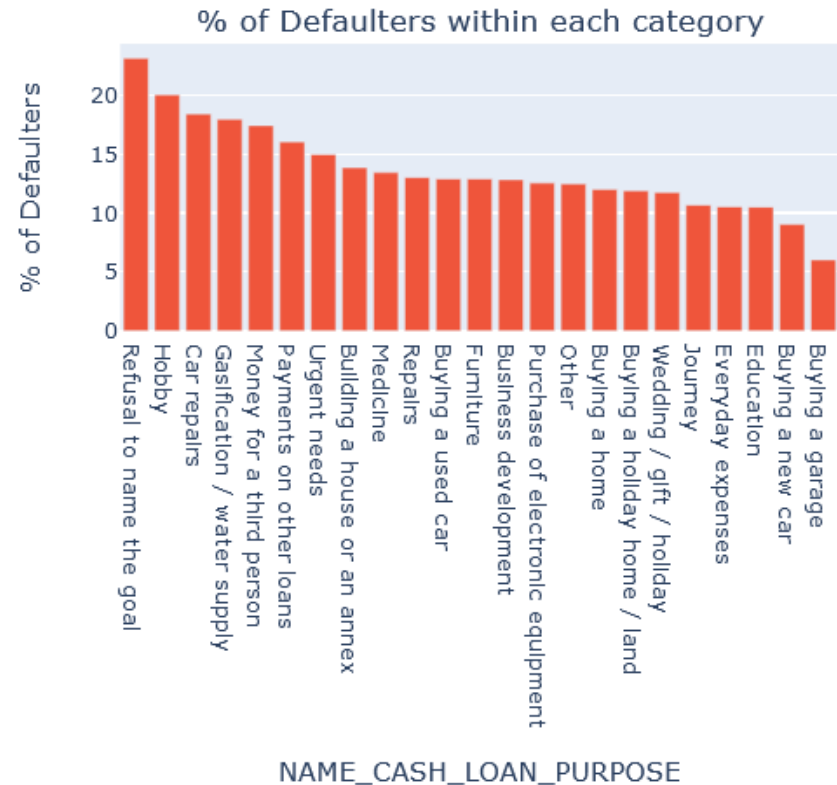
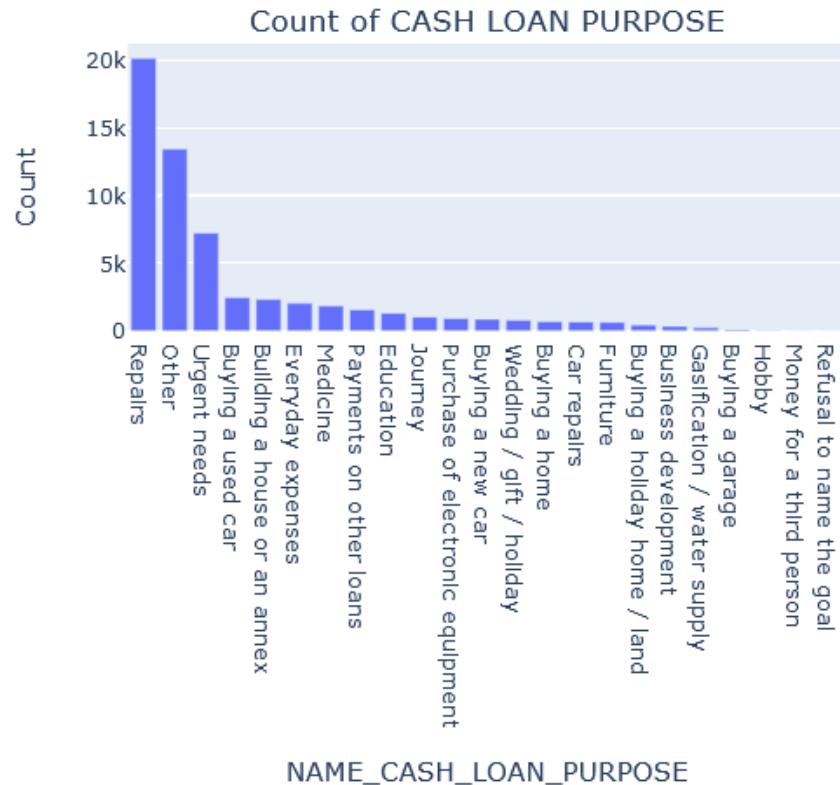
PAYMENT TYPE



- Even though maximum payments are done by cash through bank, it has less defaults than the other two methods.



# DEFAULTERS PER CASH LOAN PURPOSE

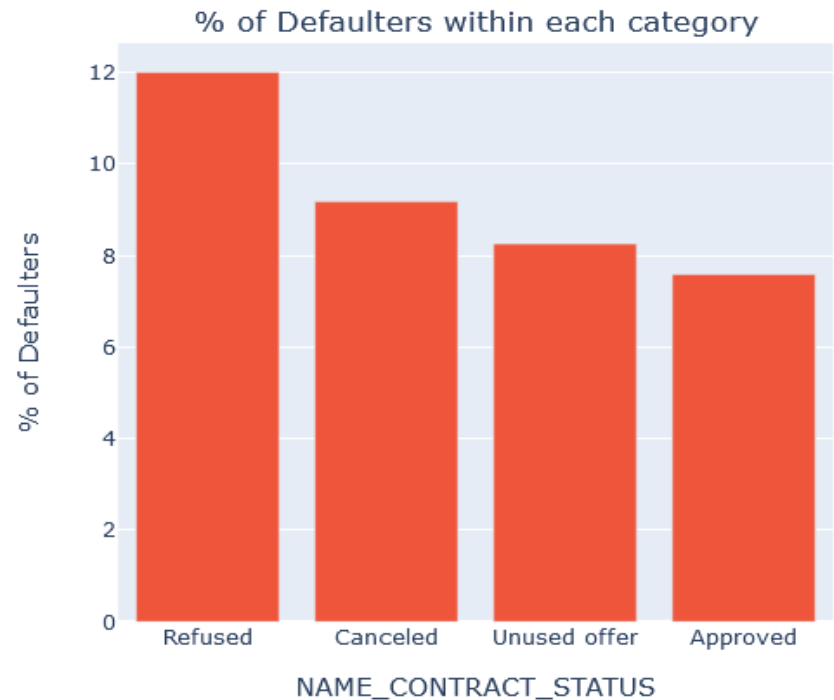
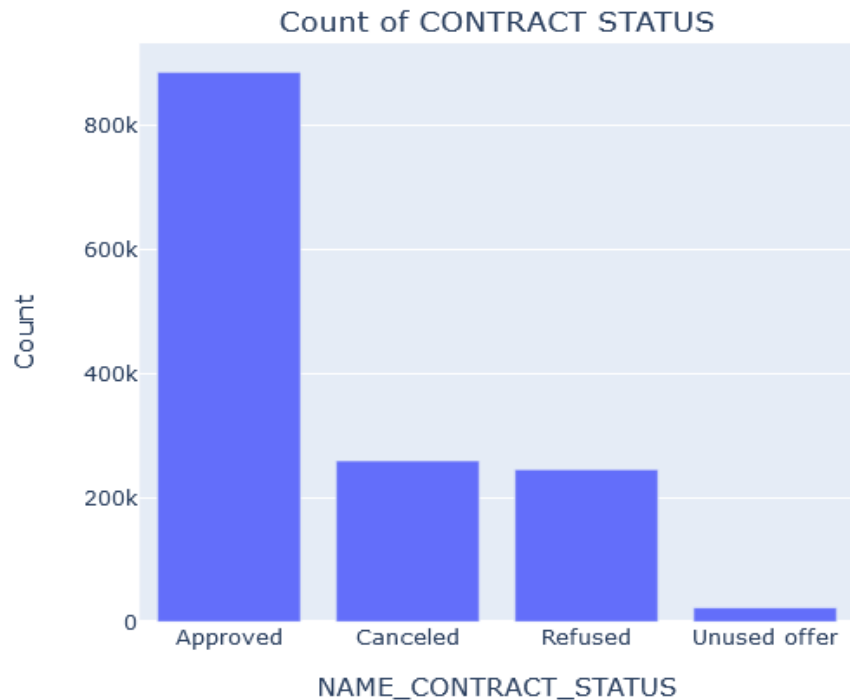


- From the first graph it can be seen that purpose of cash loan from previous data was maximum for 'Repairs'
- It can be clearly seen from the second graph that the 'Refusal to name the goal' for cash loan from previous application are the ones who have maximum % of Loan-Payment Difficulties from current application.





# DEFAULTERS PER CONTRACT STATUS

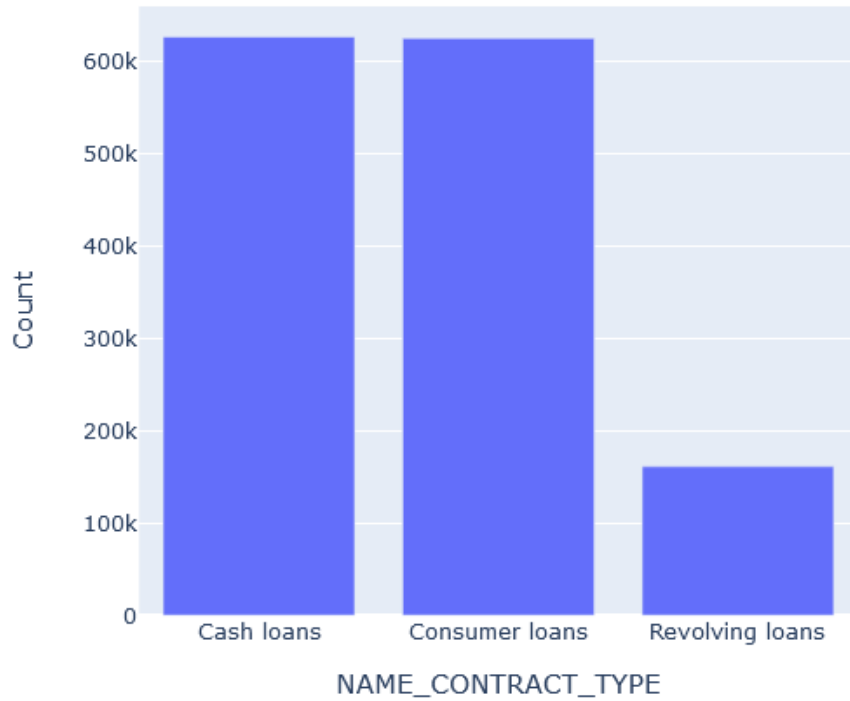


- From the first graph it can be seen that most of the contracts from previous application have been Approved.
- It can be clearly seen from the second graph that:
  - 'Refused' contracts from previous application are the ones who have maximum % of Loan-Payment Difficulties from current application.
  - 'Approved' contracts from previous application are the ones who have minimum % of Loan-Payment Difficulties from current application.

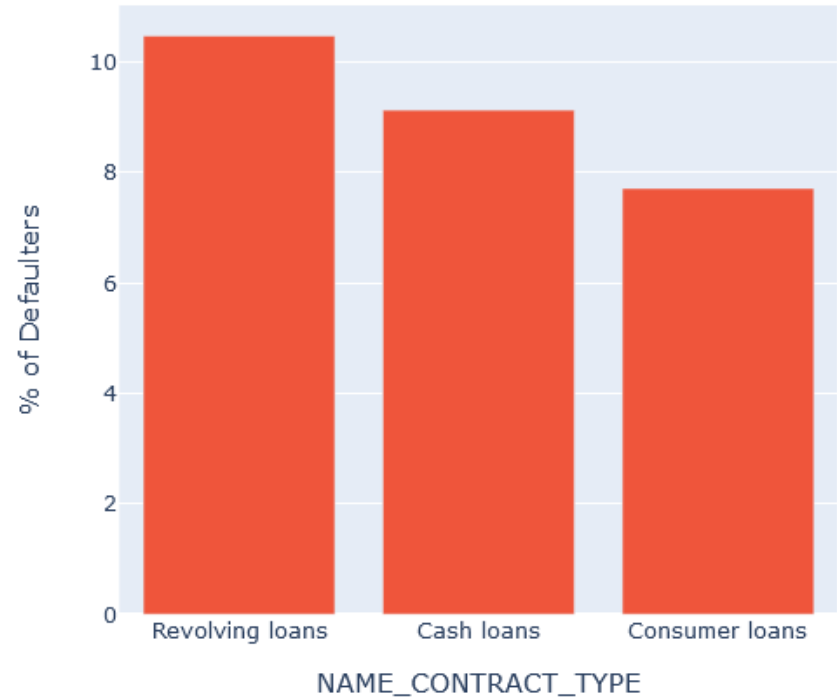


# DEFAULTERS PER CONTRACT TYPE

Count of CONTRACT TYPE



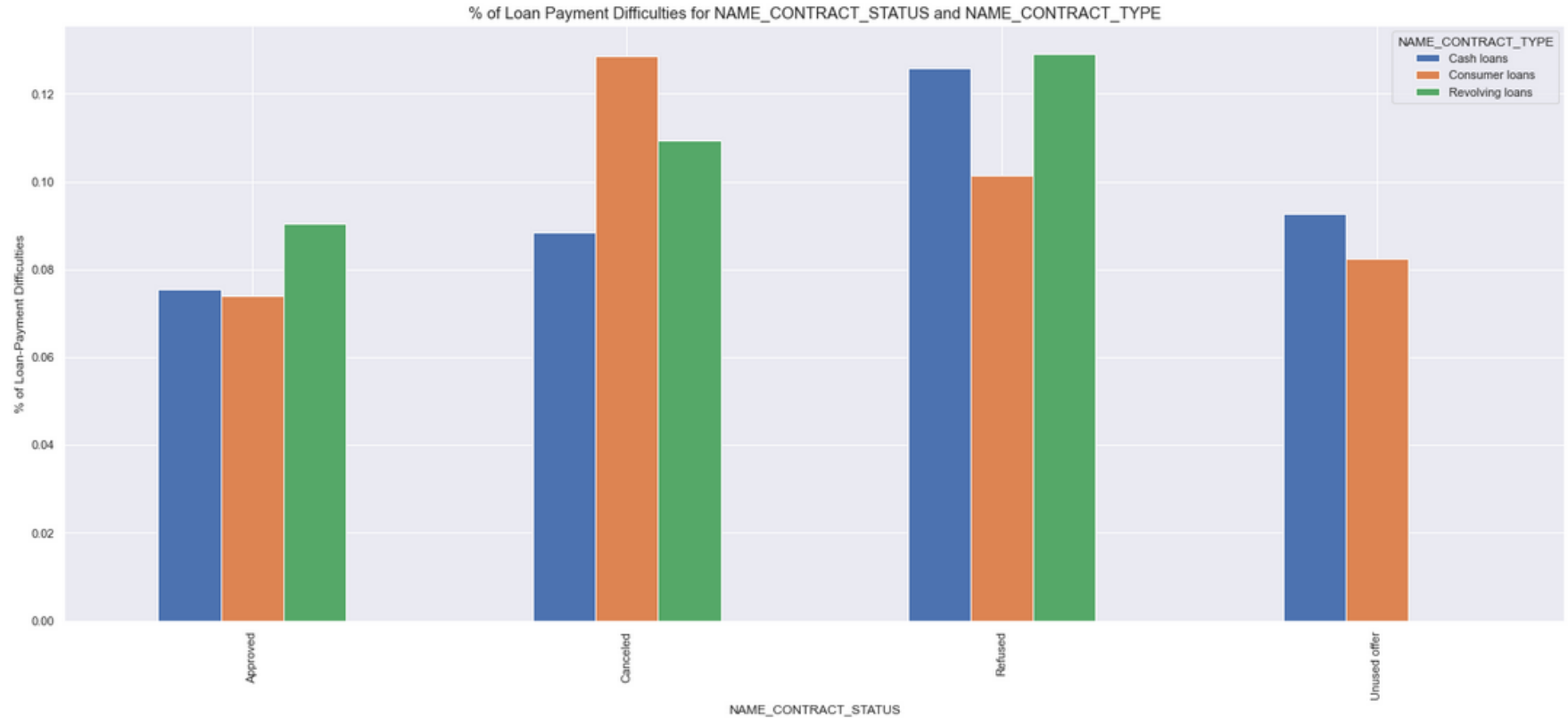
% of Defaulters within each category



- From the first graph it can be seen that most of the contract types from previous application were 'Cash loans'
- It can be clearly seen from the second graph that:
  - 'Revolving Loans' contracts from previous application are the ones who have maximum % of Loan-Payment Difficulties from current application.
  - 'Consumer loans' contracts from previous application are the ones who have minimum % of Loan-Payment Difficulties from current application.



# DEFAULTERS PER CONTRACT TYPE AND STATUS



- It can be observed from the above graph that Clients with 'Revolving loans' and with 'Refused' previous application tend to have more % of Loan-Payment Difficulties in current application.



# MAJOR INSIGHTS – NEW APPLICATION DATASET

- The count of 'Maternity Leave' in 'NAME\_INCOME\_TYPE' is very less and it also has maximum % of payment difficulties- around 40%. Hence, client with income type as 'Maternity leave' are the driving factors for Loan Defaulters.
- The count of 'Low skilled Laborers' in 'OCCUPATION\_TYPE' is comparatively very less and it also has maximum % of payment difficulties- around 17%. Hence, client with occupation type as 'Low skilled Laborers' are the driving factors for Loan Defaulters.
- The count of 'Lower Secondary' in 'NAME\_EDUCATION\_TYPE' is comparatively very less and it also has maximum % of payment difficulties- around 11%. Hence, client with education type as 'Lower Secondary' are the driving factors for Loan Defaulters.



# MAJOR INSIGHTS –PREVIOUS APPLICATION DATASET

- The count of 'Refusal to name the goal' in 'NAME\_CASH\_LOAN\_PURPOSE' is comparatively very less and it also has maximum % of payment difficulties- around 23%. Hence, clients who have 'Refused to name the goal' for cash loan in previous application are the driving factors for Loan Defaulters.
- The count of 'Refused' in 'NAME\_CONTRACT\_STATUS' is comparatively less and it also has maximum % of payment difficulties- around 12%. Hence, client with contract status as 'Refused' in previous application are the driving factors for Loan Defaulters.
- The count of 'Revolving Loans' in 'NAME\_CONTRACT\_TYPE' is comparatively very less and it also has maximum % of payment difficulties- around 10%. Hence, client with contract type as 'Revolving loans' in previous application are the driving factors for Loan Defaulters.



- Clients with 'Revolving loans' and with 'Refused' previous application tend to have more % of payment difficulties in current application. Since the count of both 'Revolving loans' and 'Refused' is comparatively less (from the graphs in previous slide), clients with 'Revolving Loans' and 'Refused' previous application are driving factors for Loan Defaulters.



**THANK YOU!**

