

# **Bank Loan Case Study**

Final Project 2

Project By: Sumit Gope

# Project Description:

The aim of this project is to apply exploratory data analysis (EDA) techniques to understand the patterns present in loan application data and identify the driving factors behind loan default. The dataset used in this project contains information about loan applications at the time of applying for the loan, as well as previous loan data for clients.

The problem statement for this project is to identify patterns that indicate if a client has difficulty paying their installments, which can be used to take actions such as denying the loan, reducing the amount of loan, or lending to risky applicants at a higher interest rate. The company wants to understand the driving factors behind loan default, i.e. the variables that are strong indicators of default, which can be used for portfolio and risk assessment.

# Approach

The analysis approach for this project involves the following steps:

- 1.Data Understanding:** This involves understanding the structure and contents of the dataset and data dictionary provided, and identifying any missing data and outliers.
- 2.Data Cleaning:** This involves dealing with missing data and outliers using appropriate methods, such as removing columns or replacing missing data with an appropriate value.
- 3.Univariate Analysis:** This involves analyzing individual variables in the dataset and identifying any trends or patterns.
- 4.Bivariate Analysis:** This involves analyzing the relationship between two variables in the dataset and identifying any trends or patterns.
- 5.Top Correlation Analysis:** This involves finding the top 10 correlations for the client with payment difficulties and all other cases, segmented by the target variable, to identify the driving factors behind loan default.

# Tech-Stack Used



**Excel 2010 and  
Microsoft 365 web**



**Power BI**



**PowerPoint**

## Overall Approach of analysis:

The bank's problem statement is to identify the major cause of bank loan default. So below are the data sets provided to analysis.

**1.application\_data.csv** contains all of the client's information at the time of application. The information related to whether or not a client is having financial issues.

**2.previous\_application.csv** provides data from the client's previous loans. It indicates if the prior application was Accepted , cancelled Refused or Unused.

Both sets of data contained many undesired columns that will not be used for risk analytics as well as many blanks. So I cleaned up the data.

Data sets contains Categorical variables as well as Numerical variables.

# Insights (Application\_data.csv)

## Data Cleaning

The dataset contained 3 lac+ rows and 161 columns of data.

I deleted numbers of columns containing more than 45% of blank data.

For this I used `=countblank()` function.

After got total blank rows of each column I divided total number of column and multiply by 100 so we got percentage of blank column.

After deleting unwanted columns we got 72 remaining columns.

## Outliers

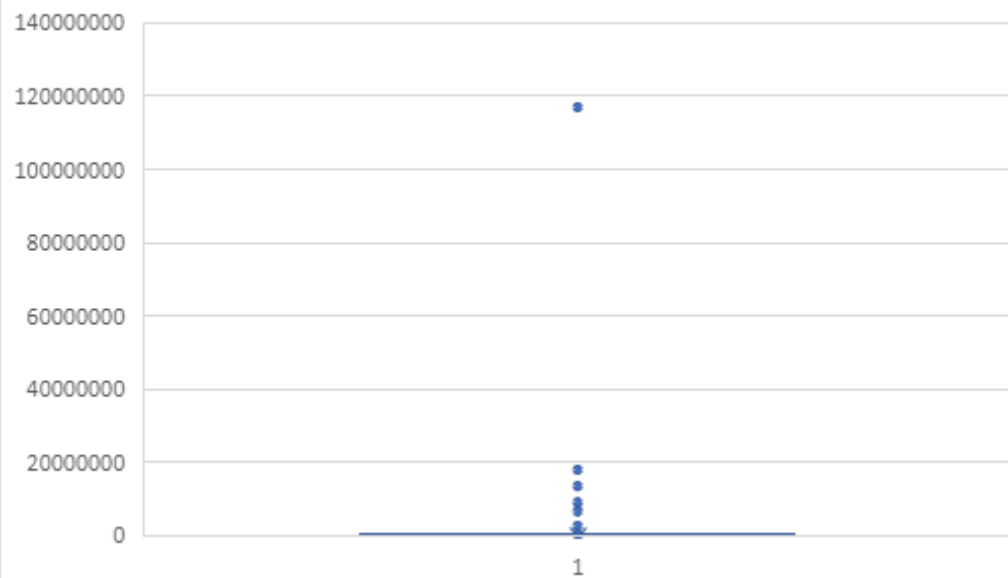
Outliers can only be identified on Numerical variable.  
**To find outlier basically I used Box and Whisker chart.**



**Statistical formula on excel that, why it is outlier.**

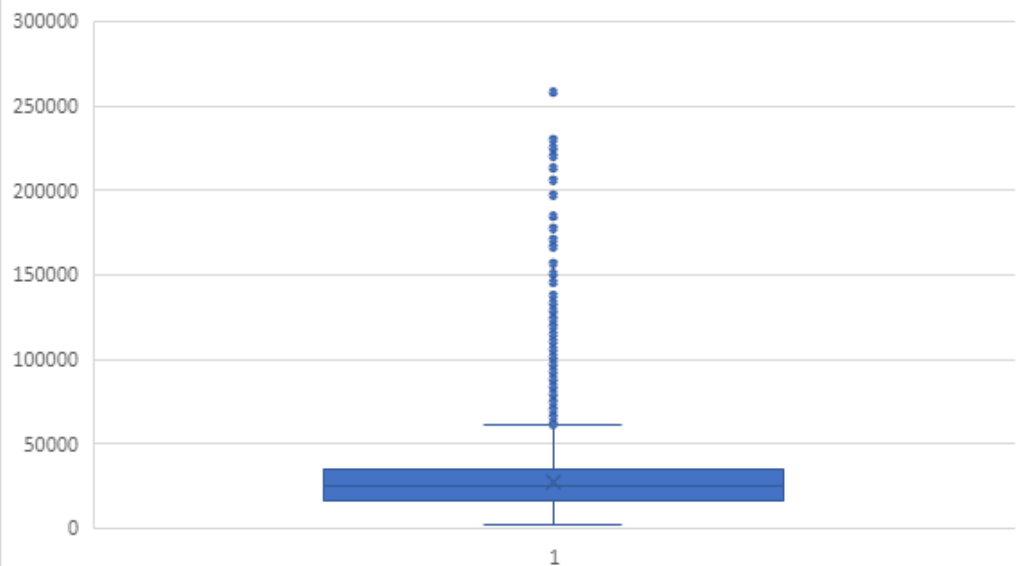
Q1	270000
Q3	808650
IQR (Q3-Q1)	538650
UPPER BOUND $[Q3+(1.5*IQR)]$	1616625
LOWER BOUND $[Q1-(1.5*IQR)]$	-942975

# AMT\_INCOME\_TOTAL



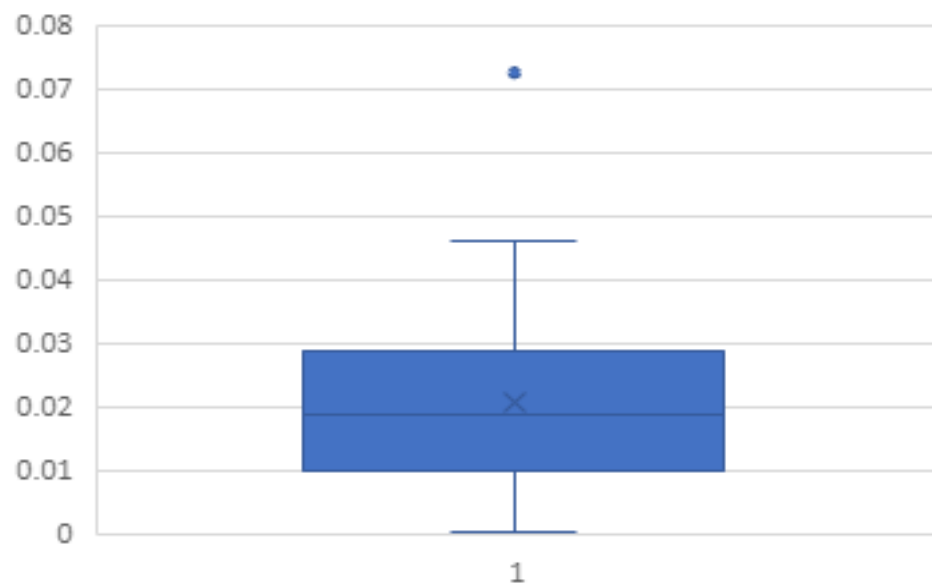
Q1	112500
Q3	202500
IQR (Q3-Q1)	90000
UPPER BOUND [Q3+(1.5*IQR)]	337500
LOWER BOUND [Q1+(1.5*IQR)]	-22500

# AMT\_ANNUITY

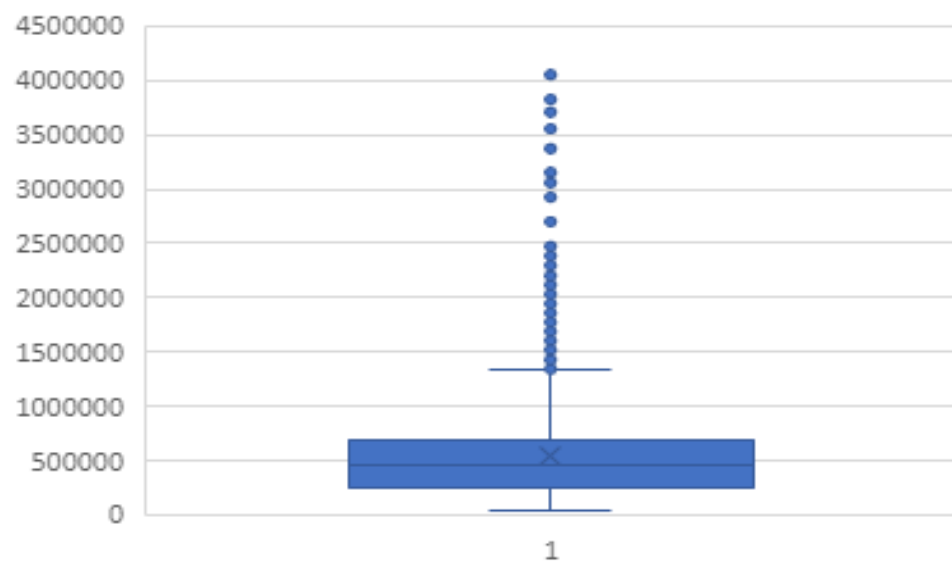




# REGION\_POPULATION\_RELATIVE



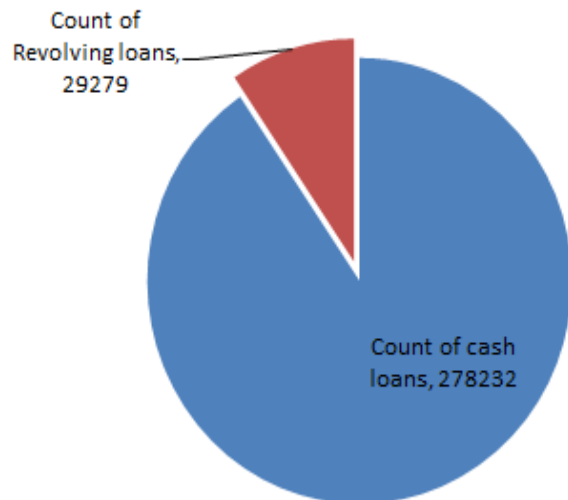
# AMT\_GOODS\_PRICE



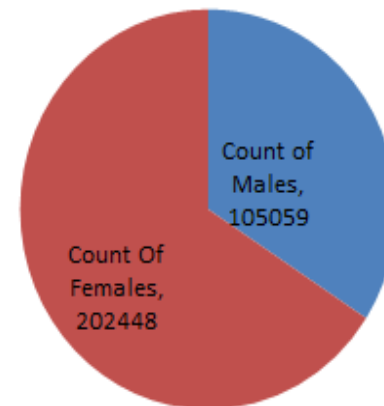
# Data Imbalance and Ratio

Data Imbalance occurs when data is disseminated in an unequal manner. I plotted data imbalance using pivot table.

Count of cash loans	Count of Revolving loans	Ratio
278232	29279	9.502783565

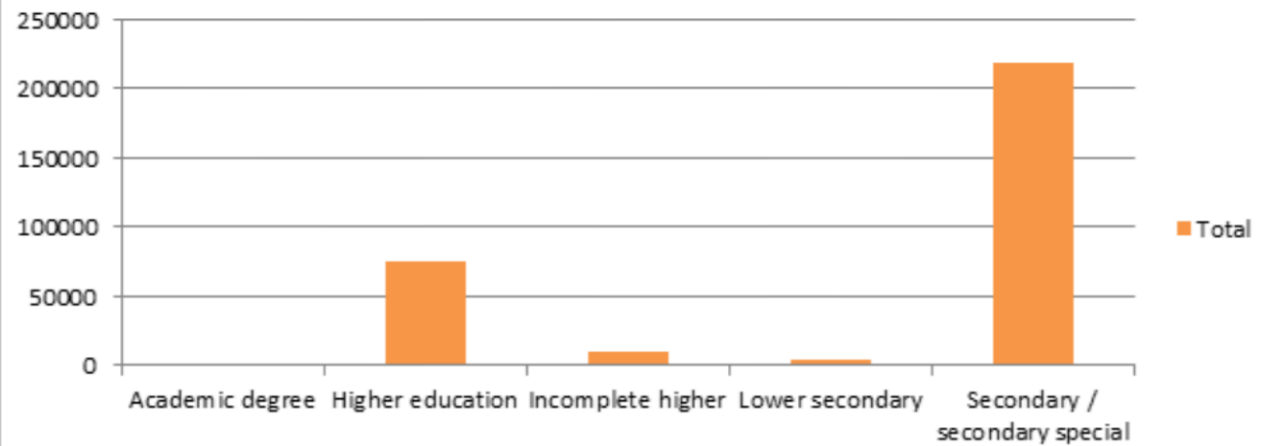


Count of Males	Count Of Females	Ratio
105059	202448	0.518943



Count of NAME\_EDUCATION\_TYPE

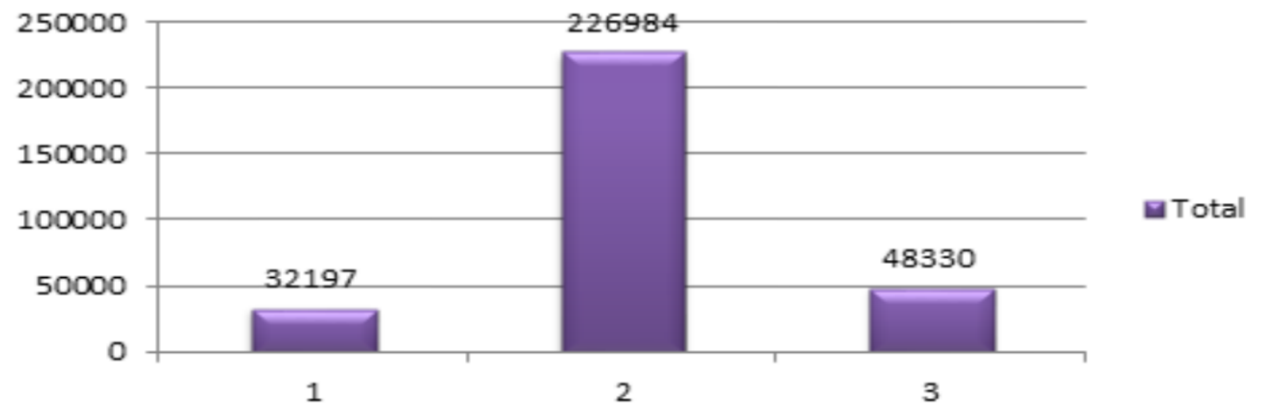
## EDUCATION TYPE



NAME\_EDUCATION\_TYPE ▼

Count of REGION\_RATING\_CLIENT

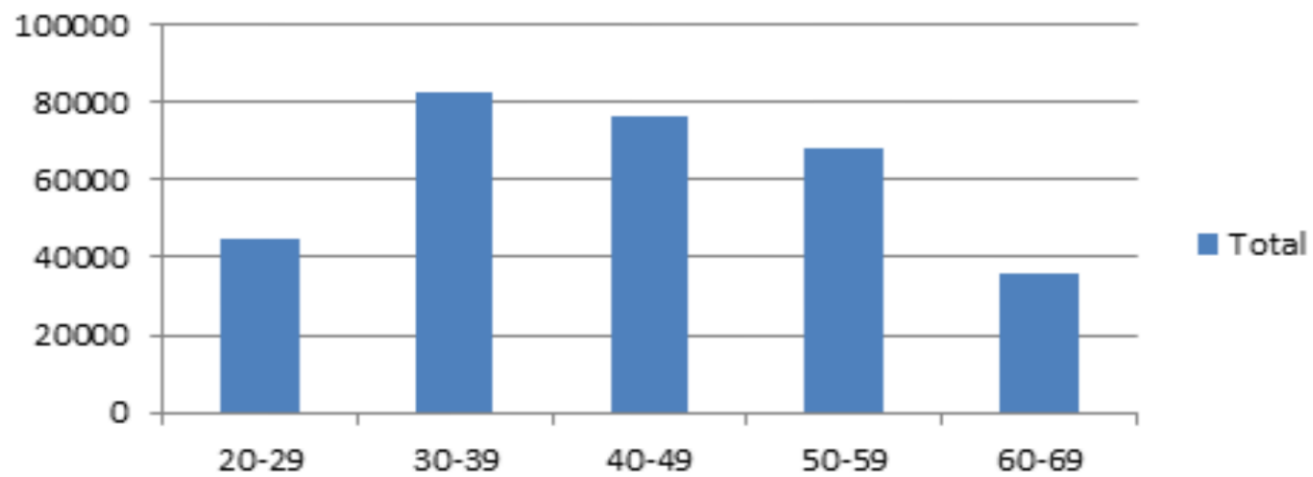
## REGION RATING



REGION\_RATING\_CLIENT ↕

Count of CLIENT'S AGE

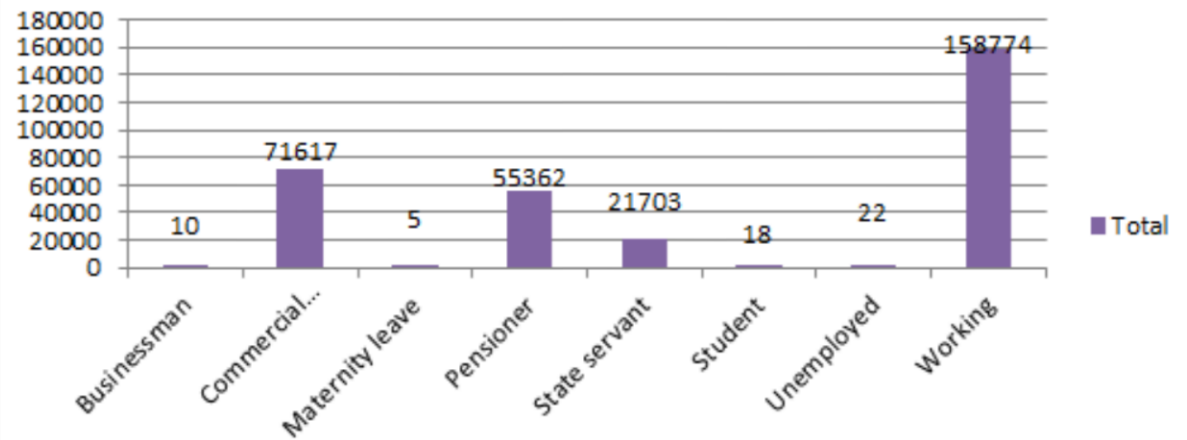
## AGE



CLIENT'S AGE ↕

Count of NAME\_INCOME\_TYPE

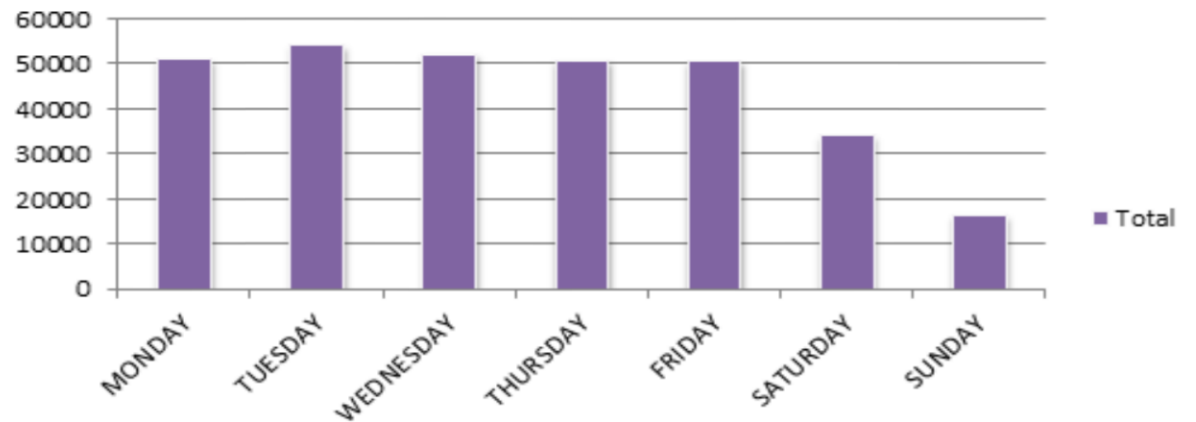
## Income Type



NAME\_INCOME\_TYPE

Count of WEEKDAY\_APPR\_PROCESS\_START

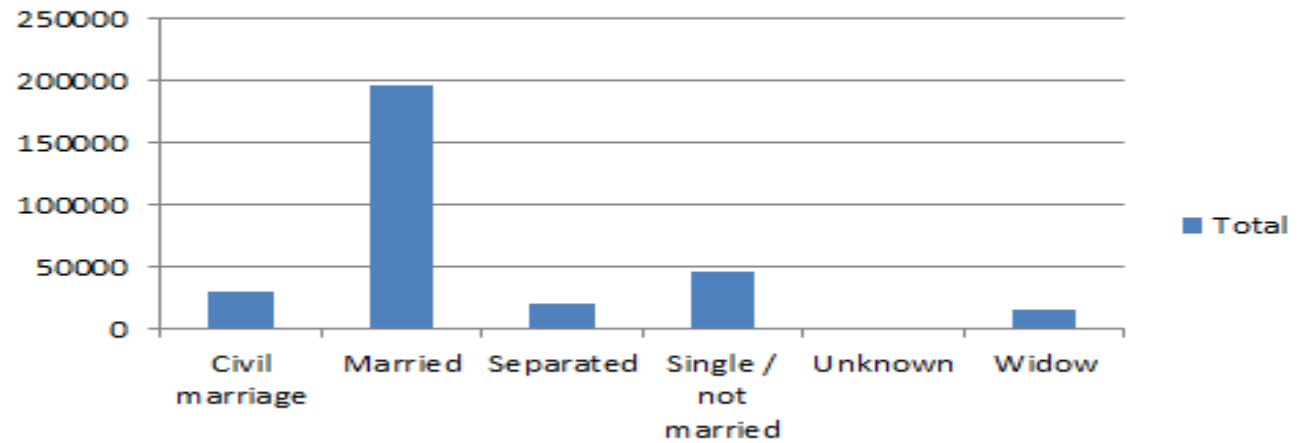
## WEEKDAY\_APPR\_PROCESS\_START



WEEKDAY\_APPR\_PROCESS\_START

Count of NAME\_FAMILY\_STATUS

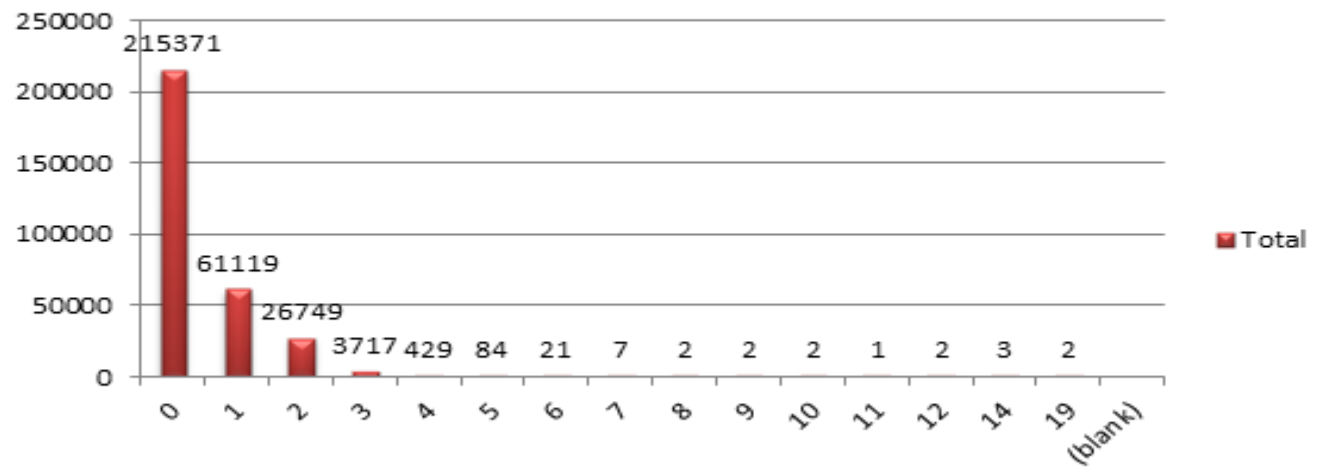
## Applicant Status



NAME\_FAMILY\_STATUS ▼

Count of CNT\_CHILDREN

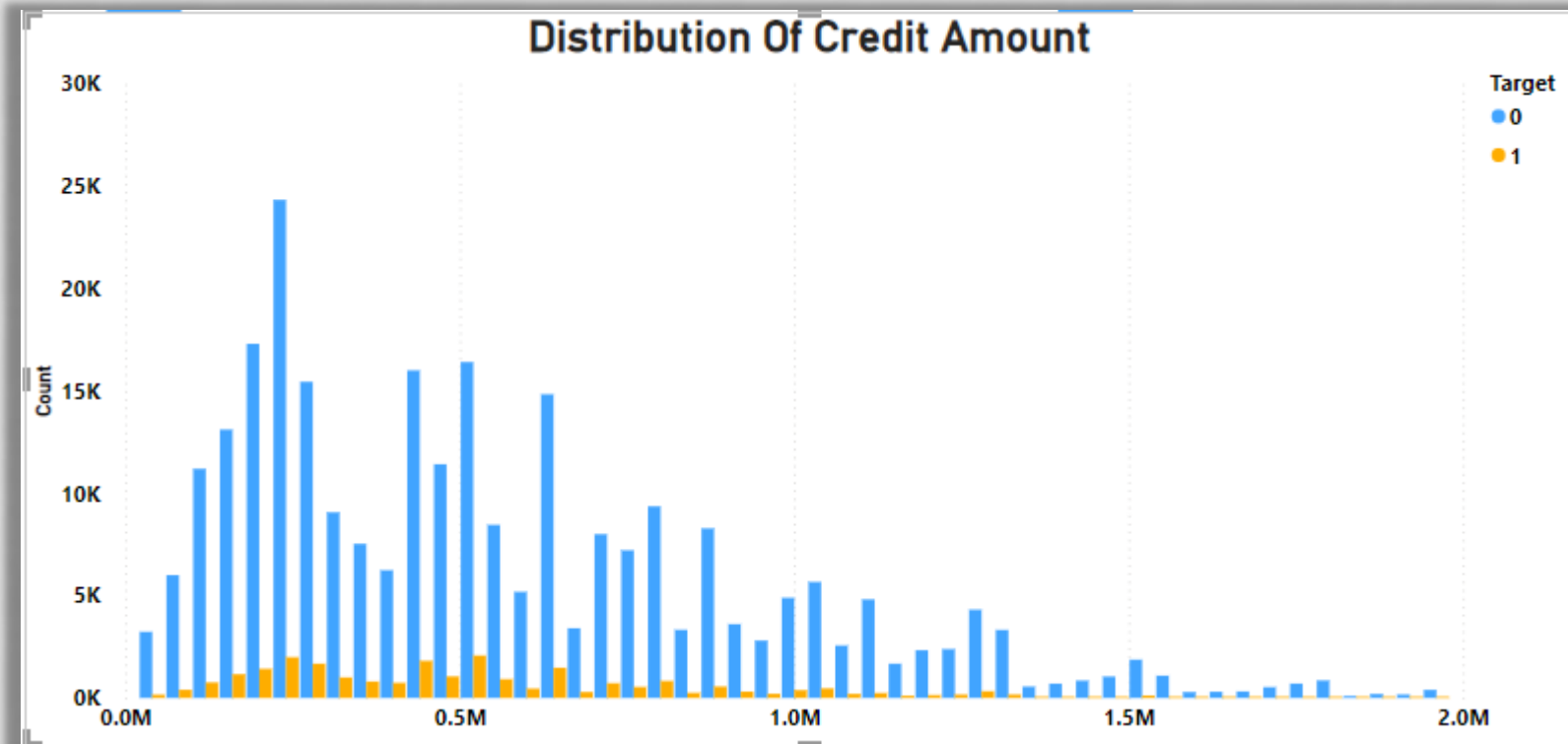
## Count Children



CNT\_CHILDREN ▼

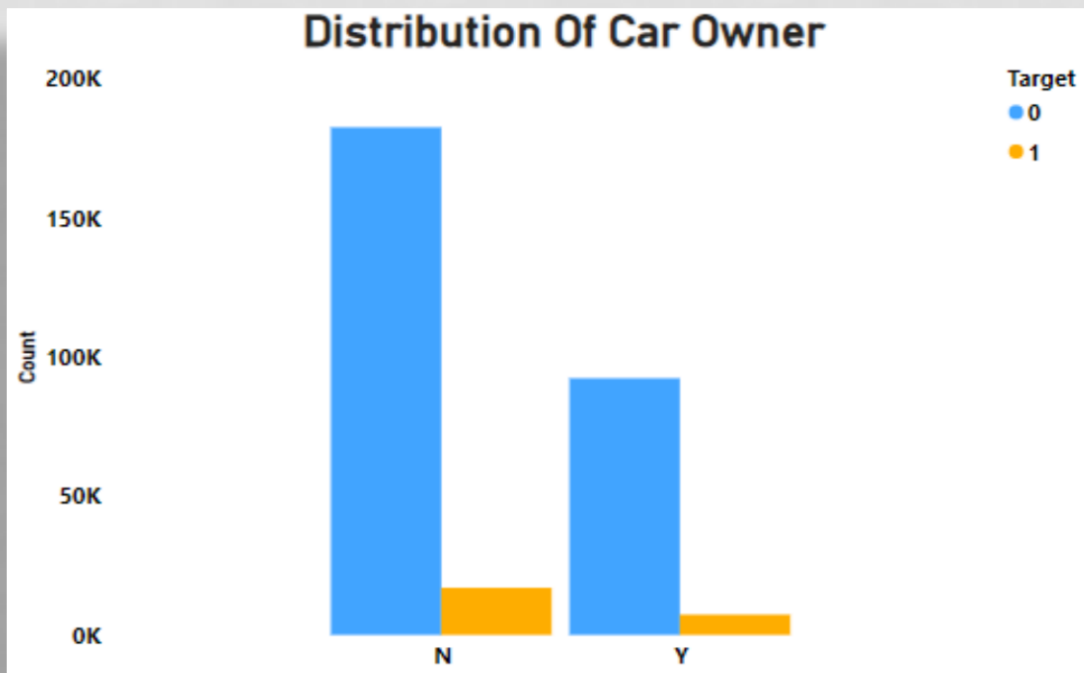
## Univariate Analysis:

So for univariate analysis I used power BI tool I drag the credit amount into x-axis and drag target variable into y-axis as well as into legend.



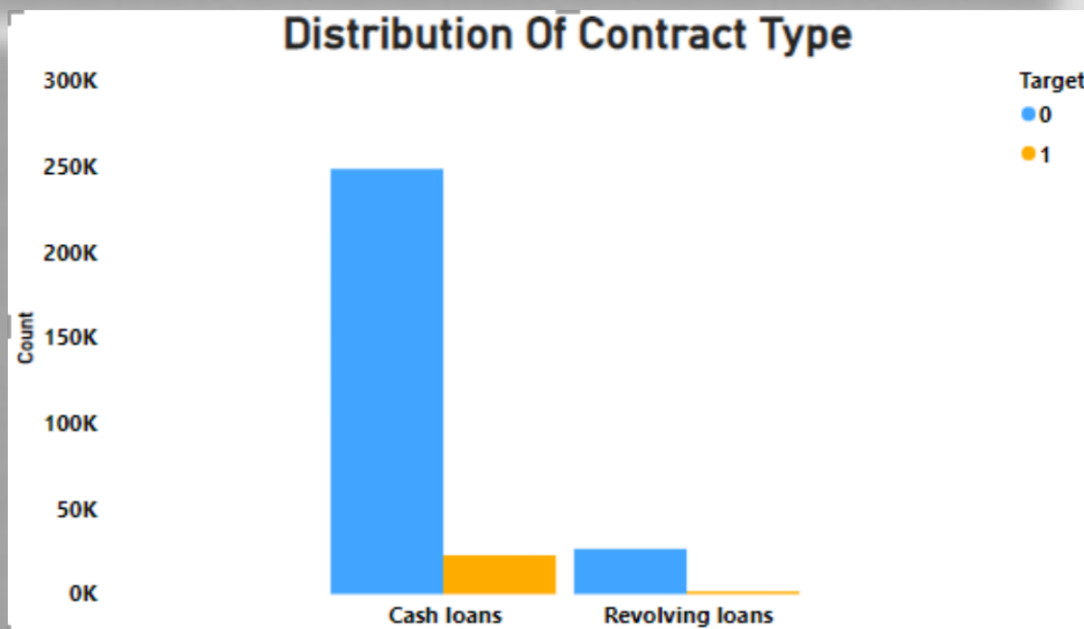
Clients with credit range lying between 0.2M-0.6M are capable of paying loan bank. Least number of clients lying in income range 1.3M-above are not capable of paying.

1. Its shows that clients without having car ownership have no payment difficulties than clients having car.



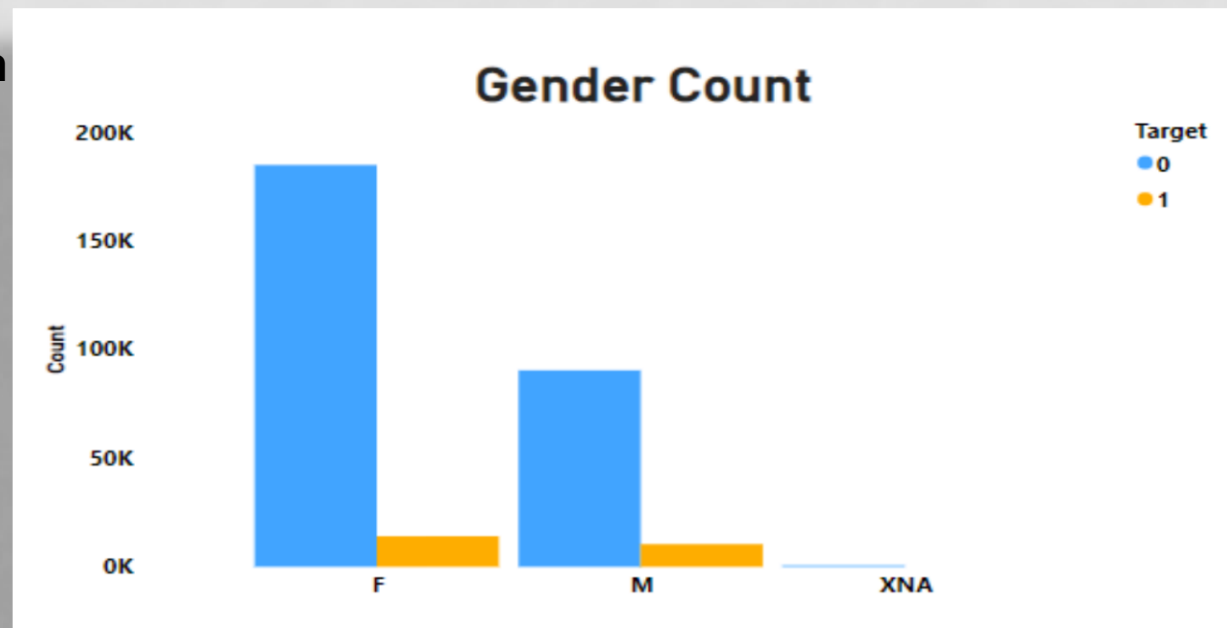
1. Most cash loans applicants have payment difficulties.

2. The same type of loans also have the most applicants with payment difficulties.

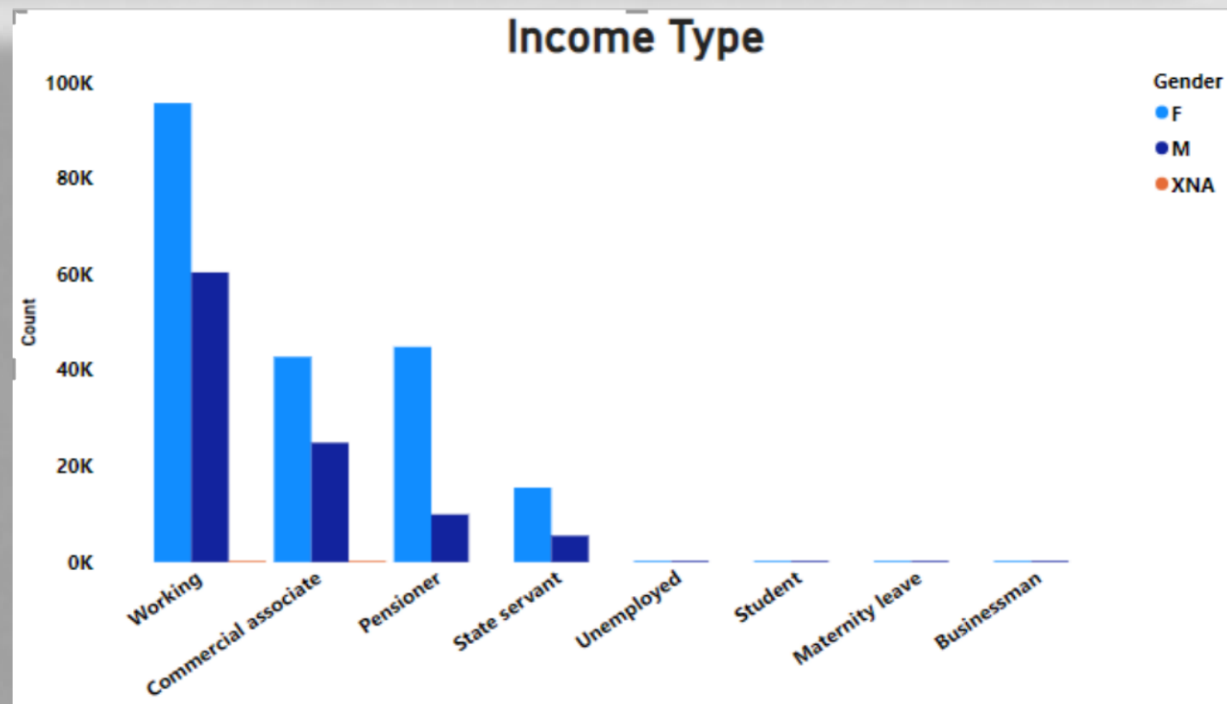




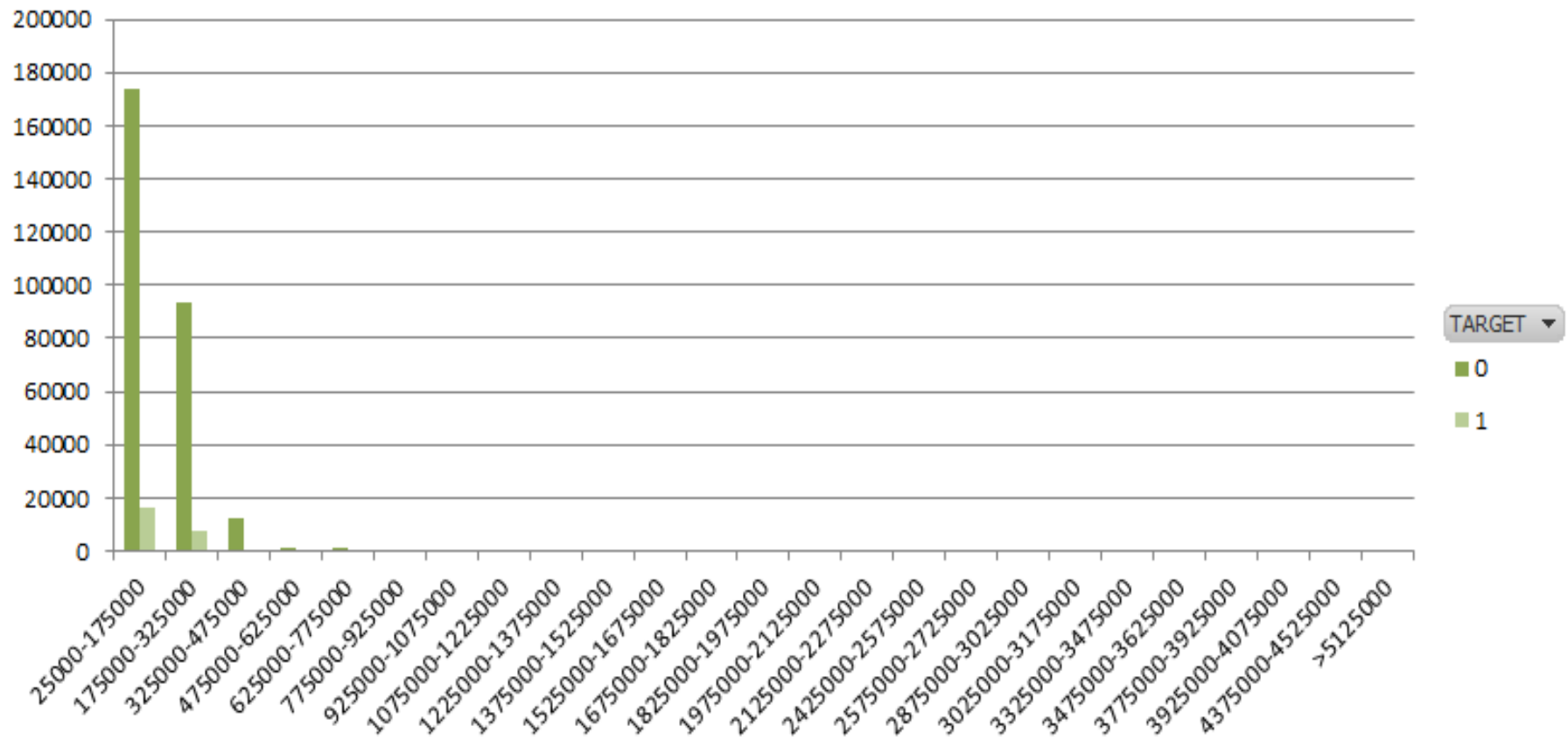
1. It can be seen that both number of males and females are same for having payment difficulties
2. It can also be seen that number of females is more than males for not having payment difficulties.



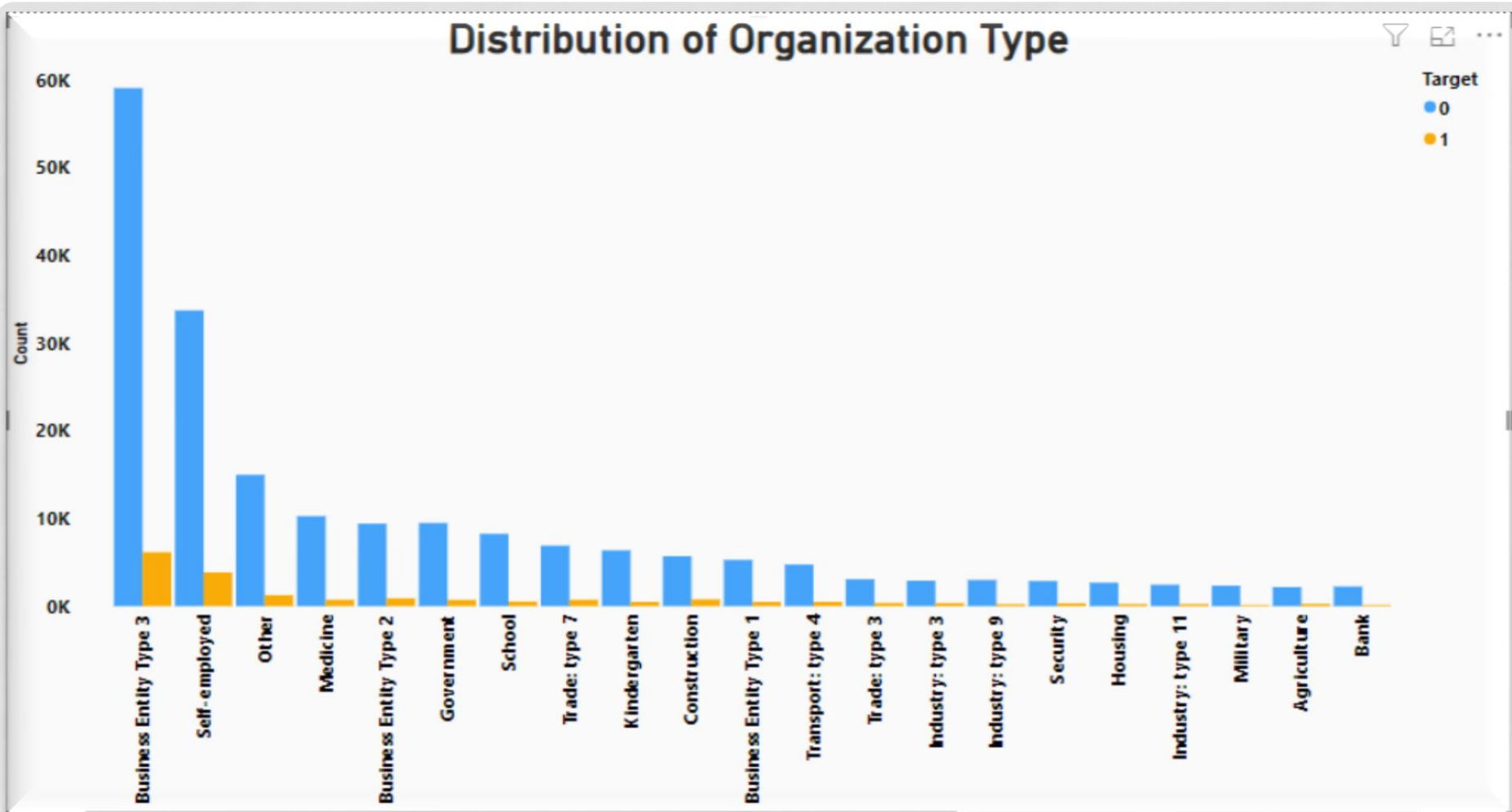
1. We can conclude that most clients who fall in working type income category are applying for loans.



# Income w.r.t Target



From the above plot we can infer that the maximum number of clients without payment difficulties lie in the income range 25K-475k. but maximum number of client with payment difficulties lie in range 25k-175k.

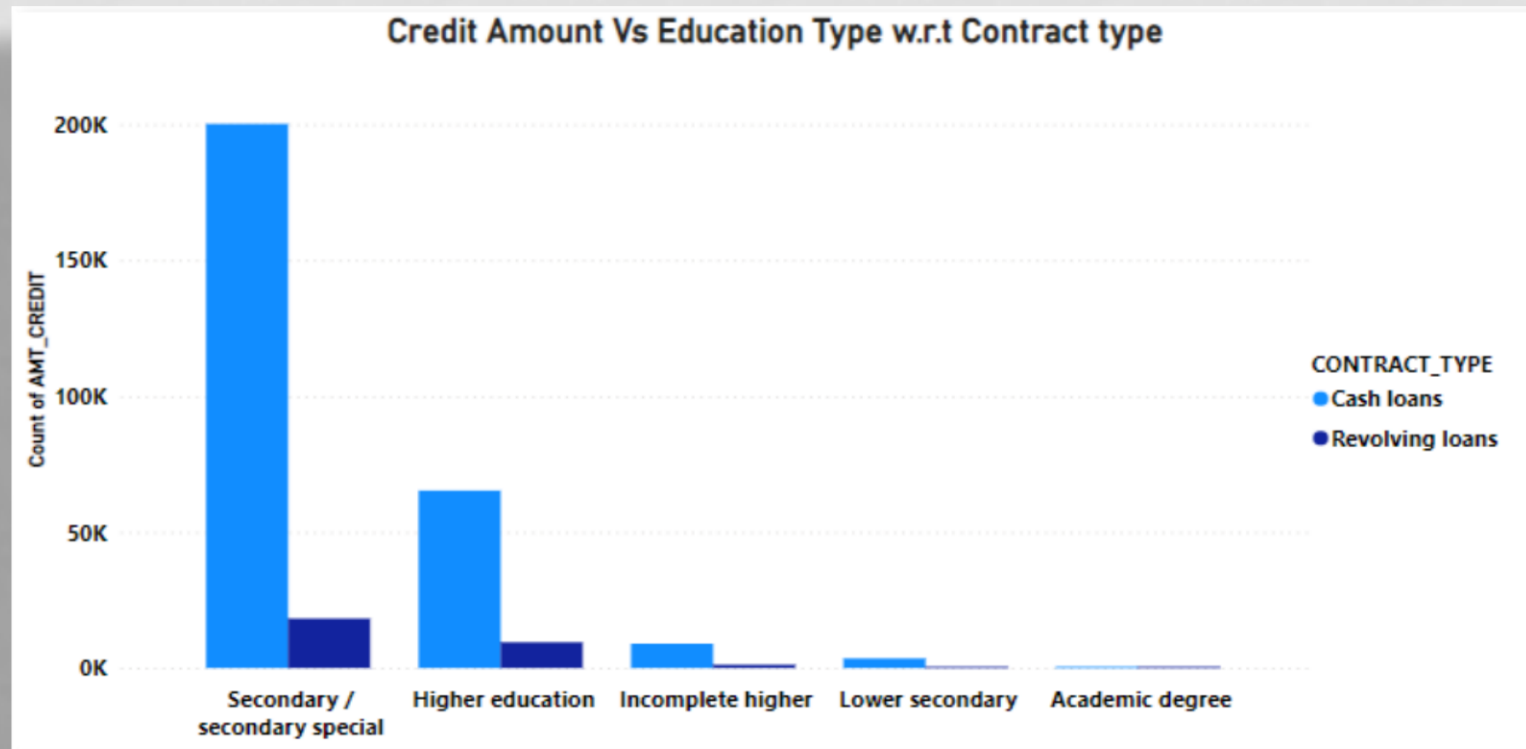


**Most clients with no payment difficulties lie in organization type named Business Entity Type 3.**

**And as you can see that after Medicine organization type all are having least number of payment difficulties clients than no payment difficulties .**

# Bivariate Analysis

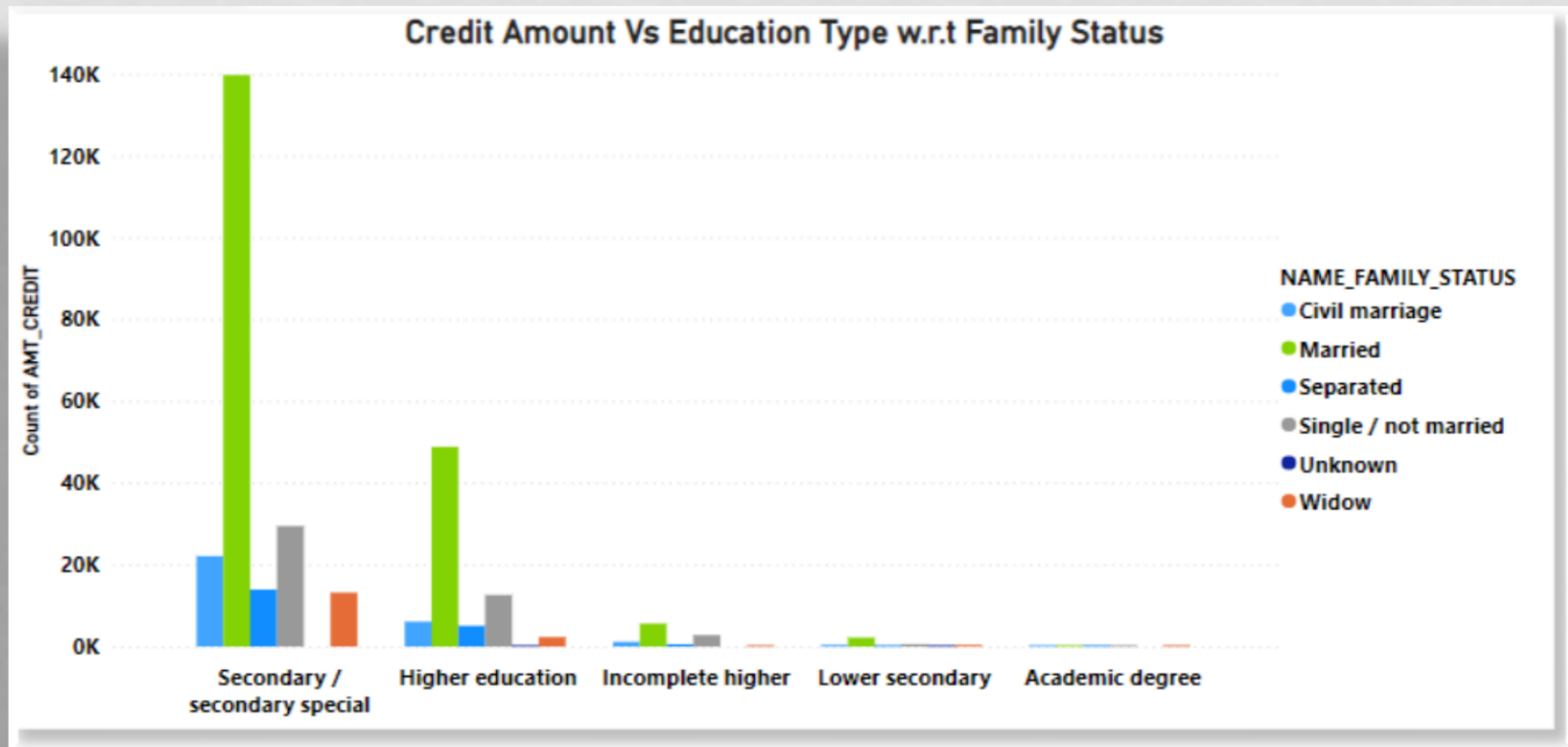
So to bivariate analysis I drag the Credit amount into value, contract type into legend and education type into x-axis.



1. Education status of higher education and secondary/secondary special have most clients for contract type cash loans

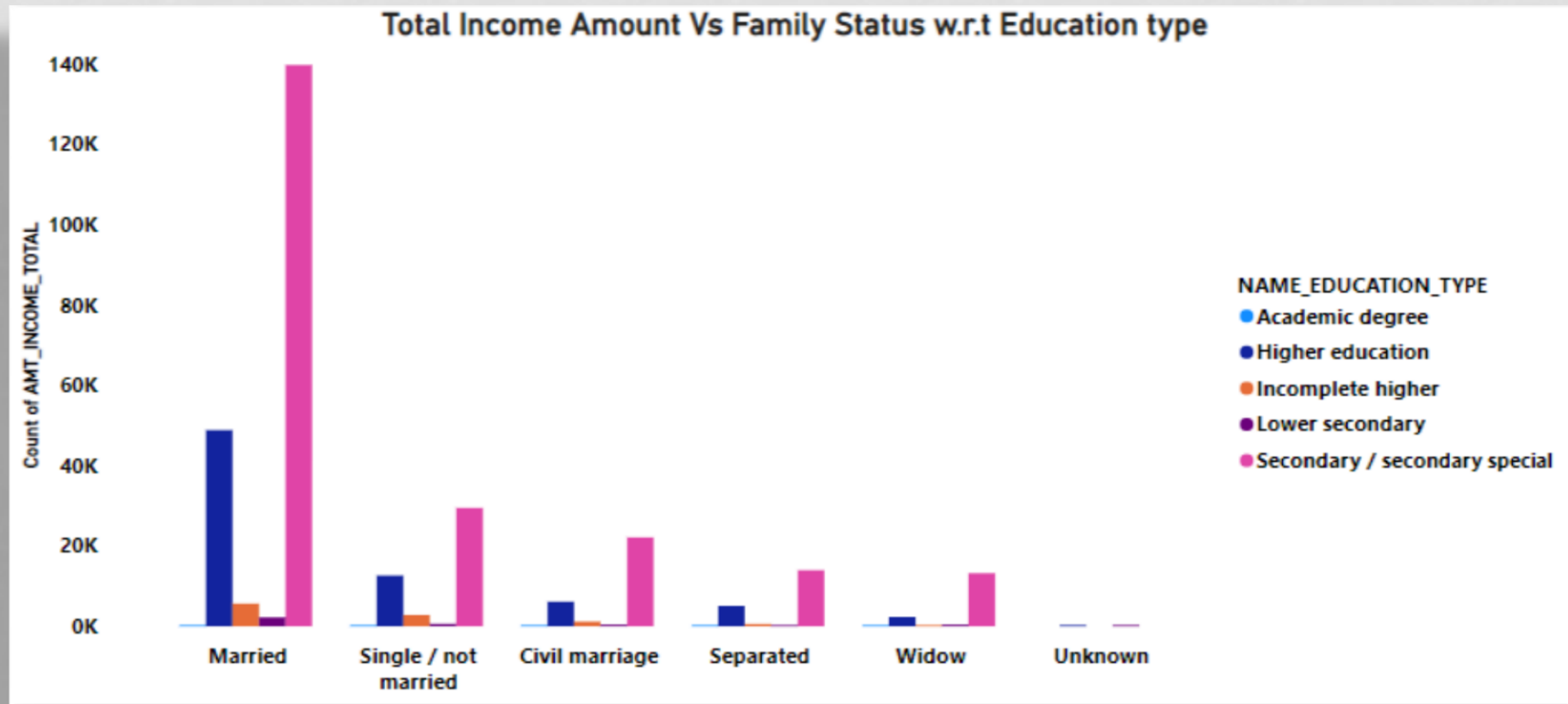
2. Most number of clients applying for revolving loans are in education status also are secondary/secondary special and higher education than other education type.

I drag the Credit amount into value, Education type into legend and family status into x-axis.



1. Family status of married of Secondary/secondary special education are having higher number of credits than others, then family status of civil marriage and separated of secondary/ secondary special education type having higher number of credits than others.

I drag the Credit amount into value, family status into legend and education type into x-axis.



1. Family status of 'civil marriage', 'marriage' and 'Single/not married' of secondary/ secondary special are having higher number of income than others.

2. All family status of academic degree education type having minority of income rate. It means bank having less loan customers in education type academic degree. As we see earlier in contract type.

# Correlations (target 0)

The below chart provides correlation of different variables for clients with no payment difficulties.

For Correlation I filter the target variable to 0 and selected the all variables then clicked the data bar then click on Data analysis.

As we can see top 10 correlations.

1. Family member & children count
2. Goods price& credit
3. Annuity & credit
4. Goods price& annuity
5. Annuity &total income
6. Credit &total income
7. Goods price& total income
8. Day birth & children count
9. Day registration & days birth
10. Days publish & days birth

	CNT_CHIL DREN	AMT_INC OME_TOT AL	AMT_CRE DIT	AMT_AN NUITY	AMT_GO ODS_PRIC E	DAYS_BIR TH	DAYS_EM PLOYED	DAYS_RE GISTRATI ON	DAYS_ID_ PUBLISH	CNT_FAM _MEMBE RS
CNT_CHILDREN	1									
AMT_INCOME_TOTAL	0.0274	1								
AMT_CREDIT	0.00308	0.3428	1							
AMT_ANNUITY	0.02091	0.41895	0.77131	1						
AMT_GOODS_PRICE	-0.00052	0.34946	0.98725	0.77669	1					
DAYS_BIRTH	0.33697	0.06261	-0.04738	0.01226	-0.04456	1				
DAYS_EMPLOYED	-0.24336	-0.14125	-0.07251	-0.10642	-0.07105	-0.61805	1			
DAYS_REGISTRATION	0.18579	0.06494	0.01348	0.03944	0.01592	0.33315	-0.21019	1		
DAYS_ID_PUBLISH	-0.02875	0.0229	-0.00146	0.01411	-0.00365	0.27131	-0.27429	0.10024	1	
CNT_FAM_MEMBERS	0.87857	0.03426	0.06454	0.07579	0.06281	0.28582	-0.23741	0.17563	-0.02046	1

# Correlation ( target 1)

Same as before I filter the target variable to 1 and selected the variables then clicked the data bar then click on Data analysis.

So these are the top 10 correlations

1. Family mem & children count
2. Goods price & credit
3. Annuity & credit
4. Goods price & annuity
5. Days birth & family member
6. Days registration & days birth
7. Days publish & days birth
8. Days registration & family member
9. Days registration & children count
10. Days id publish & days registration

	CNT_CHIL DREN	AMT_INC OME_TO TAL	AMT_CRE DIT	AMT_AN NUITY	AMT_GO ODS_PRI CE	CNT_FA M_MEM BERS	DAYS_BI RTH	DAYS_E MPLOYE D	DAYS_RE GISTRATI ON	DAYS_ID _PUBLISH
CNT_CHILDR EN	1									
AMT_INCO ME_TOTAL	0.004796	1								
AMT_CREDIT	-0.00167	0.038131	1							
AMT_ANNUITY	0.031257	0.046421	0.752195	1						
AMT_GOOD S_PRICE	-0.00811	0.037583	0.983103	0.752699	1					
CNT_FAM_ MEMBERS	0.885484	0.006654	0.051224	0.075711	0.047388	1				
DAYS_BIRTH	0.259109	0.003096	-0.13532	-0.0143	-0.13581	0.203267	1			
DAYS_EMPL OYED	-0.19194	-0.01498	-0.00097	-0.08255	0.003587	-0.18656	-0.5751	1		
DAYS_REGIS TRATION	0.149154	0.000158	-0.02585	0.034279	-0.02568	0.145828	0.289114	-0.18893	1	
DAYS_ID_PU BLISH	-0.0323	-0.00421	-0.05233	-0.01677	-0.05609	-0.03178	0.252863	-0.22647	0.096833	1

The above chart provides correlation of different variables for clients with payment difficulties



## Insights (previous\_application.csv)

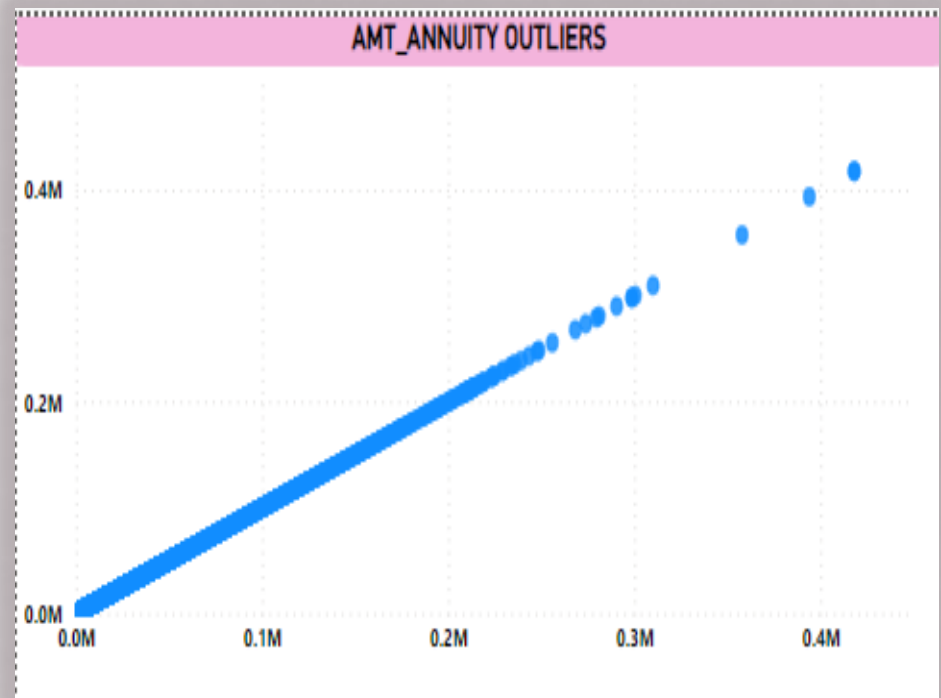
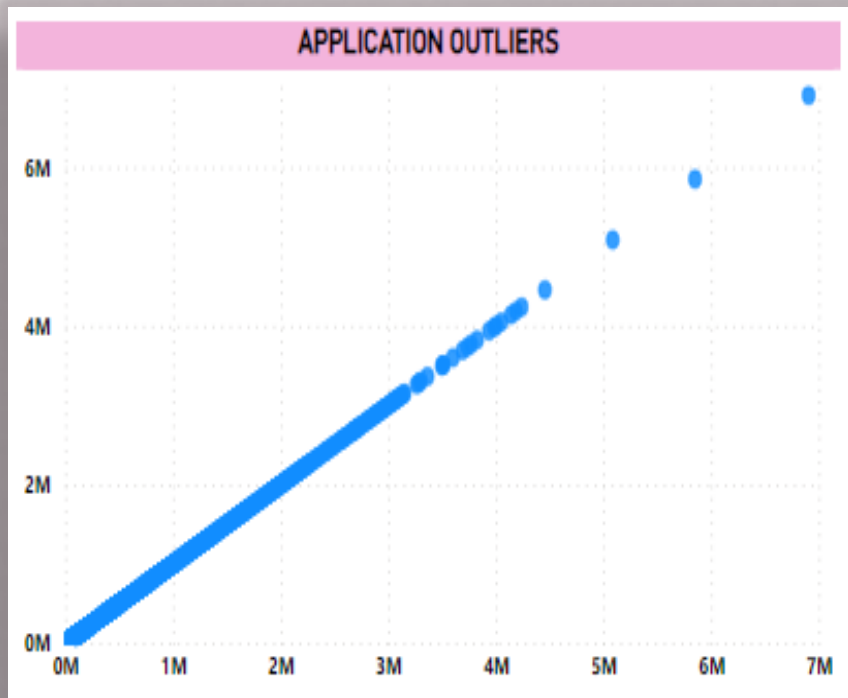
**Data Cleaning:** The dataset contained more than 1.6 million rows and 37 columns of data.

So I used Power BI software because there are more than million rows of data which cannot be import in Excel, after importing the dataset into power BI basically I transform the data by removing unwanted column contained 45% of blank data.

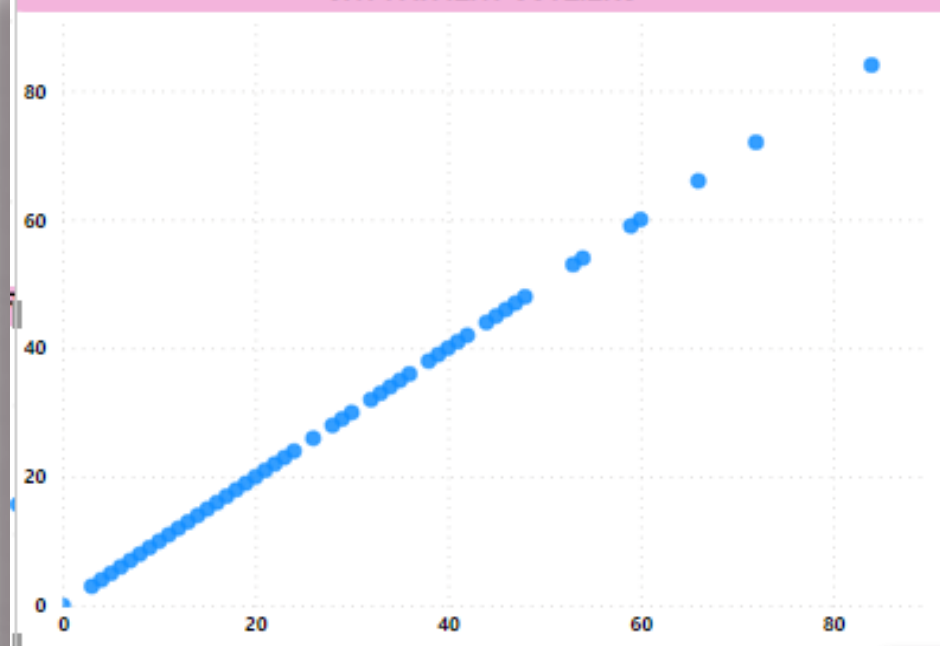
After deleting unwanted columns we got 34 remaining columns.

# Outliers

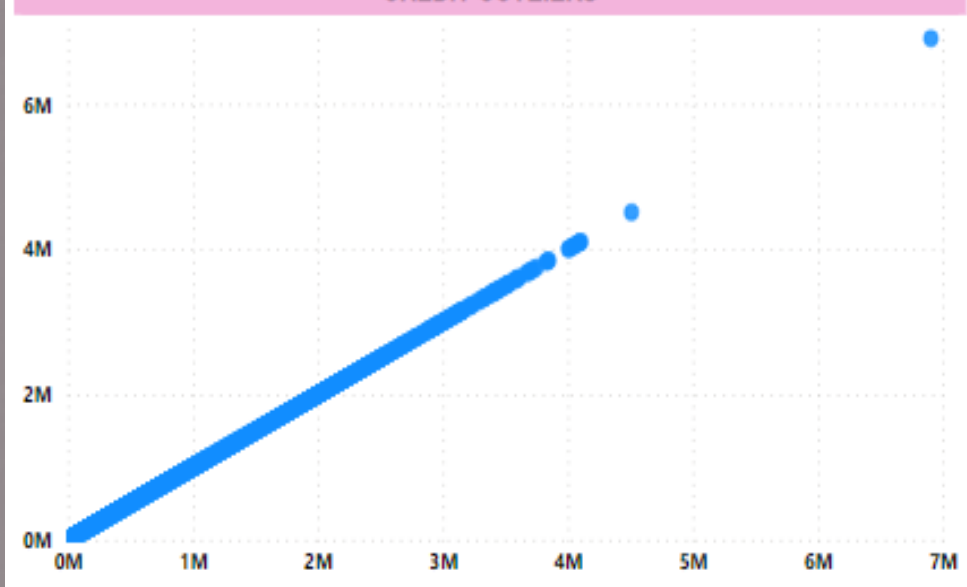
Outliers can only be identified on Numerical variable.  
To find outlier basically I used Scatter plot chart.

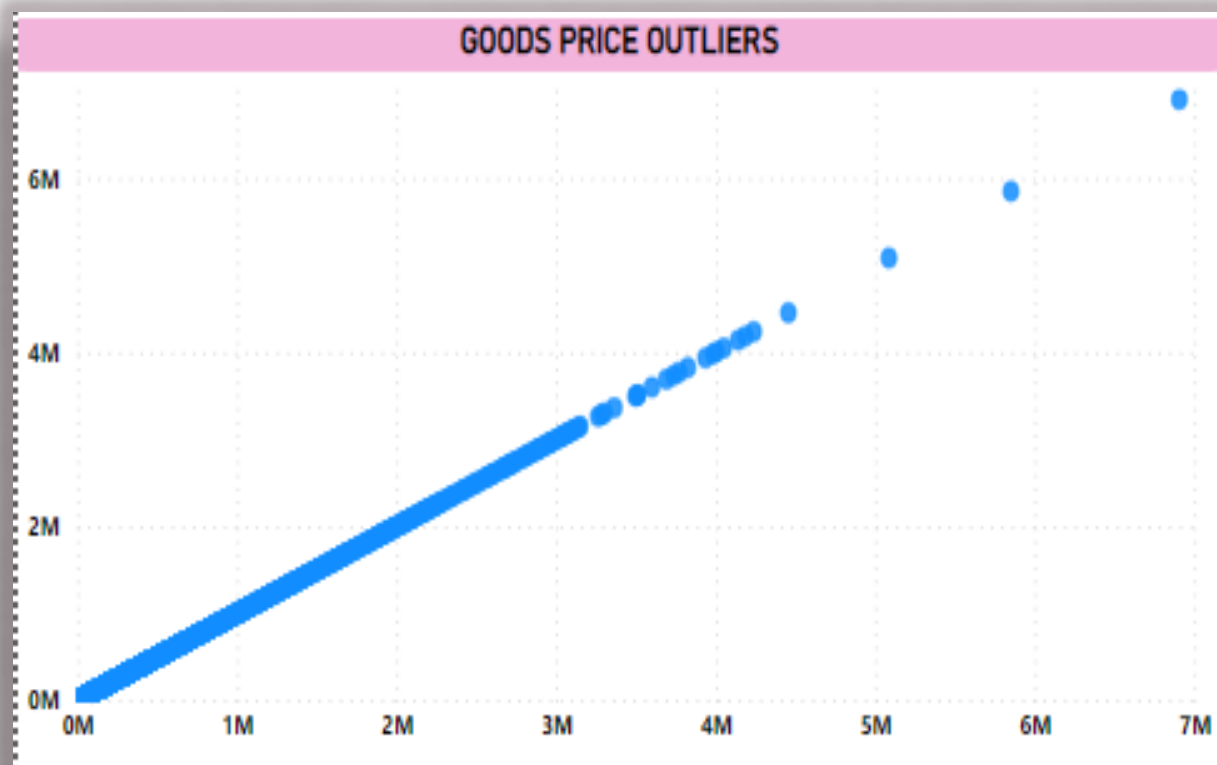


CNT PAYMENT OUTLIERS



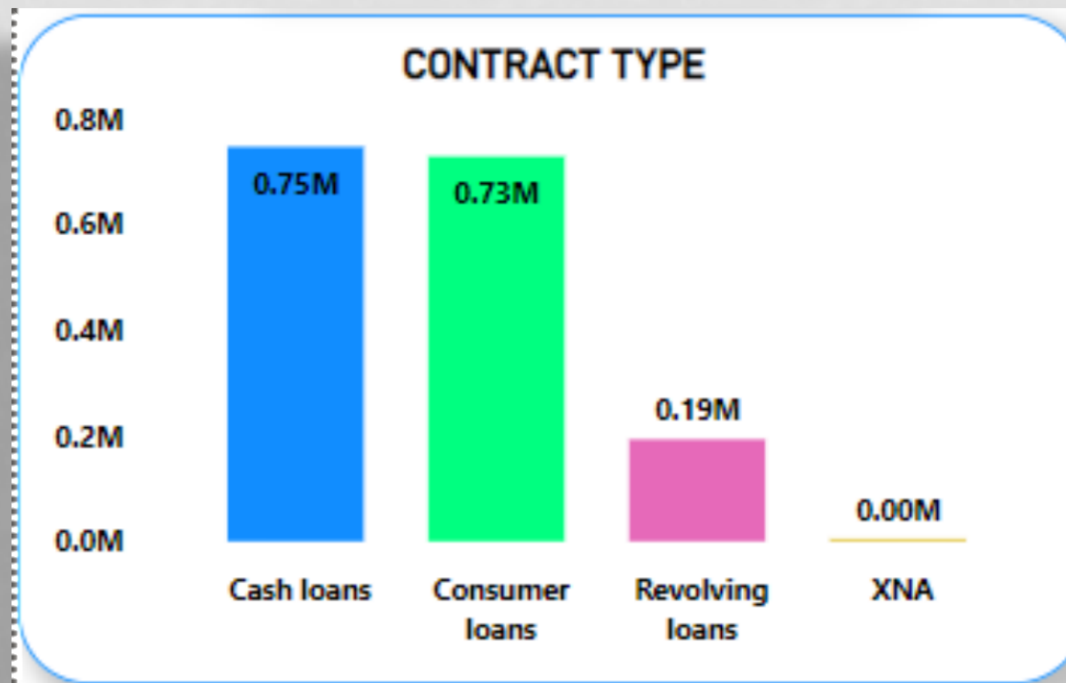
CREDIT OUTLIERS

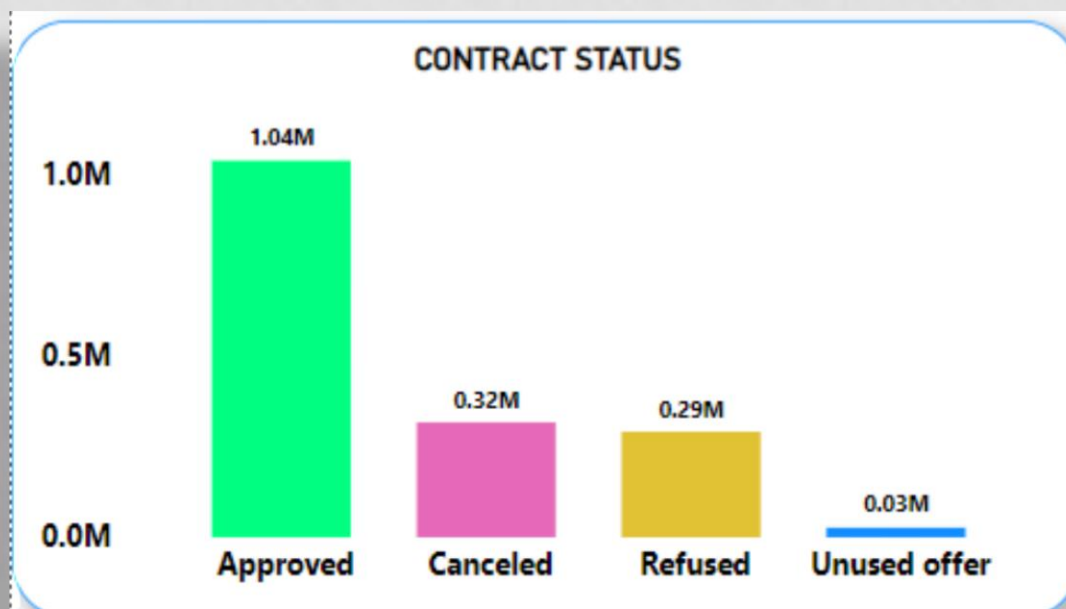
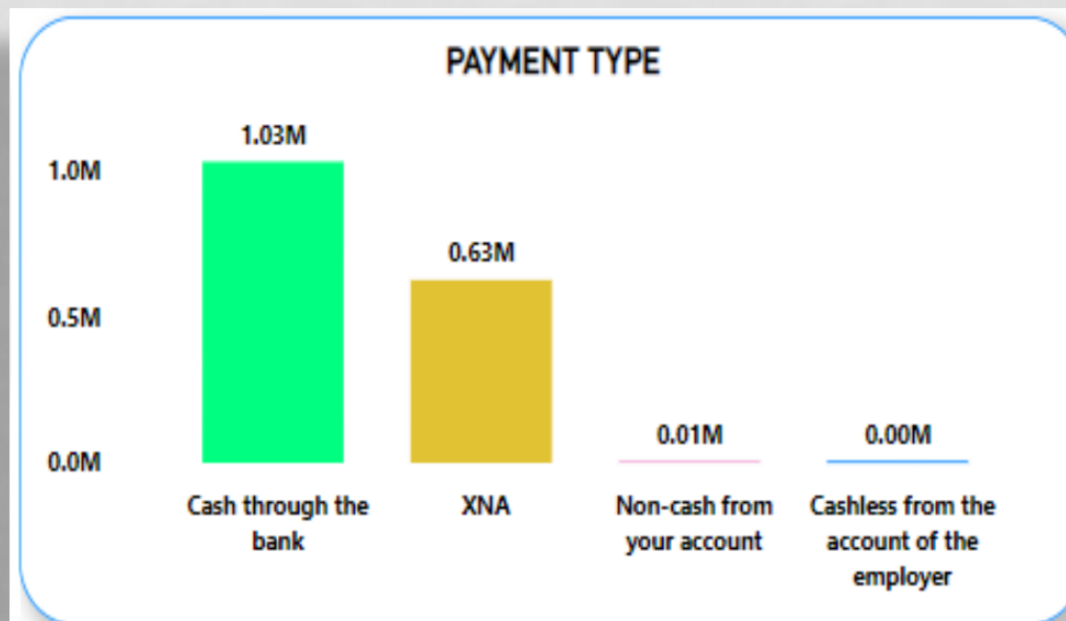


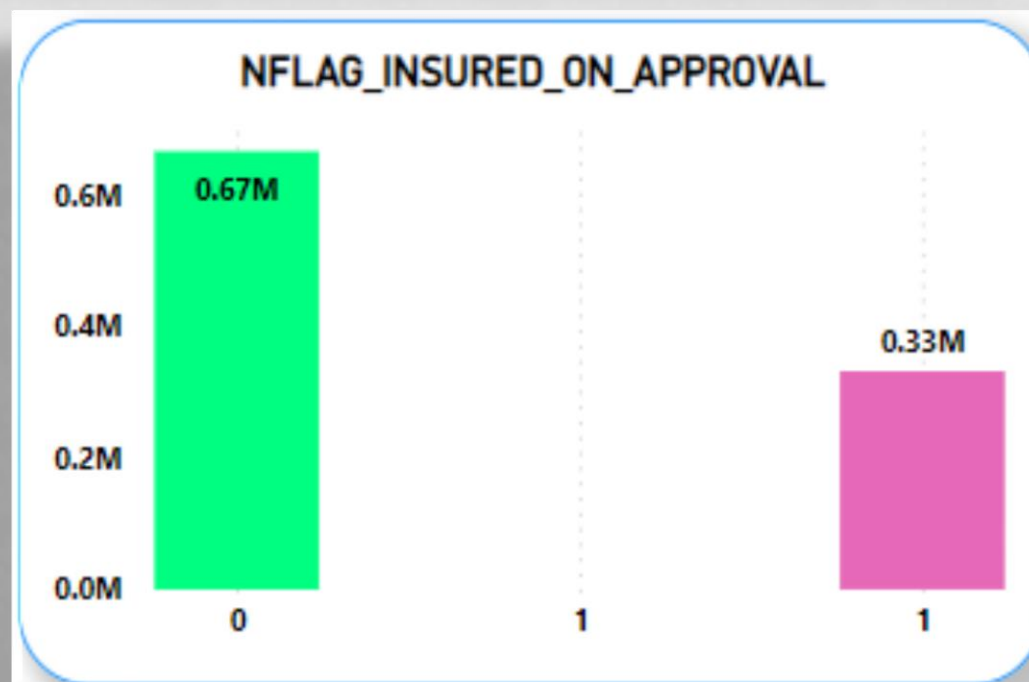
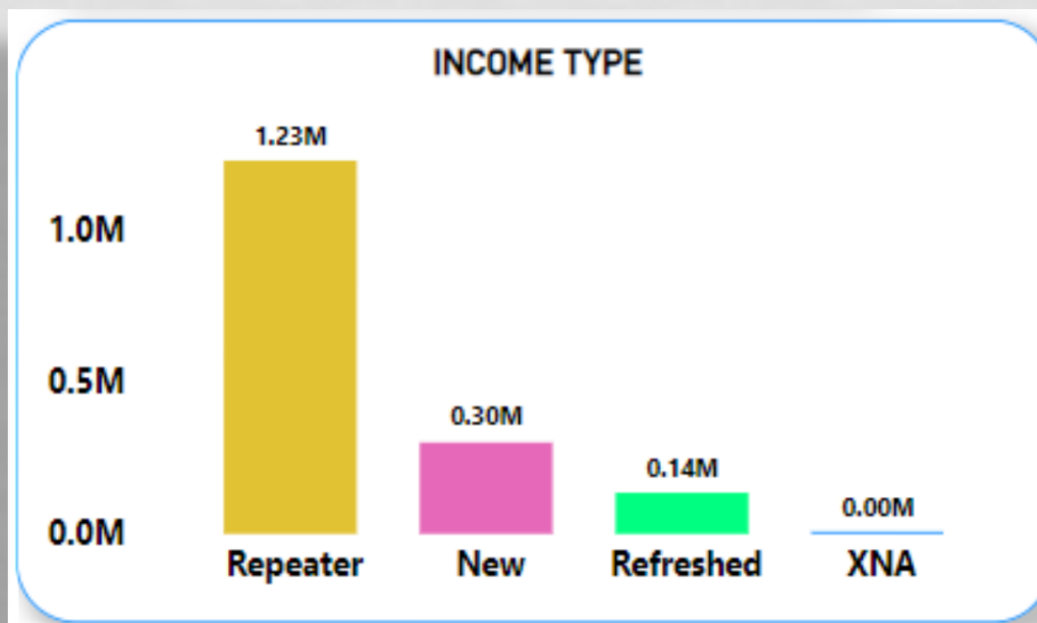


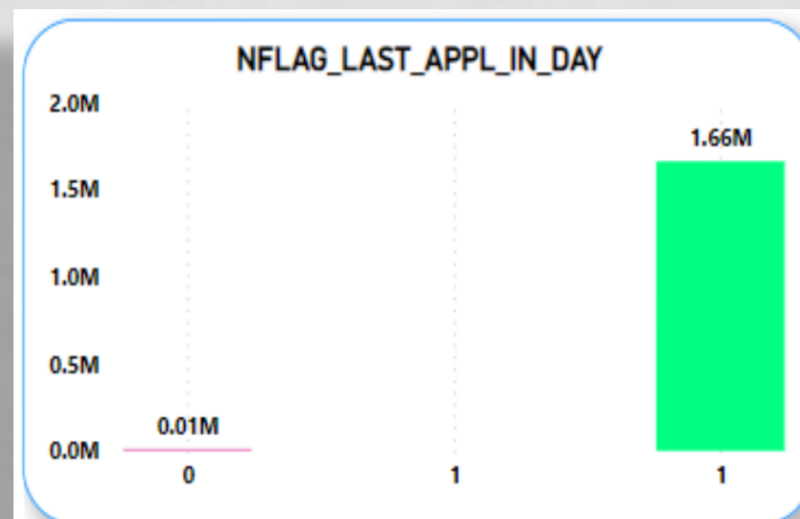
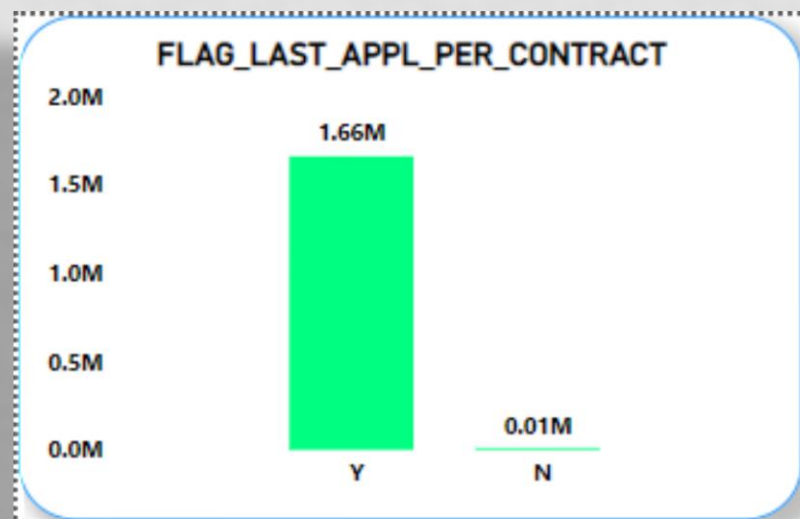
## Data Imbalance:

Below are the columns where data is unevenly distributed





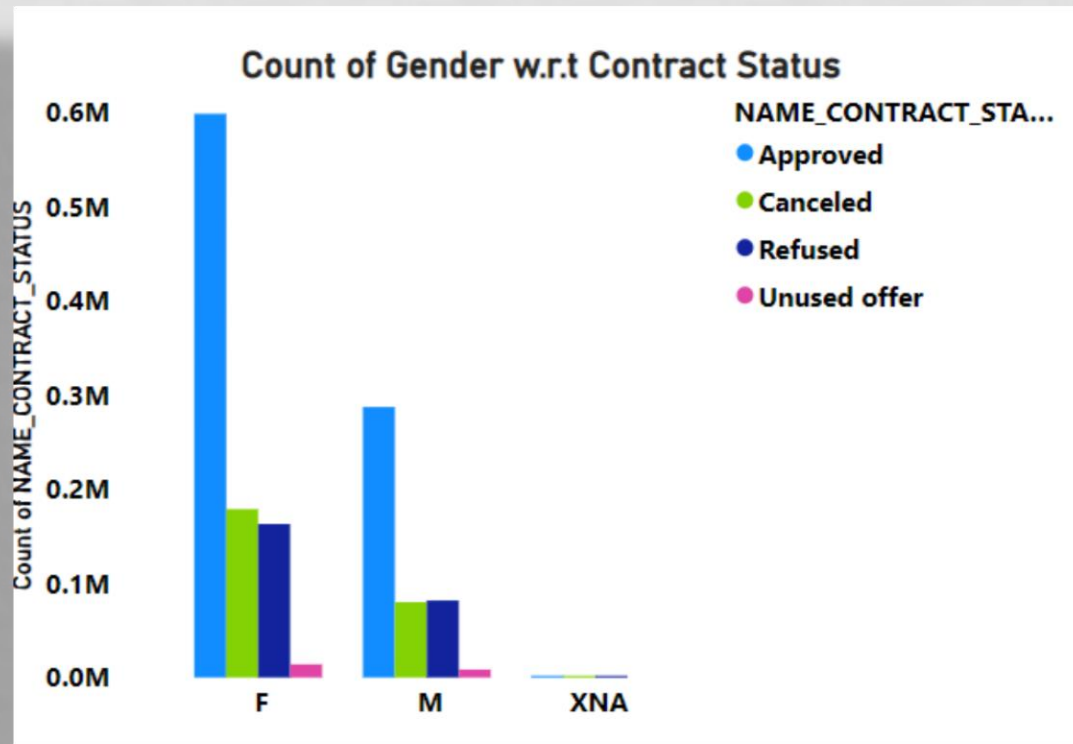




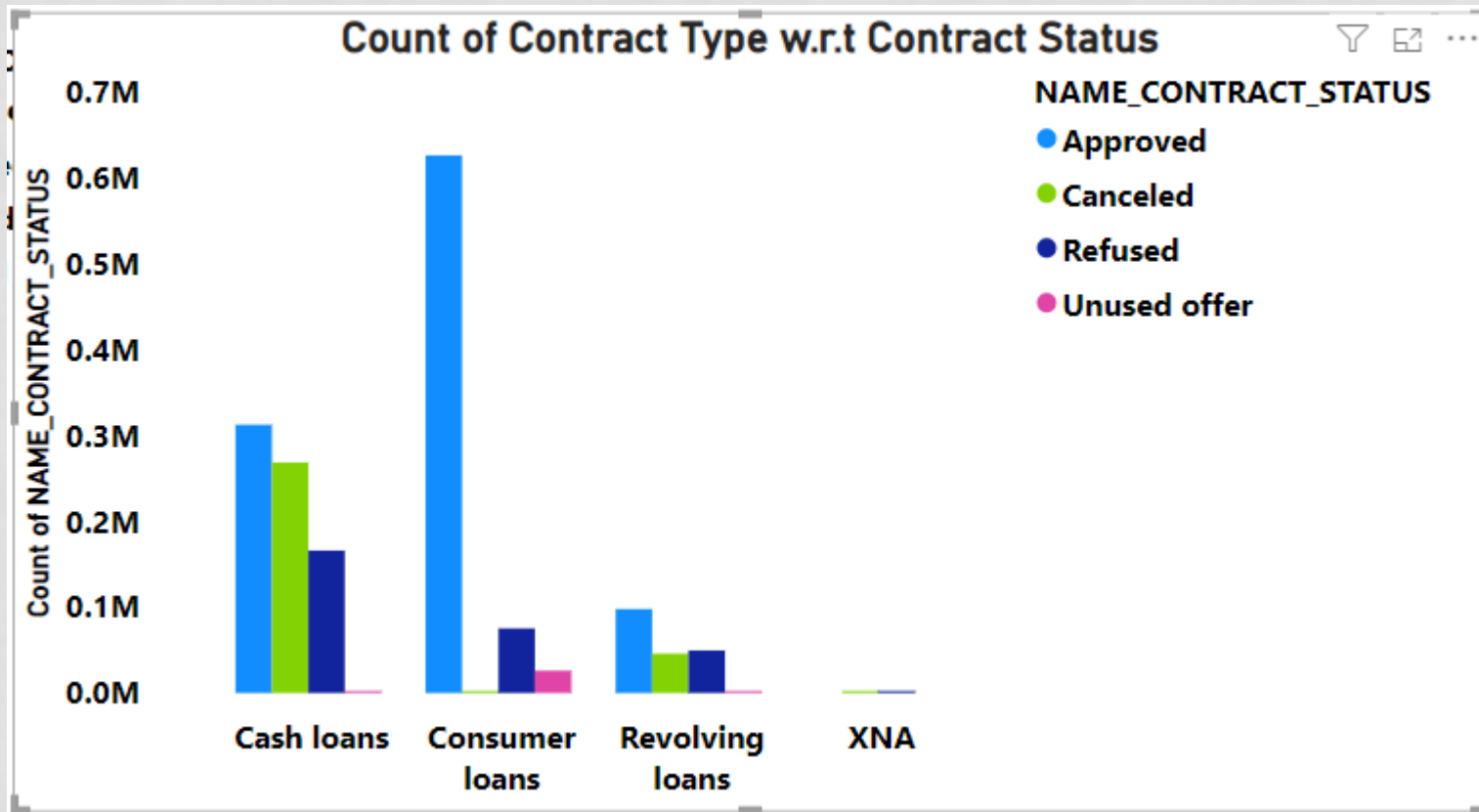


# Univariate Analysis

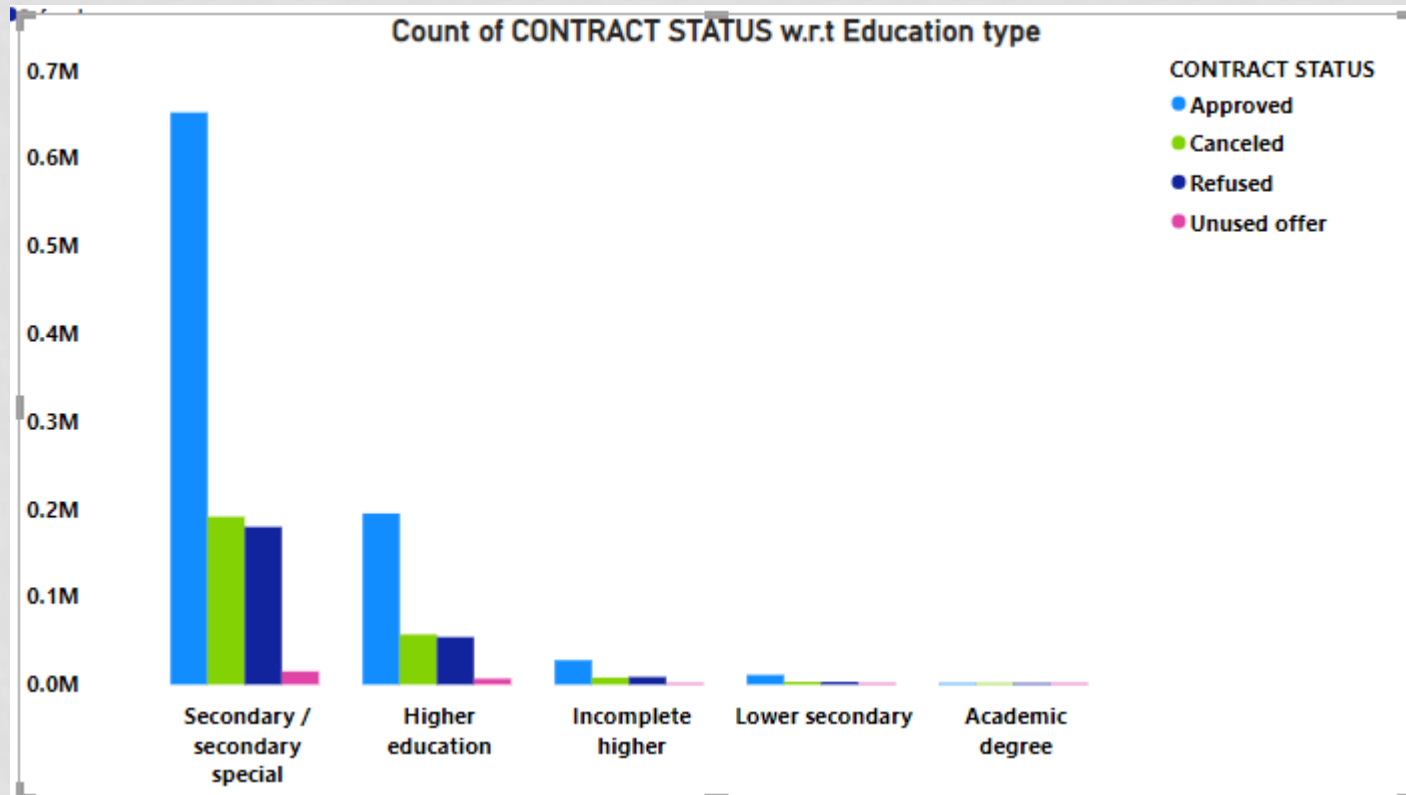
So for univariate analysis I did same as application\_data.csv dataset. And after this I create relationship between both dataset table.



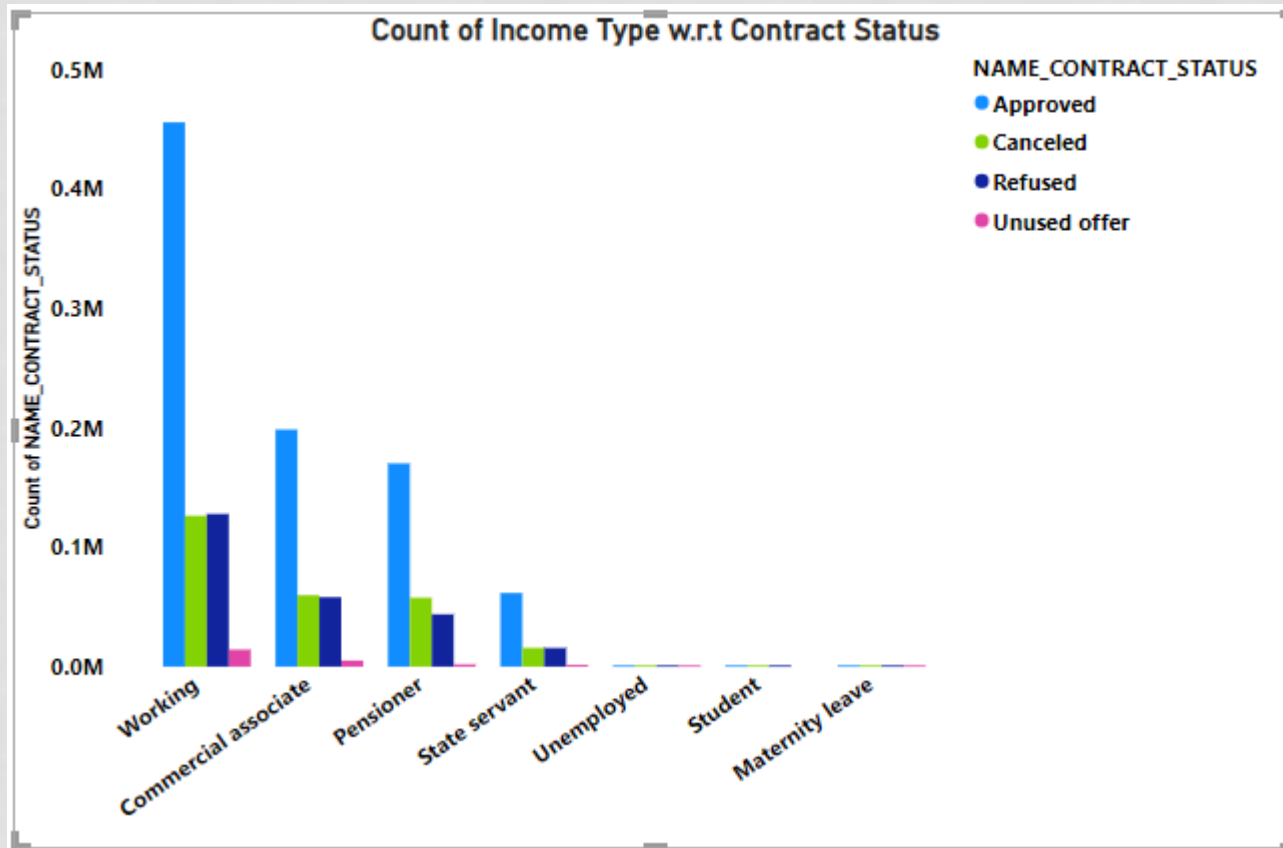
- Maximum approved loans are for female clients
- It can also be observed that male clients use most of the offers of loans as unused offers are very less for male clients than that of female clients.



- Contract type cash loans are maximum in number where all kinds of contract statuses are more than contract type revolving loans
- Consumer loans are most approved loans.

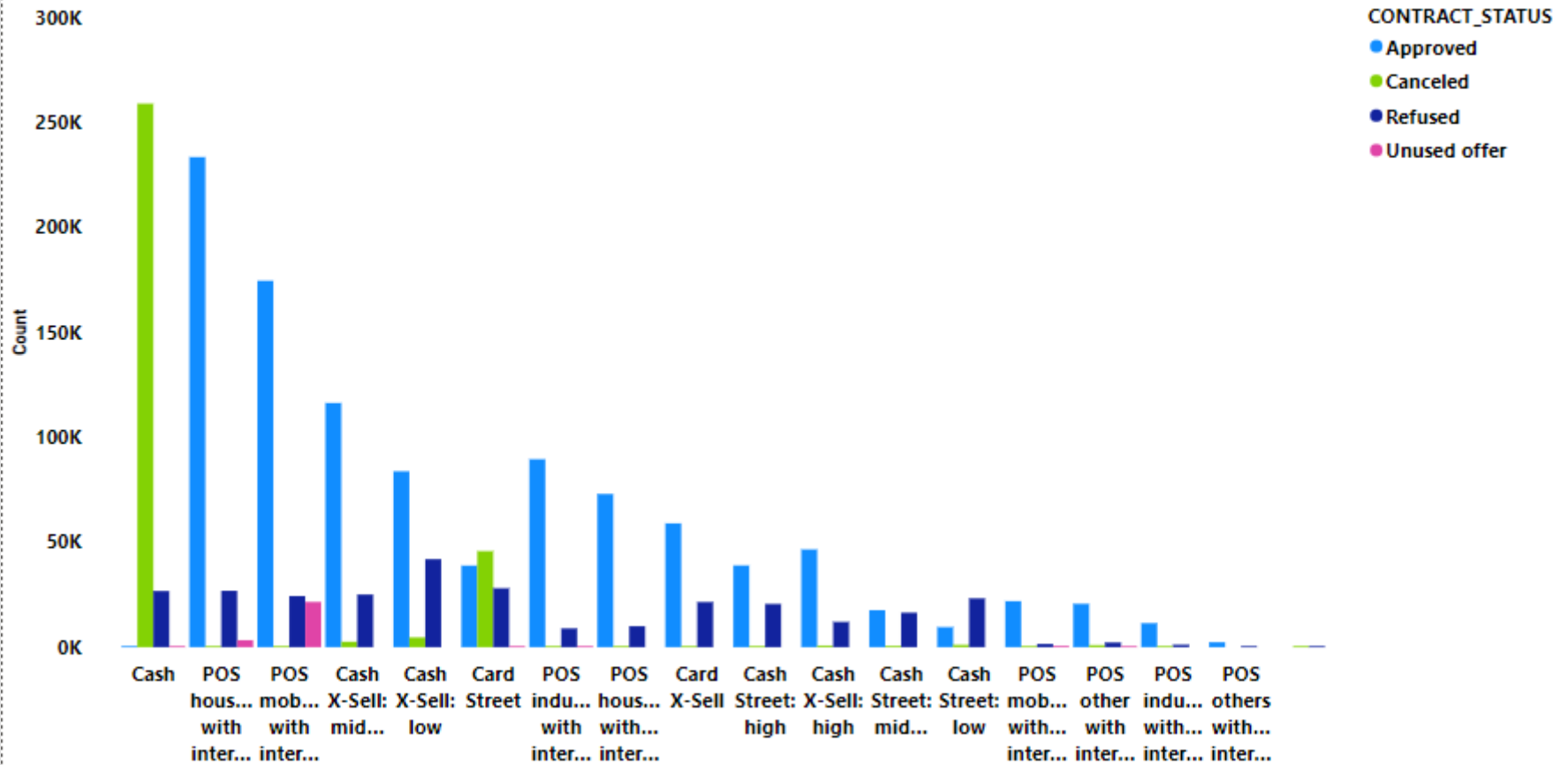


**Most clients with all kinds of education types have approved rate is better than canceled, refused rate.**



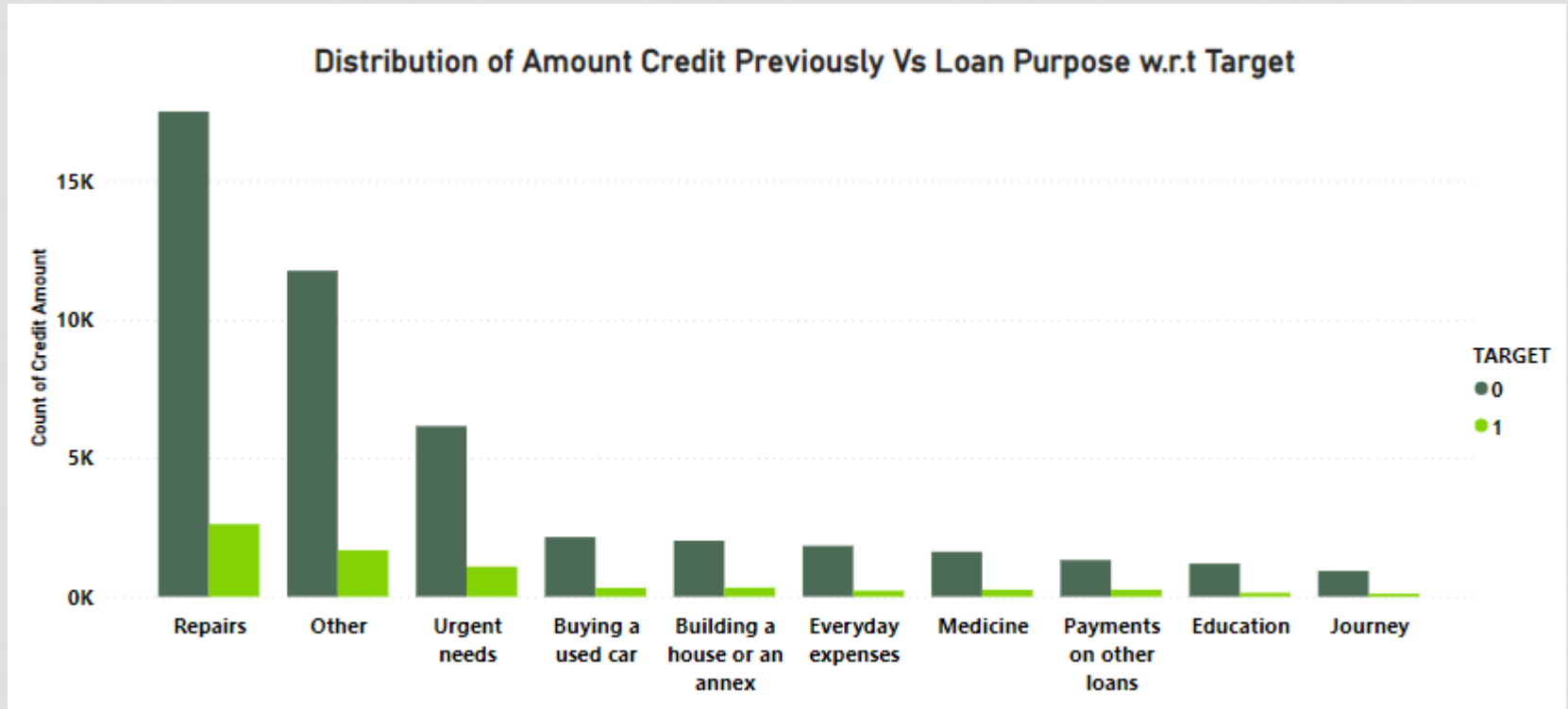
- There are very less rate of getting approved and unused offer for students and for maternity leave income type.
- The number of approved loans for state servants is almost equal to the refusal or canceled loans for Commercial associates.
- Maximum unused offers is by working clients and more approved is by working clients.

## Count of Product Combinations w.r.t Contract Status



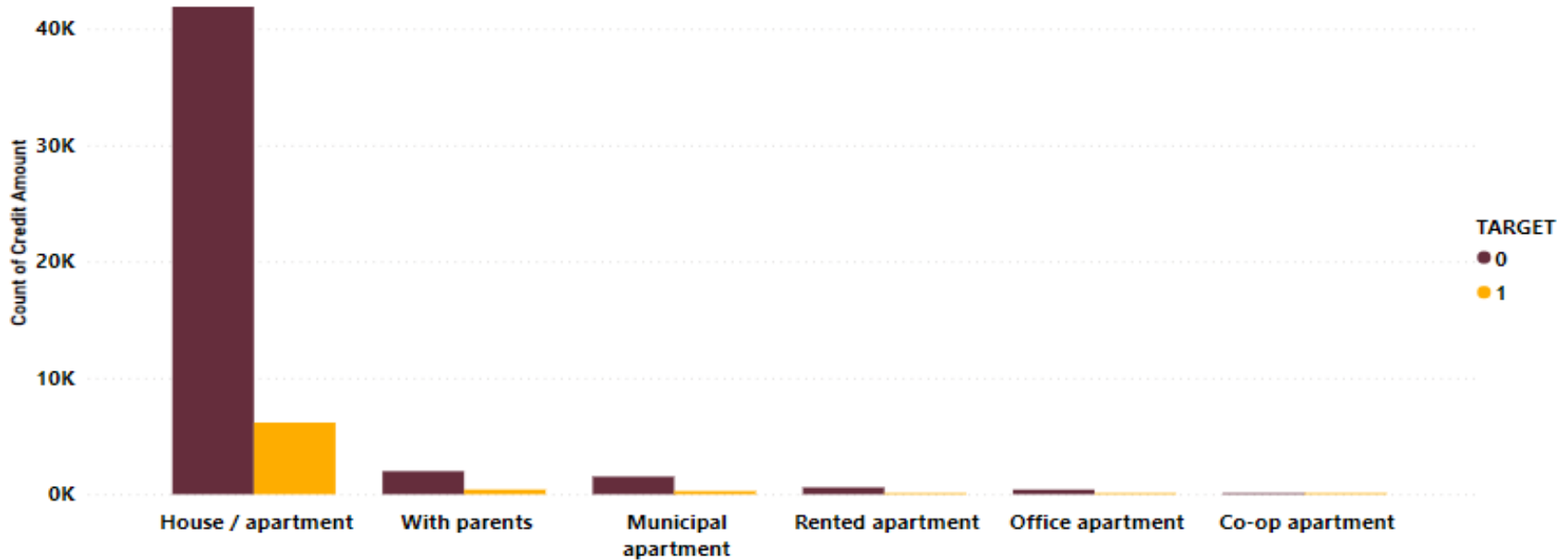
- Most canceled loans are of the product combination, Cash.
- Most refused loans are of product combination Cash X-Sell: low
- Some product combinations are null incase of unused offers as well as cancelled loans
  - POS industry without interest
  - POS others without interest

# Bivariate Analysis for previous\_application.csv



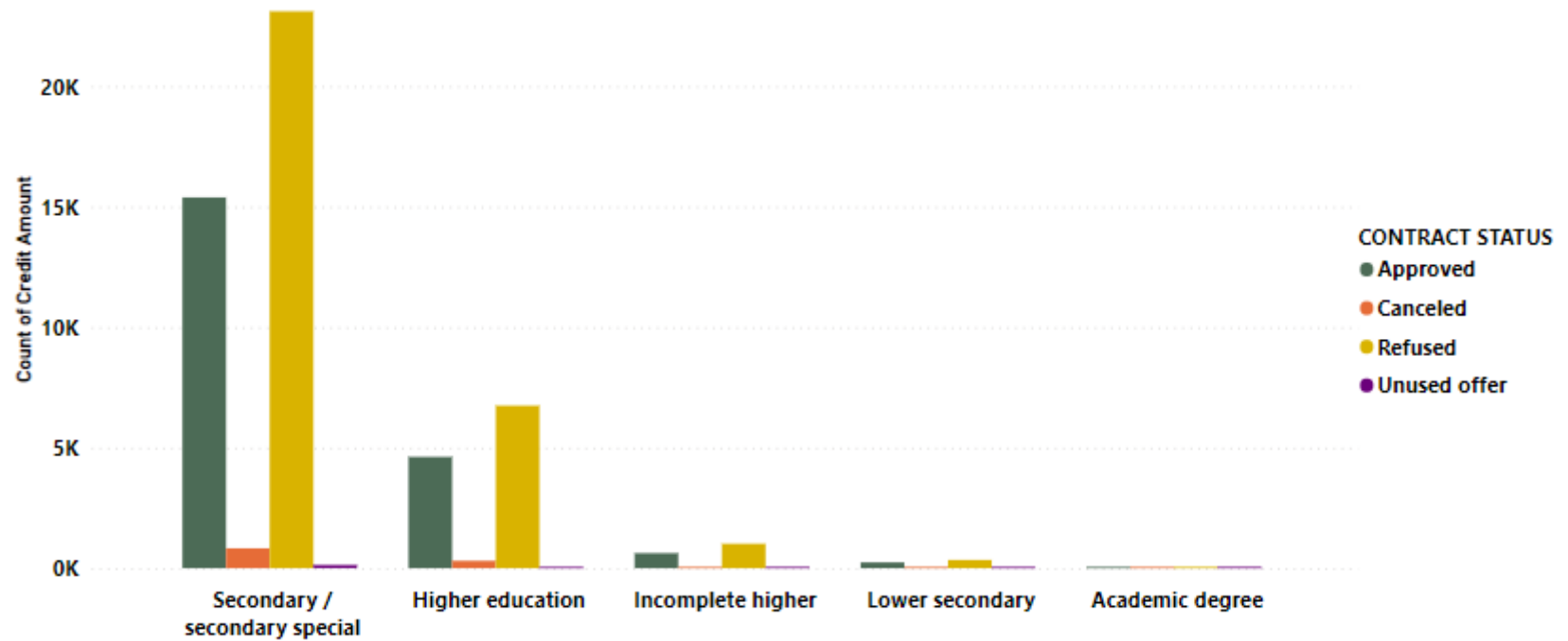
- The credit amount of Loan purposes like 'urgent needs', 'buying a used car' and 'building a house' are higher and also they are having no payment difficulties.
- Repairs loan purpose have maximum number of payment difficulties clients and maximum number of no payment difficulties clients.

**Distribution of Amount Credit Previously Vs Housing type w.r.t Target**



- Here for Housing type, house/apartment is having higher credit of target 1.
- So, we can conclude that bank should avoid giving loans to the housing type of office apartment , rented apartment and co-op apartment as they are having difficulties in payment.
- Bank can focus mostly on housing type with parents or House or apartment or municipal apartment for successful payments.

Distribution of Amount Credit Previously Vs Education Type w.r.t Contract Status



It can be seen that number of refusals are for clients has nothing to do with their education levels

Almost all education level clients have equal unused offers.



# Results

- Bank should focus on education type higher education and secondary/secondary special who took cash loan or revolving loan are capable of successful payment on time.
- All type of organization don't have any payment difficulties. Except business entity type 3 and self employed have some payment difficulties but it is negligible.
- Contract type revolving loans have less client with payment difficulties.
- Bank should focus on income type working, commercial associate and pensioner they are capable of successful payment back to bank.
- Housing type house/apartment, with parents and municipal apartment have more clients with no payment difficulties than payment difficulties.

**The end**