

# EDA ASSIGNMENT

Rushmi Roshina M

# Data Handling and Data cleaning

- Columns with more than 40% NULL values are dropped
- Few other Null values are substituted with most suitable values
- Inconsistent data(XNA/XAP) are substituted with appropriate values
- Data with errors(Days columns) are converted to their correct and convenient value.

# Forming Derived Metrics and Bins

- **Age Binning**

*[18,27,40,50,60,100]=[['Young','Young\_Adult', 'Middle\_Aged','Old','Very\_Old']]*

- **Income Binning**

*[min(income),50000,100000,300000,1000000,3000000,max(income)]=['Below\_Poverty','Poverty','Lower\_Middle\_Class','Upper\_Middle\_Class','Rich','Very\_Rich']*

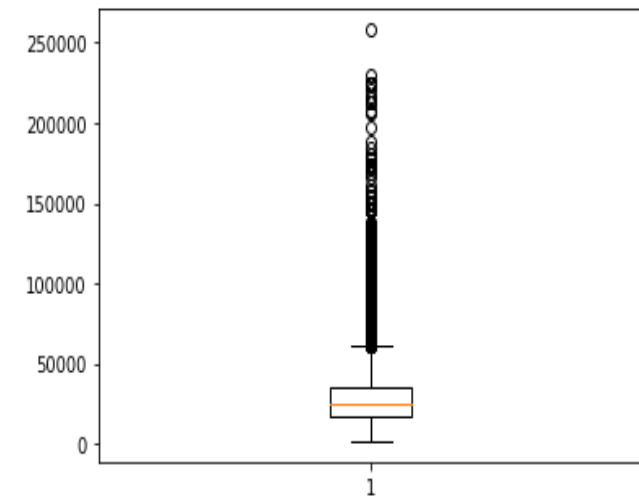
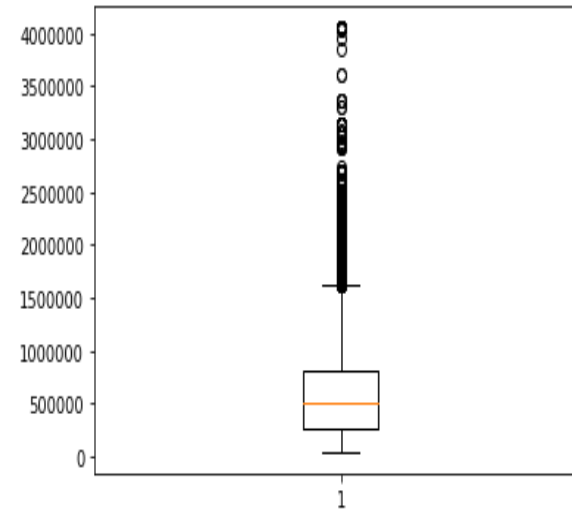
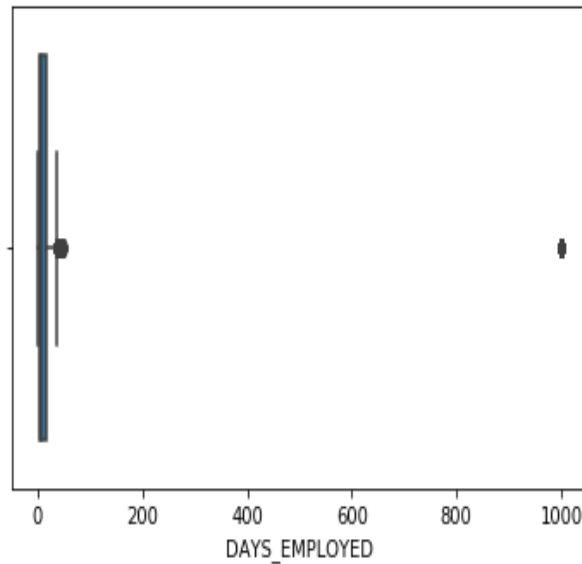
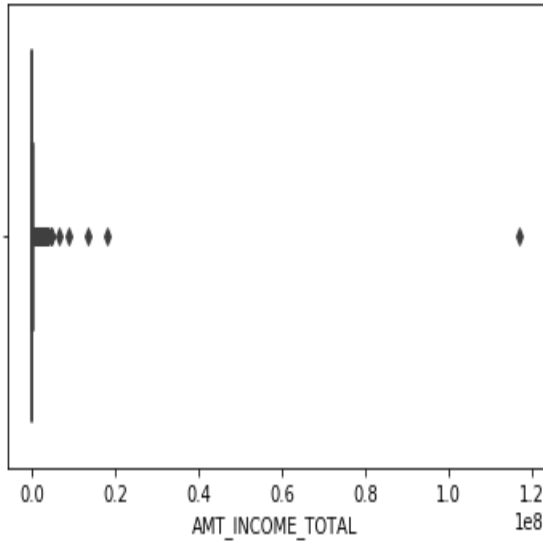
- **CRED\_INC\_RATIO** = *Ratio of Credit\_Amount and Income*

- **CRED\_GOODS\_RATIO** = *Ratio of Credit\_Amount and Goods*

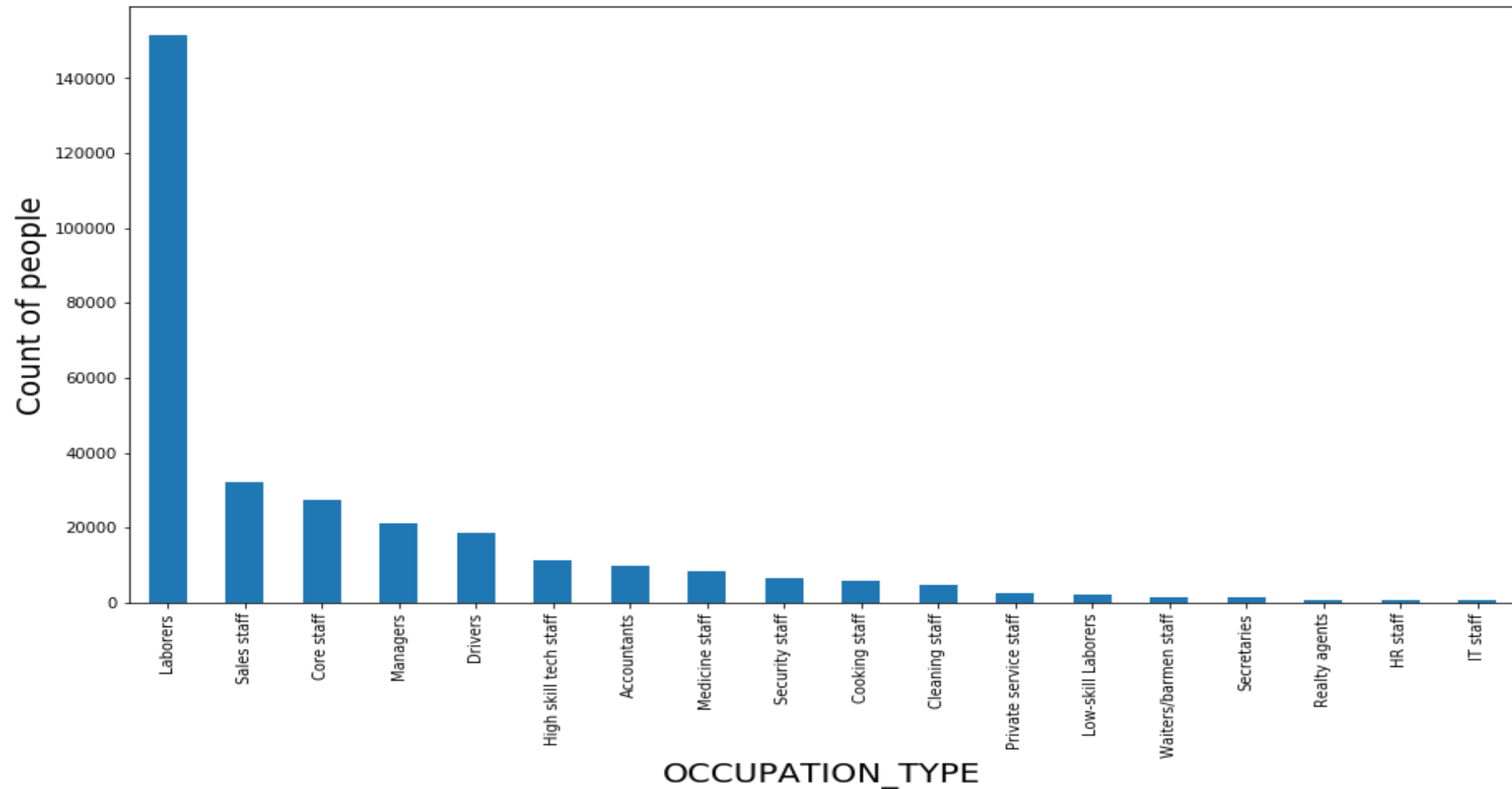
The above Metrics/Bins are created for Application Dataset

# Outlier Analysis

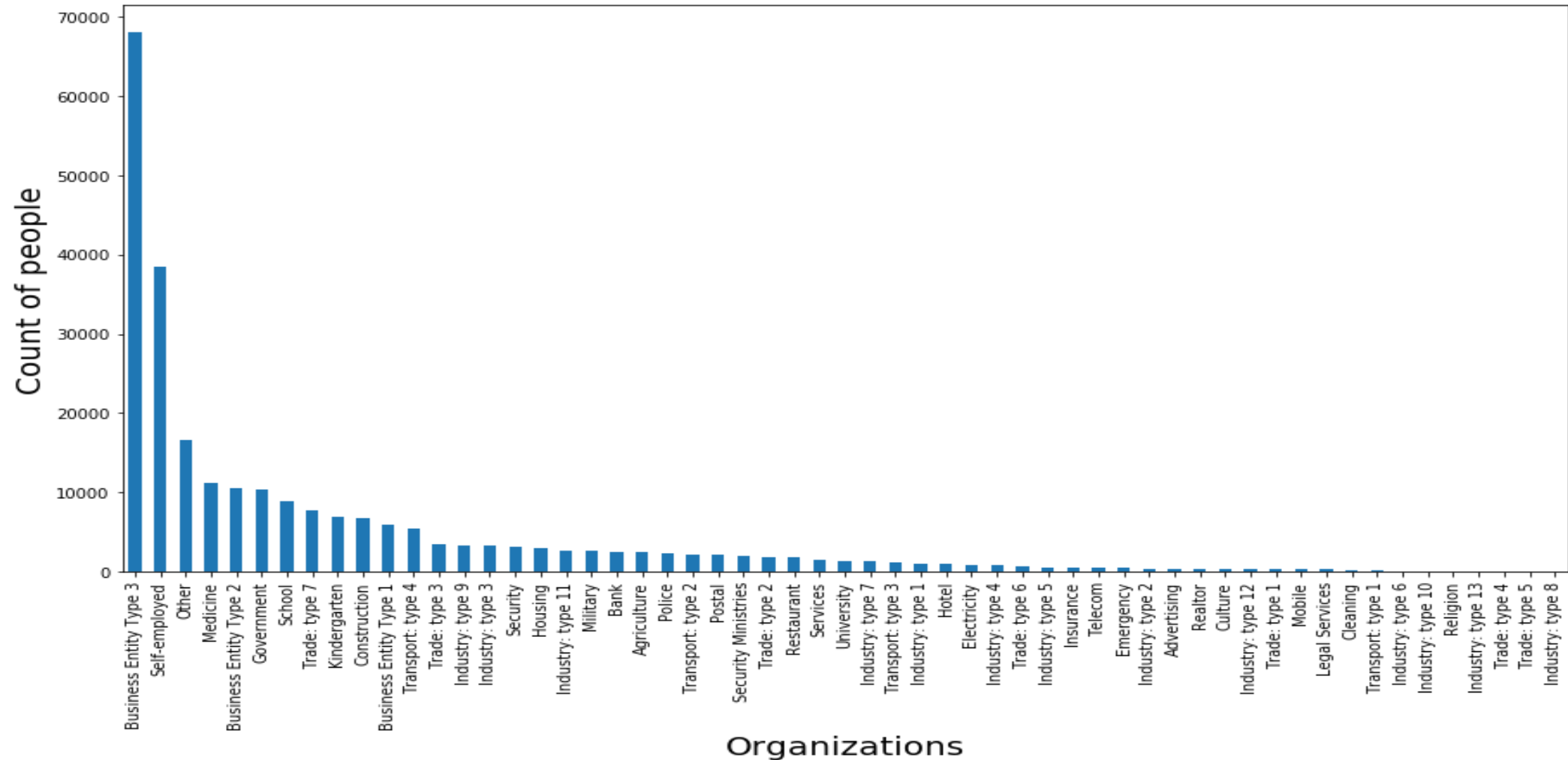
Outliers are Analyzed for AMT\_INCOME\_TOTAL,DAYS\_EMPLOYED,AMT\_CREDIT,AMT\_ANNUITY.They seem to have no impact on the data.Hence it is not taken care off.



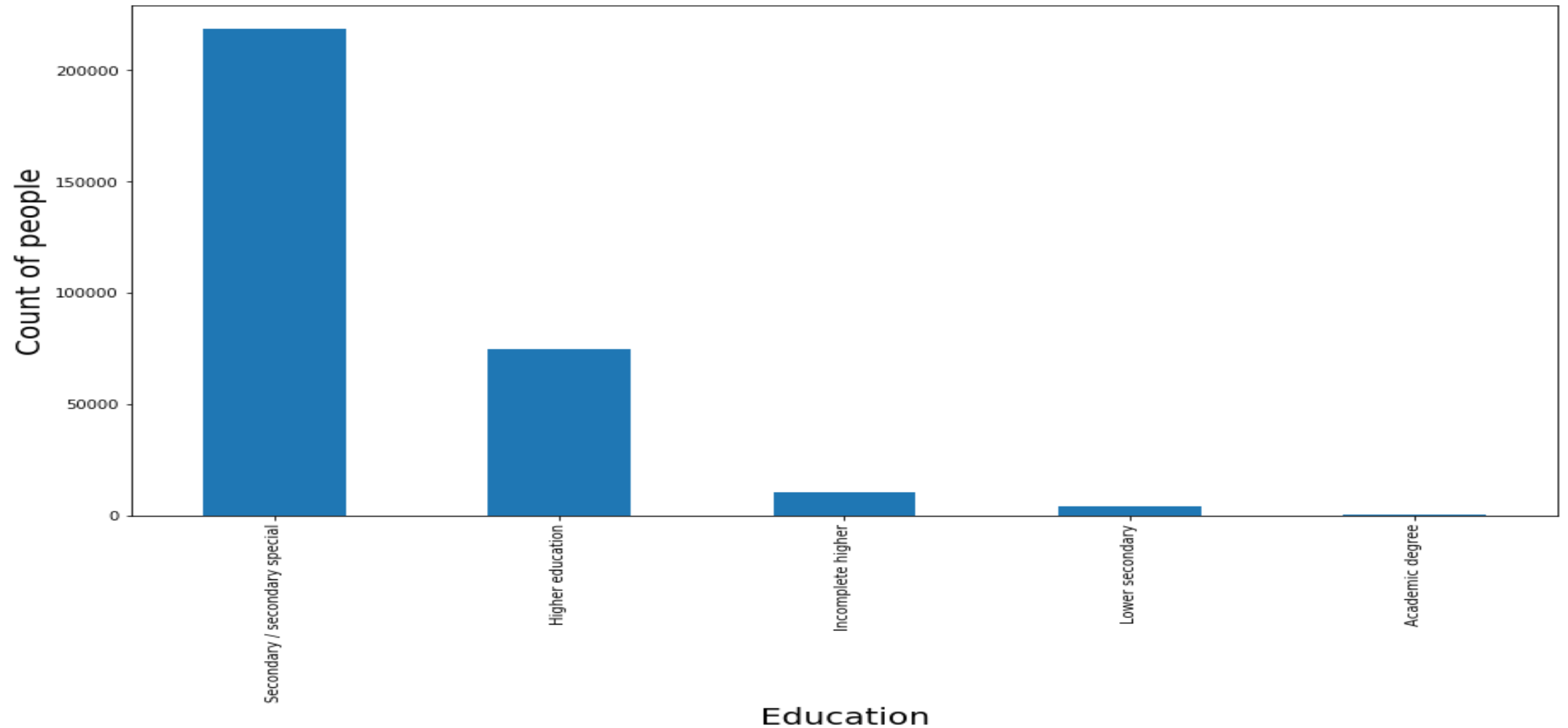
# Data Distribution on Occupation type



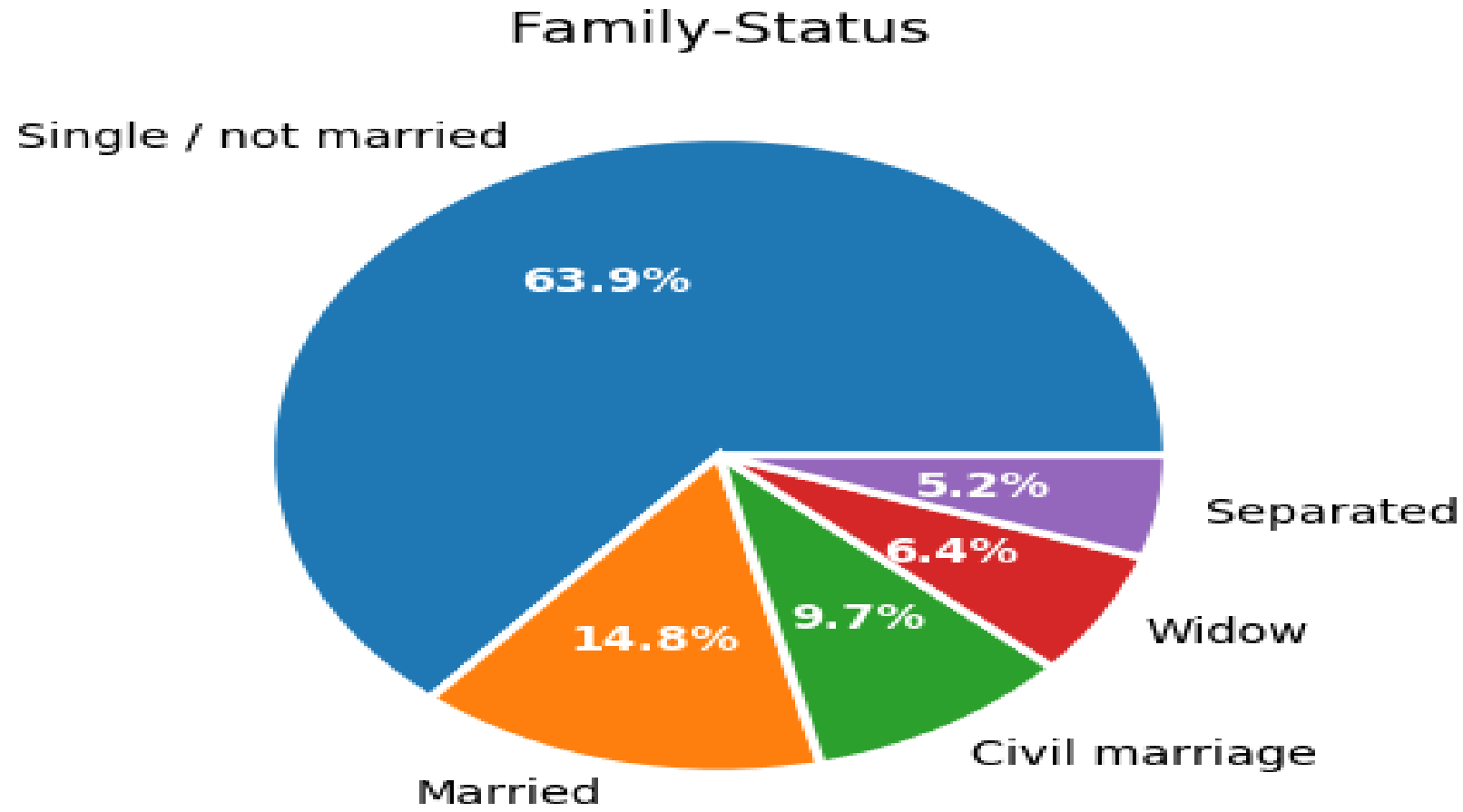
# Data Distribution on Organization type



# Data Distribution on Education type



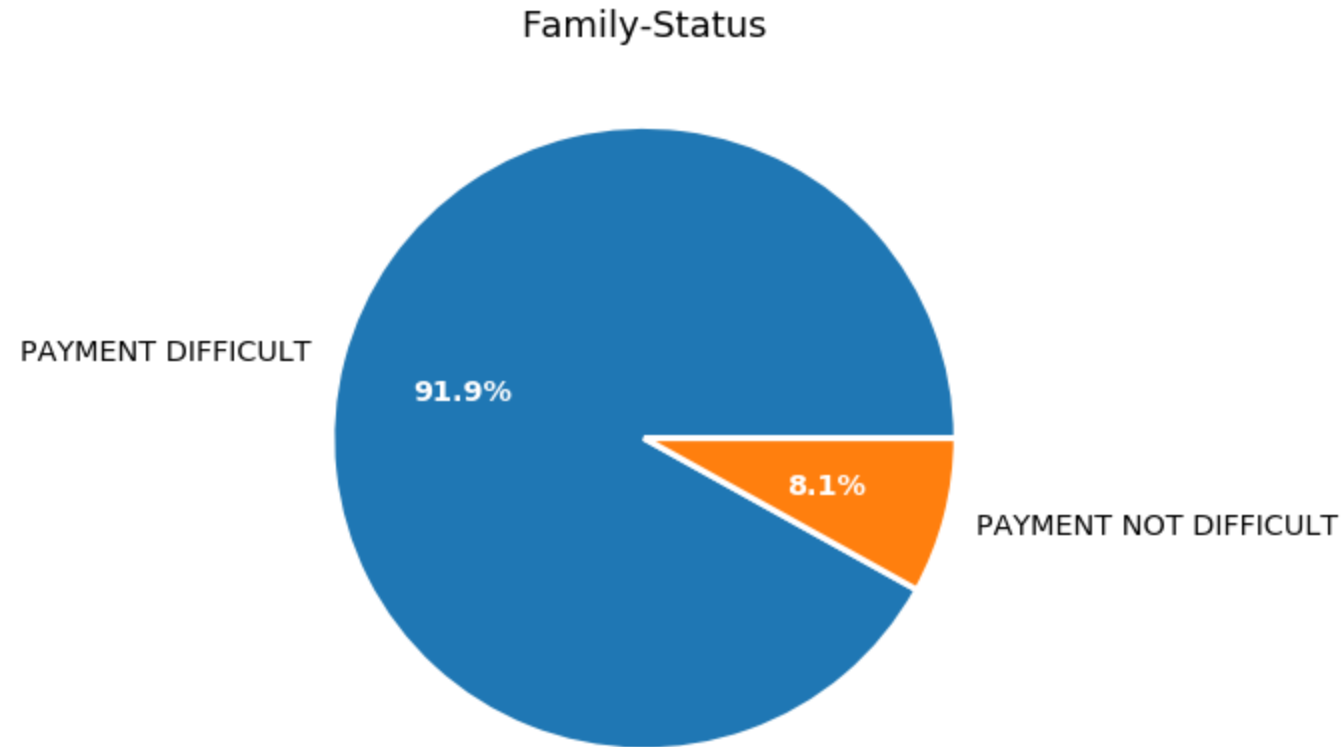
# Proportion of Family-status





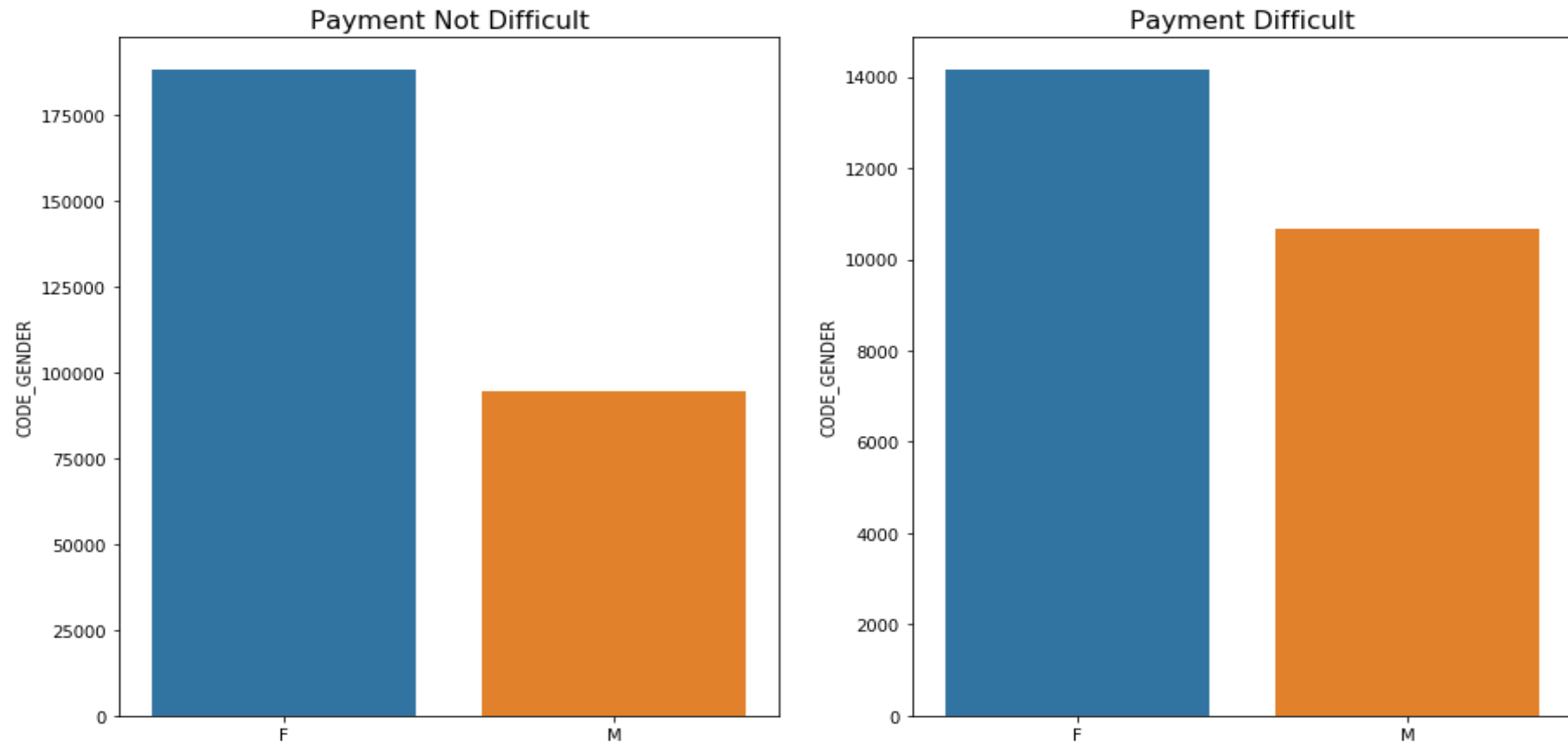
# Data Imbalance

Most of the Applications have difficulties with payment

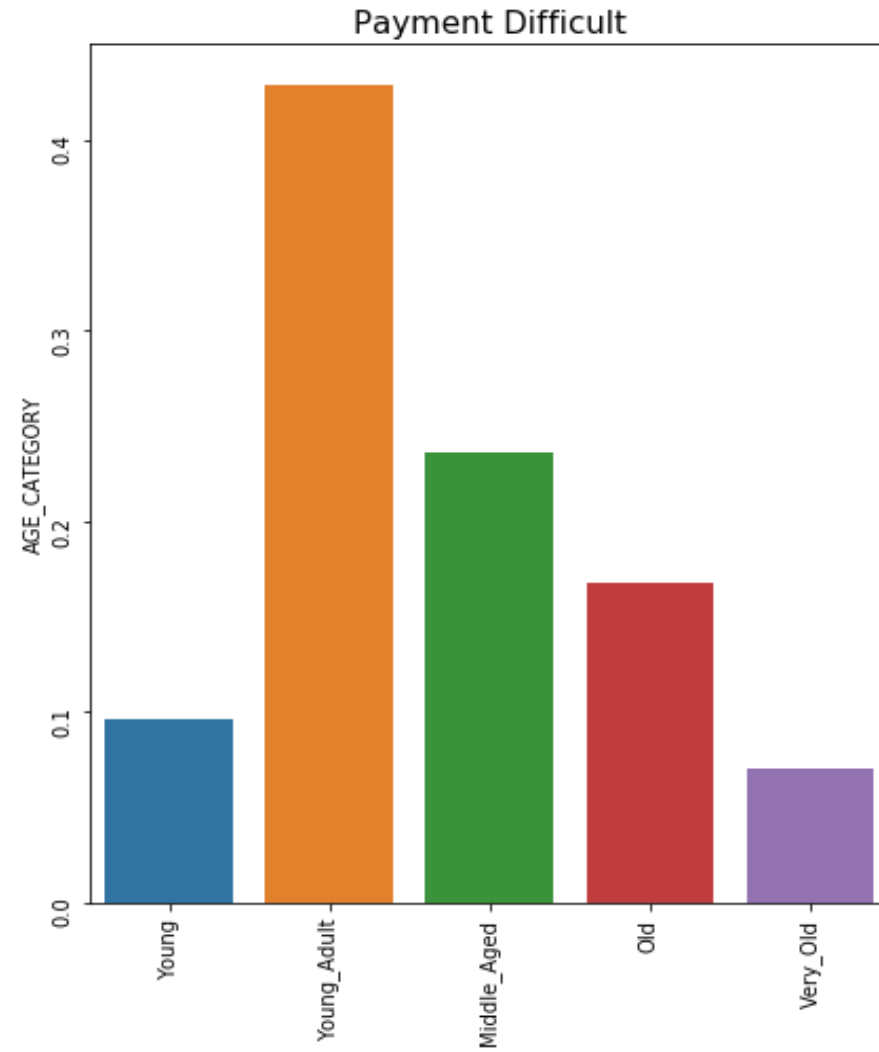
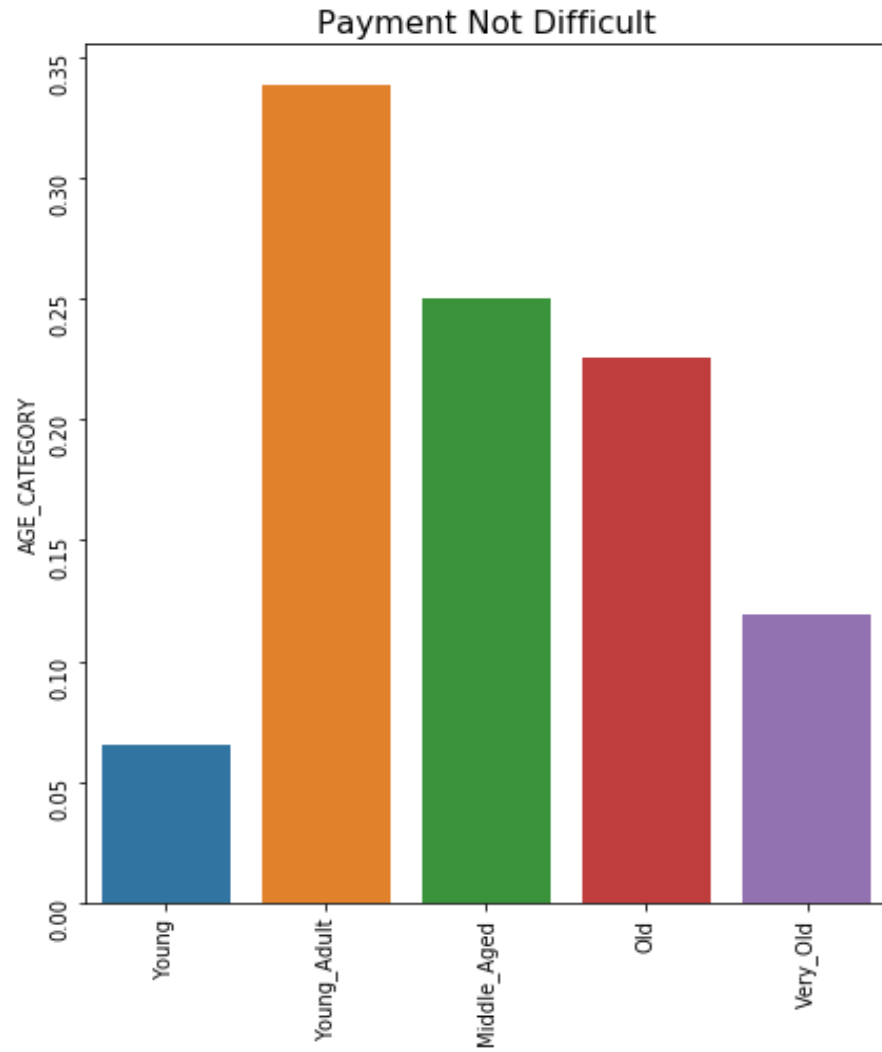


# Bi-variate Analysis

## Gender Analysis

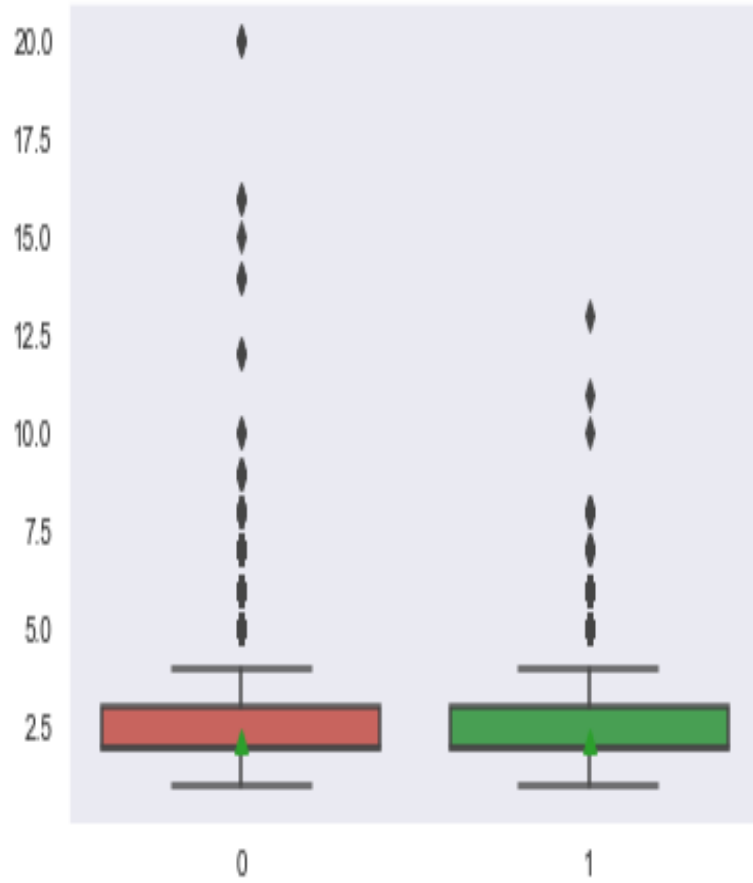


## AGE CATEGORY ANALYSIS

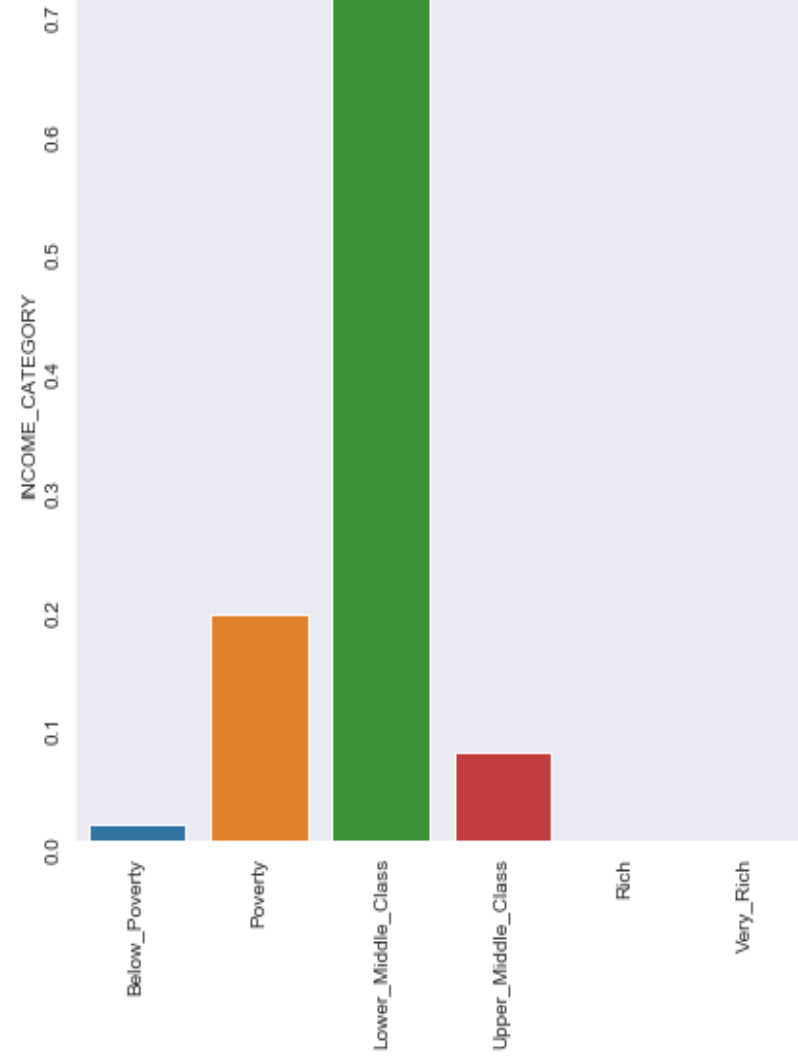


## INCOME\_RANGE distribution

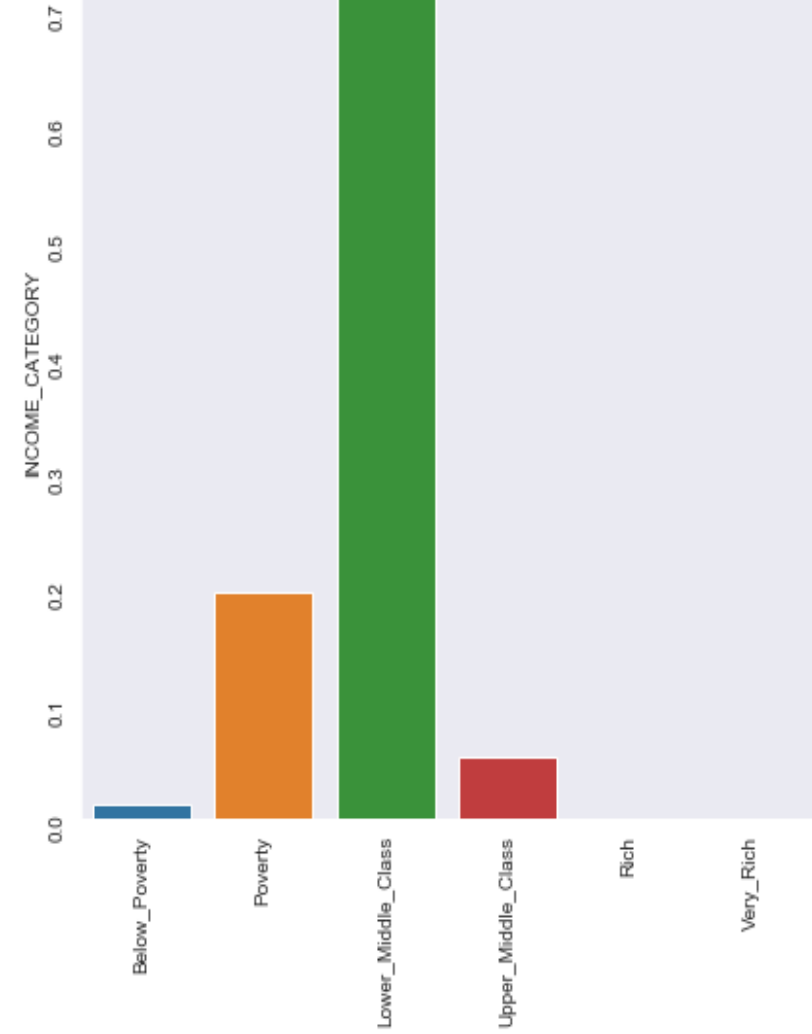
### Family Member Count Analysis



### Payment Not Difficult

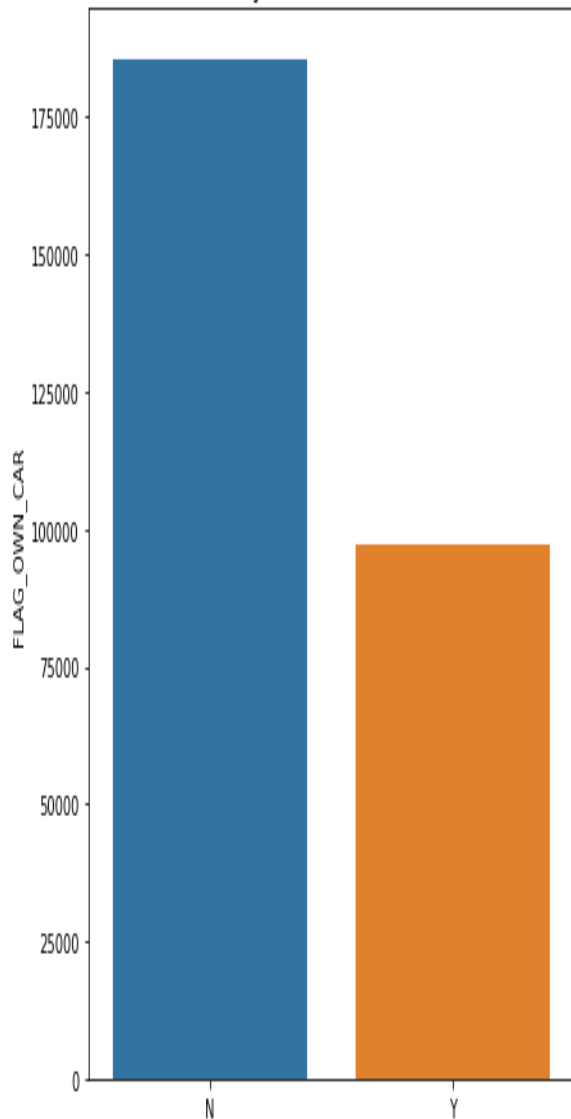


### Payment Difficult

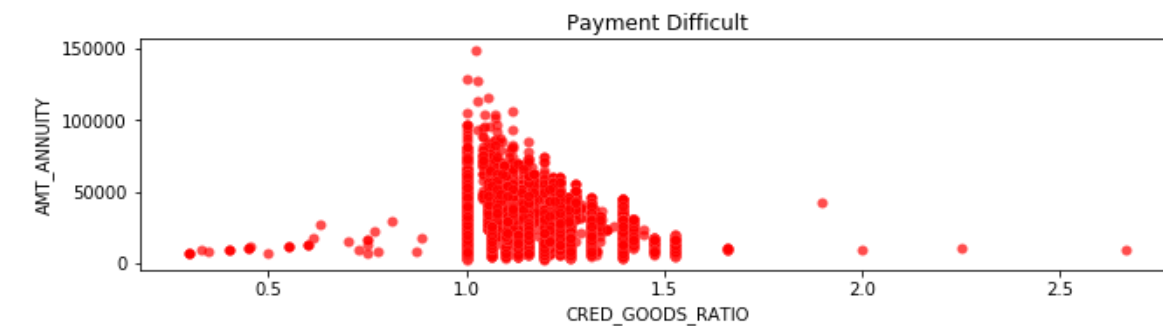
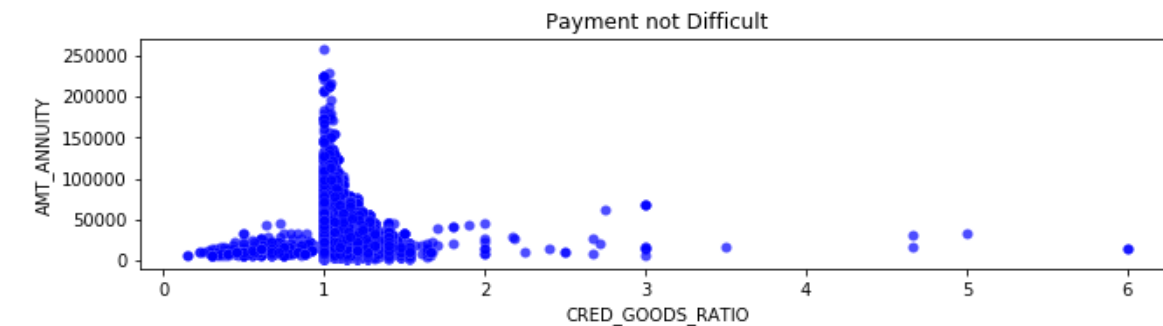
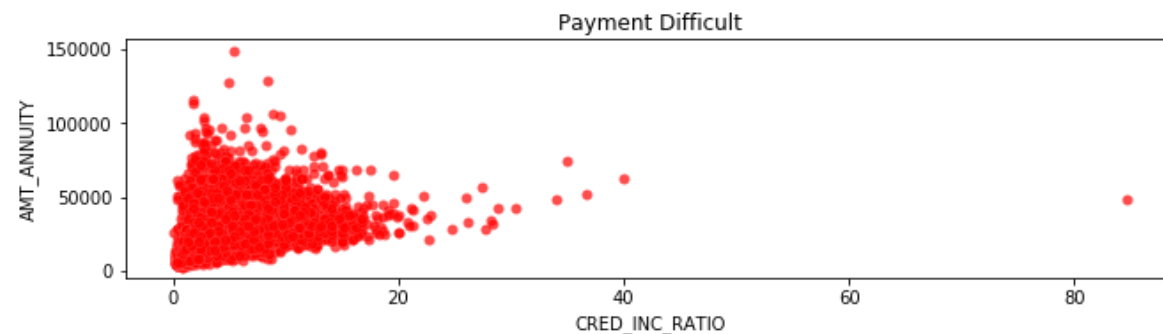
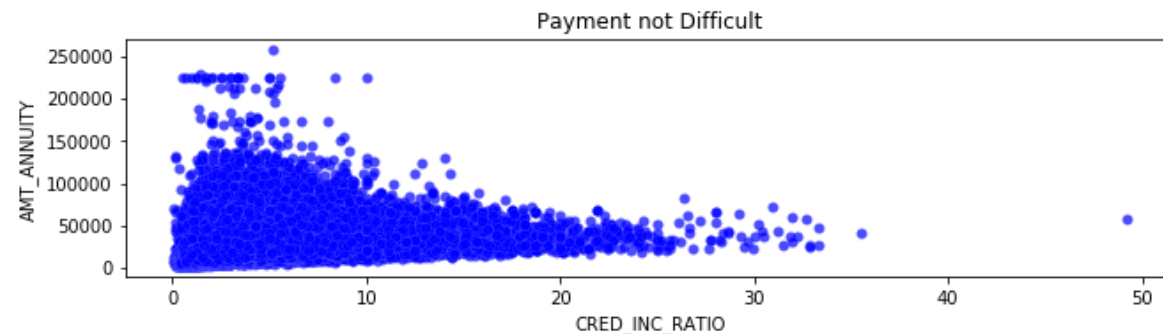
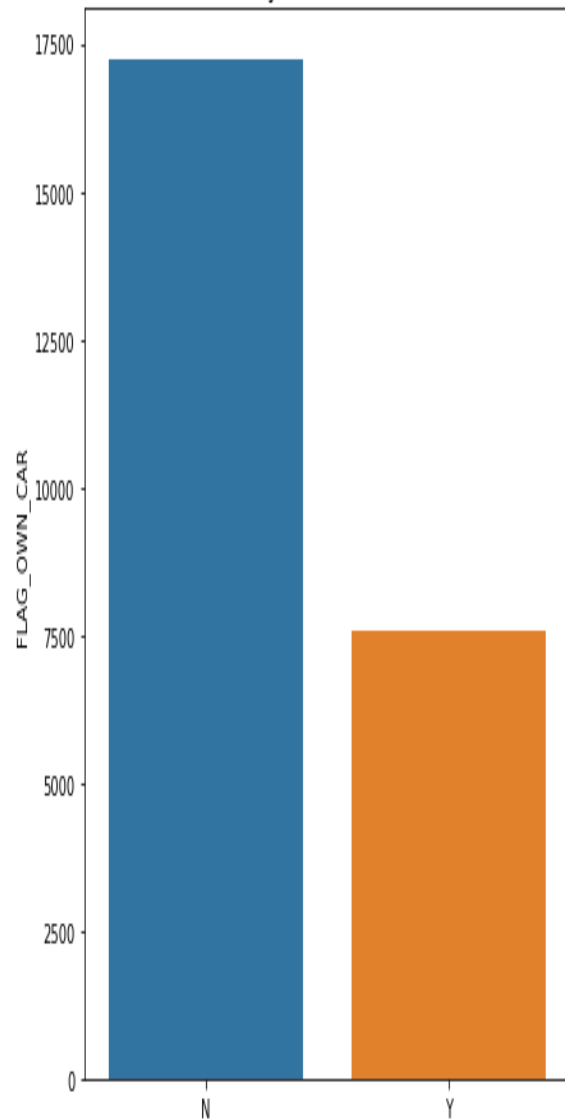


# Owens Car Analysis

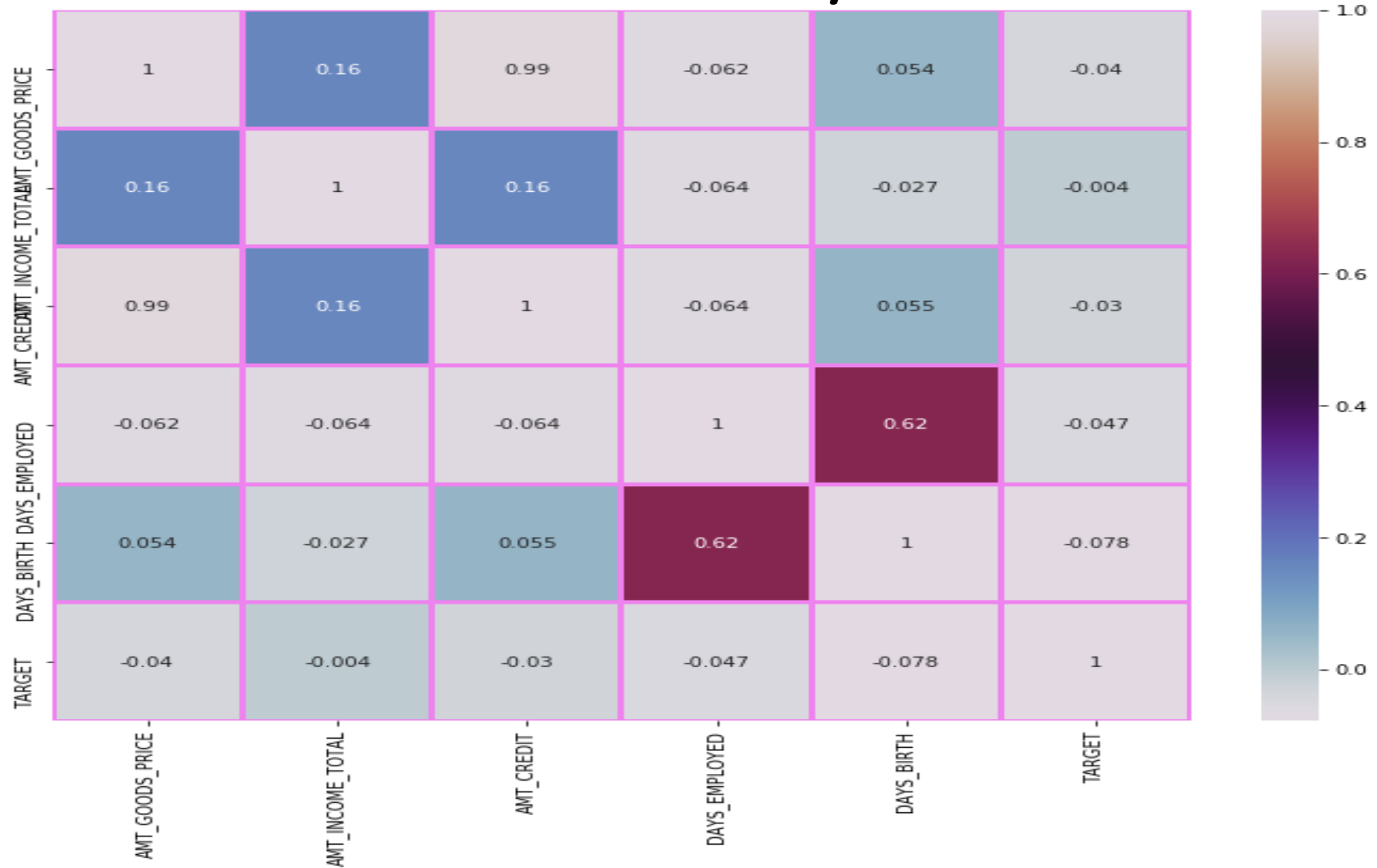
Payment Not Difficult



Payment Difficult



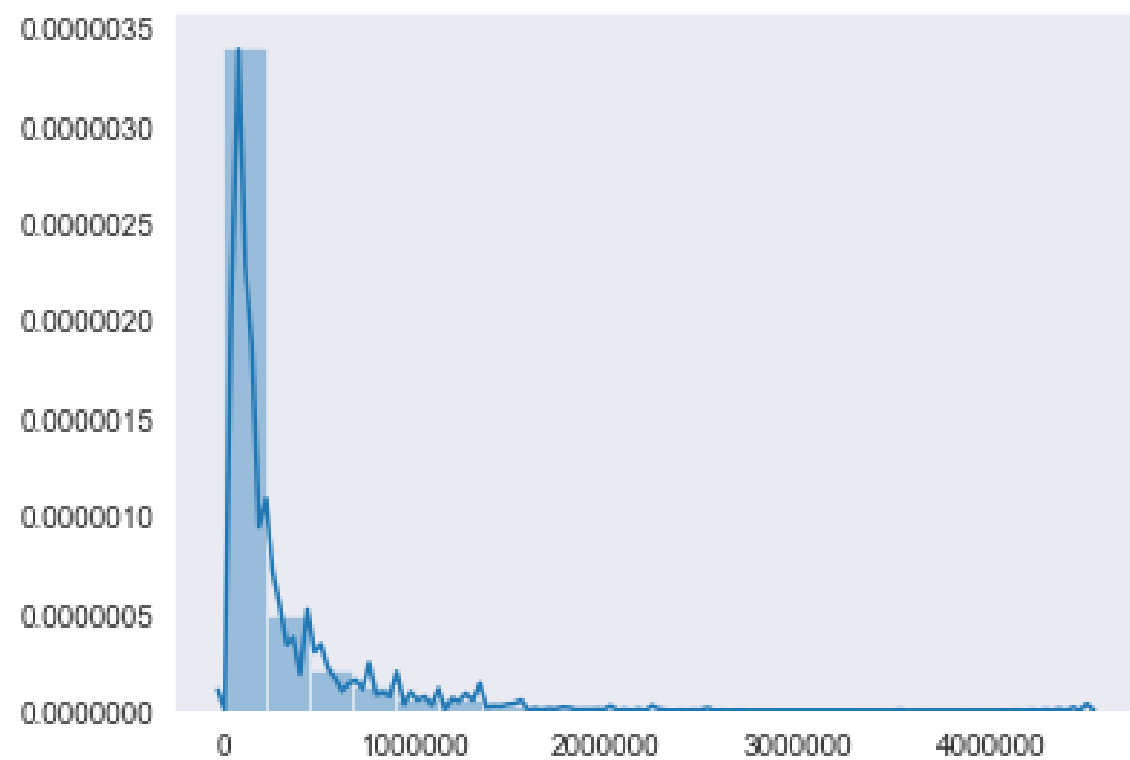
# Multivariate Analysis



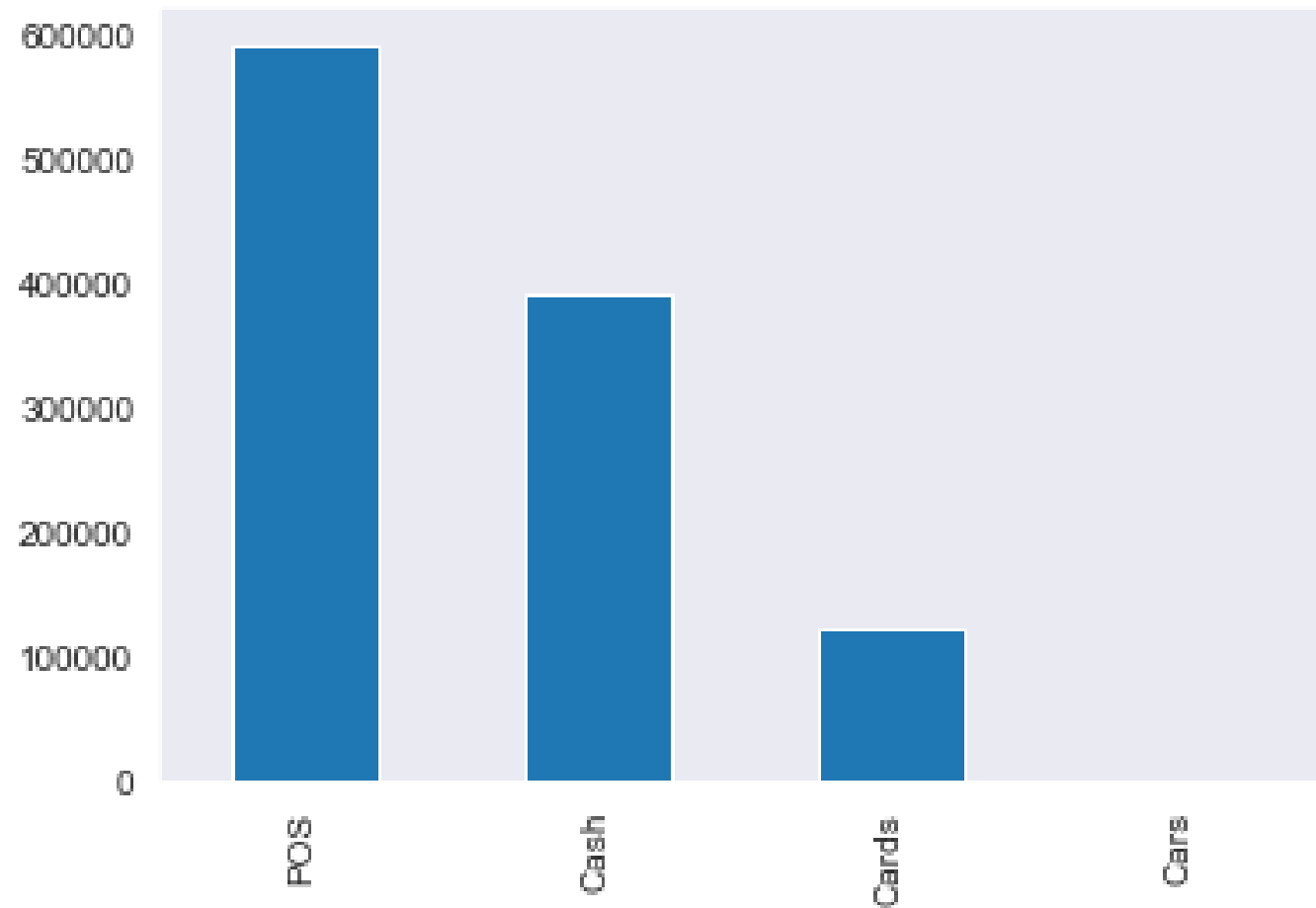
# Insights on Application Data

- Most of the Male Applicants have Payment Difficulty
- People of age less than 40 have more payment difficulty
- Family Member count has no big impact on the Target variable
- All the Income\_Category people are equally distributed on the Target variable
- People who don't own car have more difficulty with payment than the people who own car

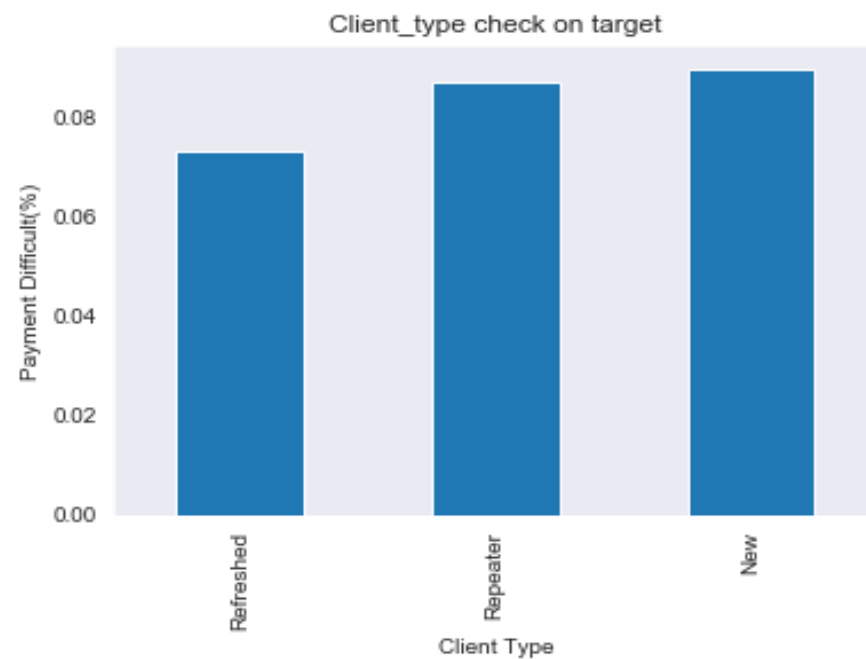
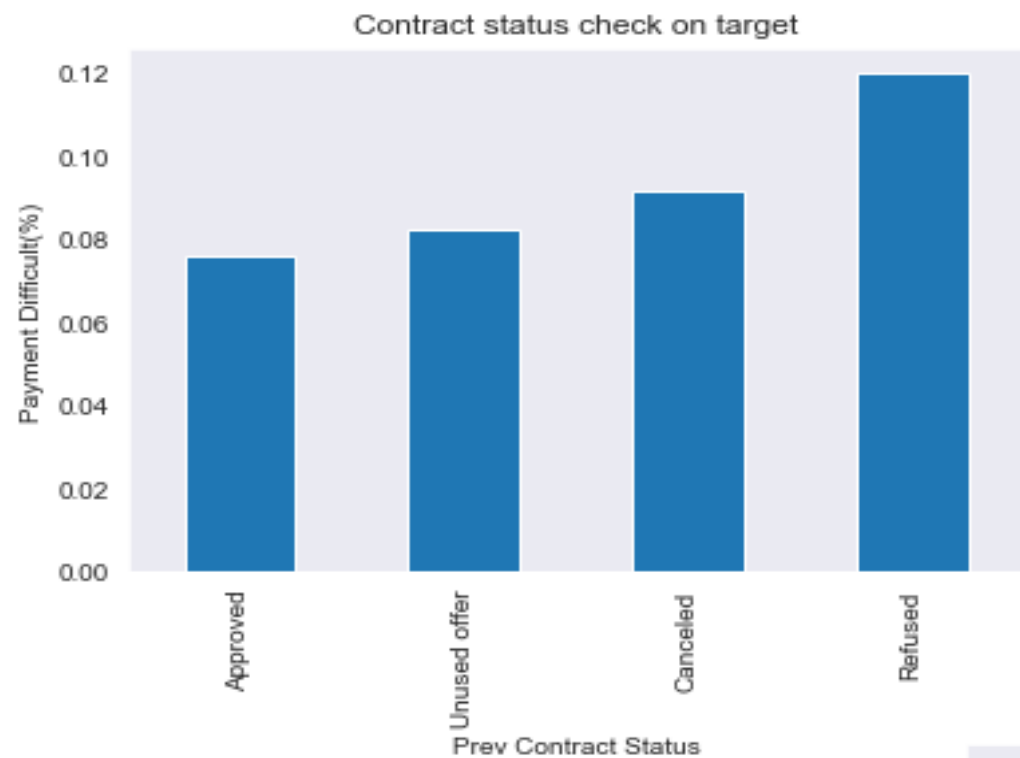
Previous Application loan Amount



NAME\_PORTFOLIO Analysis





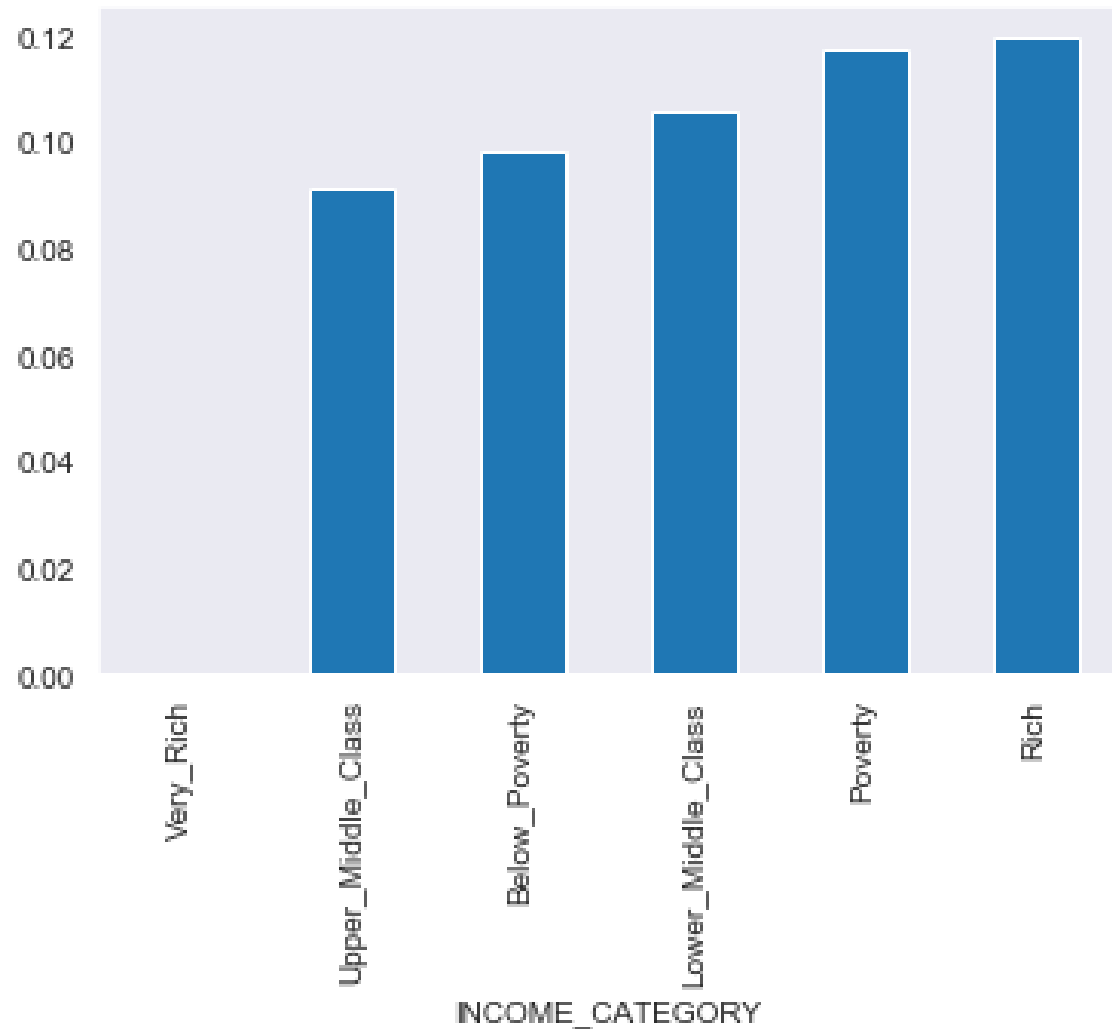


# Insights on Previous Application Dataset

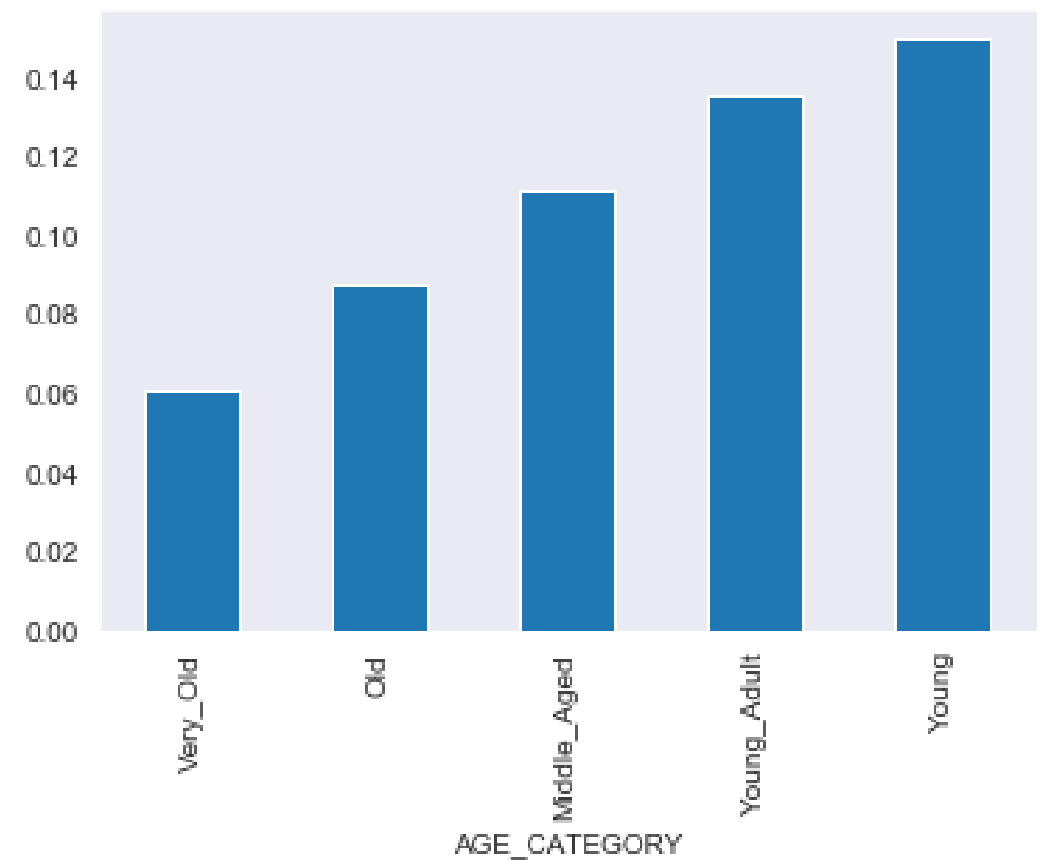
- Most of the Applicant who has difficulties with payment, was previously Refuse loan
- The highest rejected reason for people with payment difficulties was SCOFR
- Most Applicant's loan\_amount was below 1000000
- Most of the loan was through POS
- People with payment Difficulties are mostly new loan Applicants

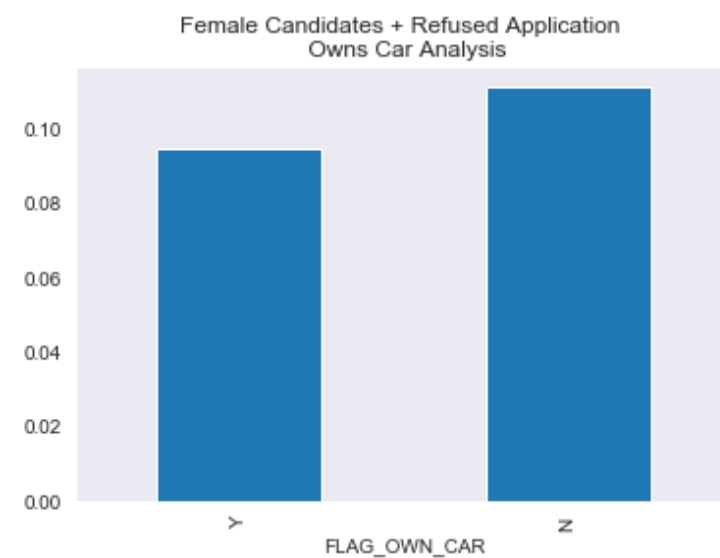
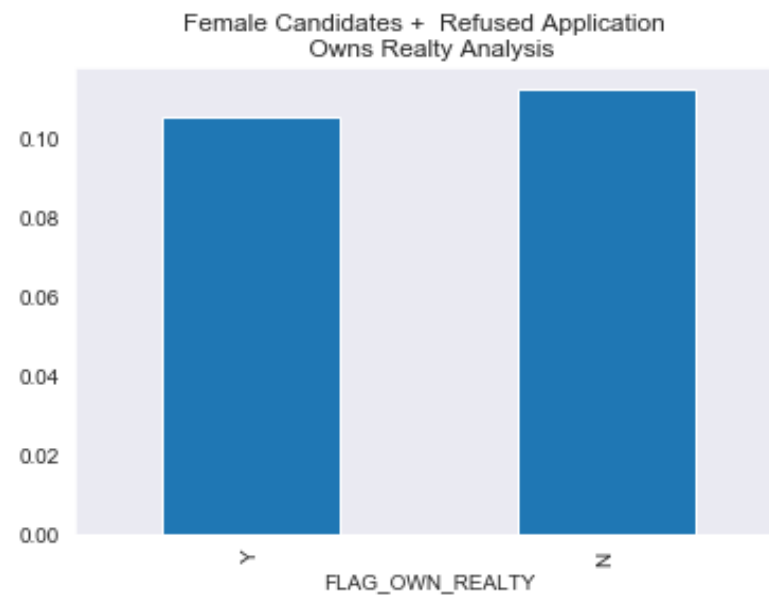
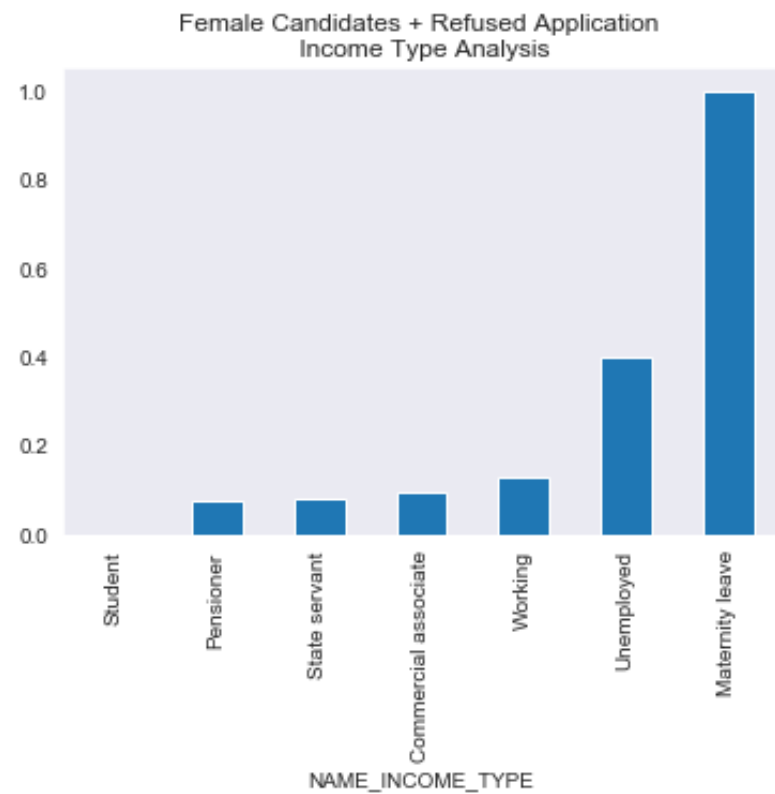
# Combined Analysis(Female)

Female Candidates + Refused Application  
Income Category Analysis



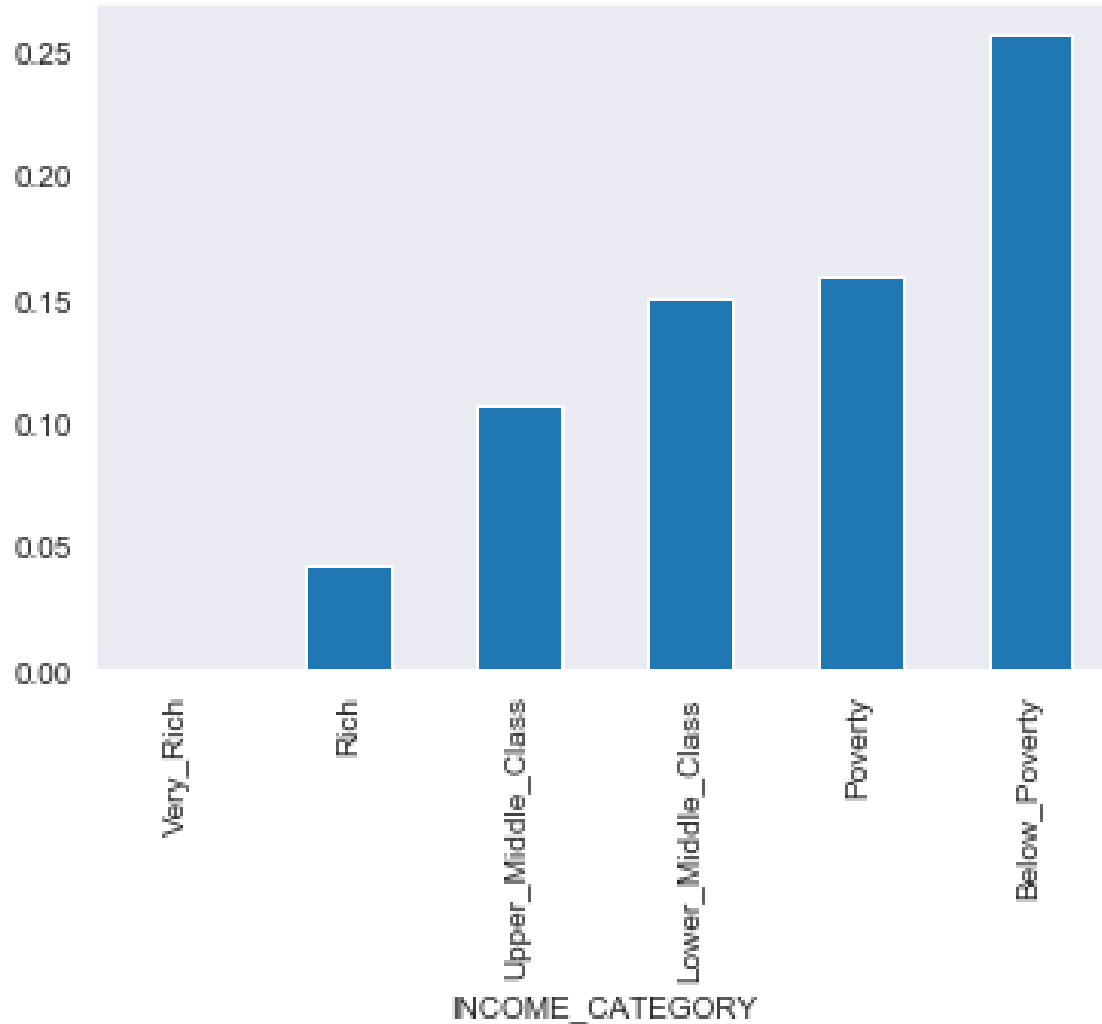
Female Candidates + Refused Application  
Age Category Analysis



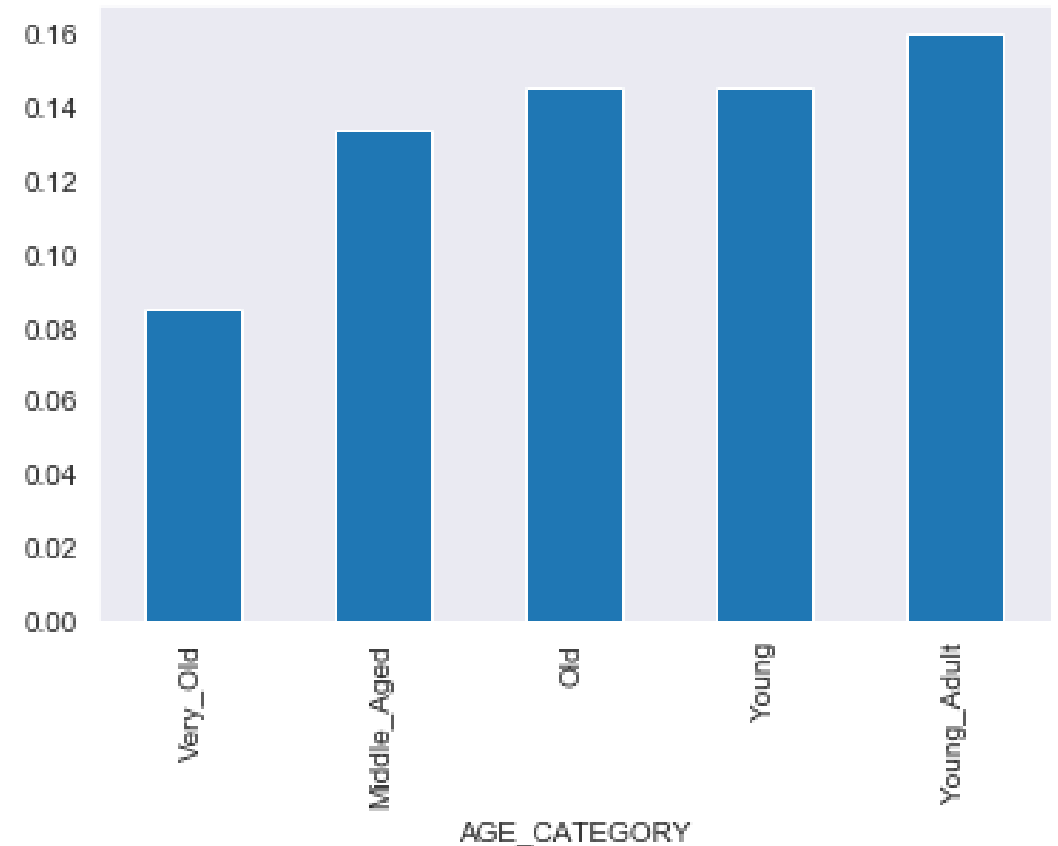


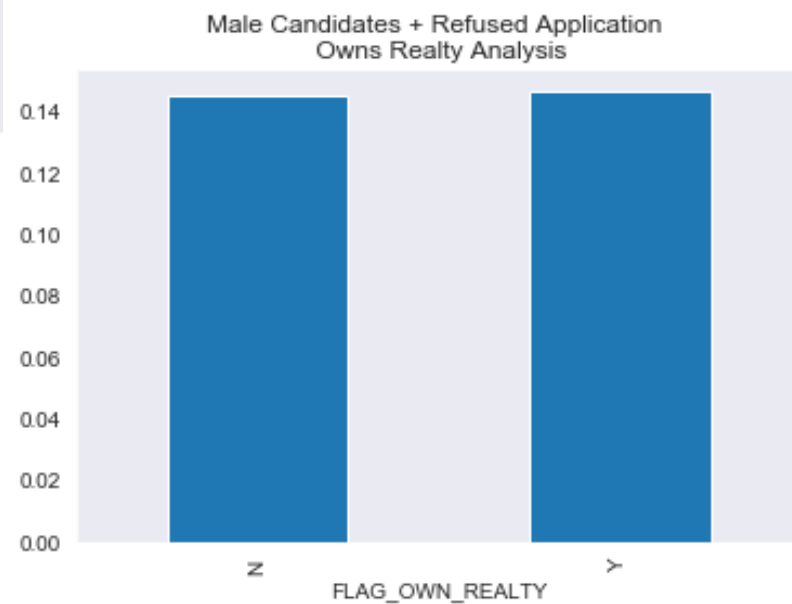
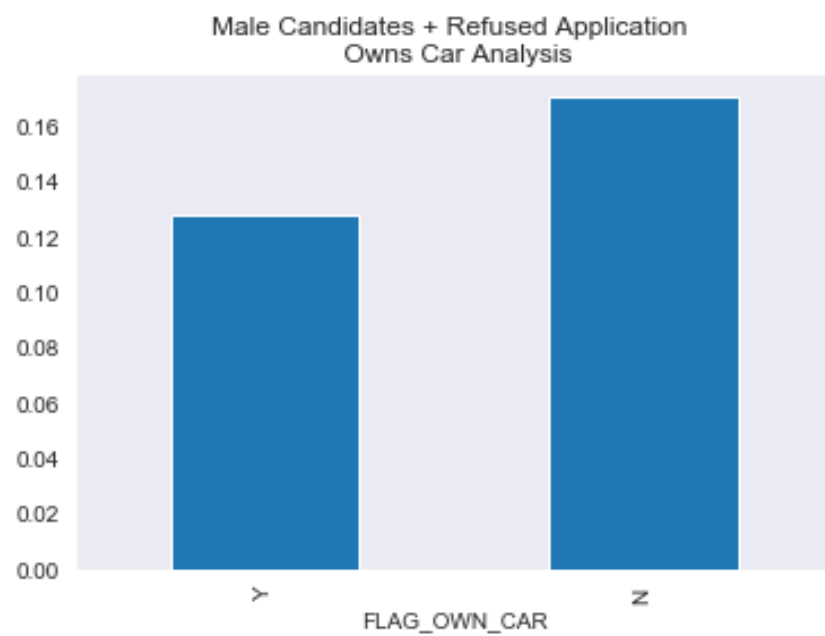
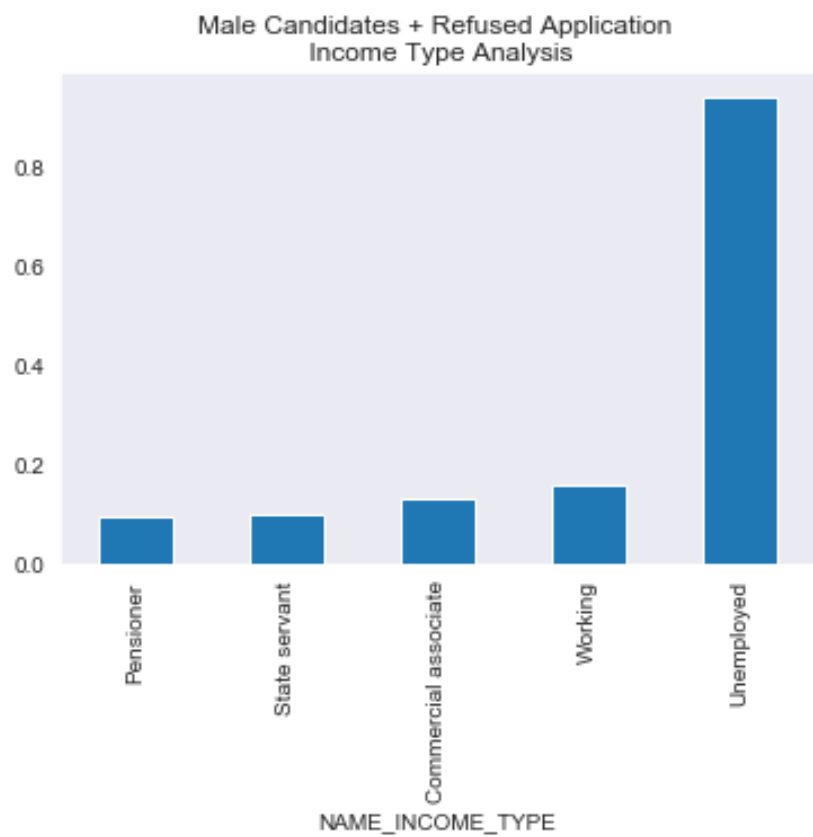
# Combined Analysis(Male)

Male Candidates + Refused Application  
Income Category Analysis

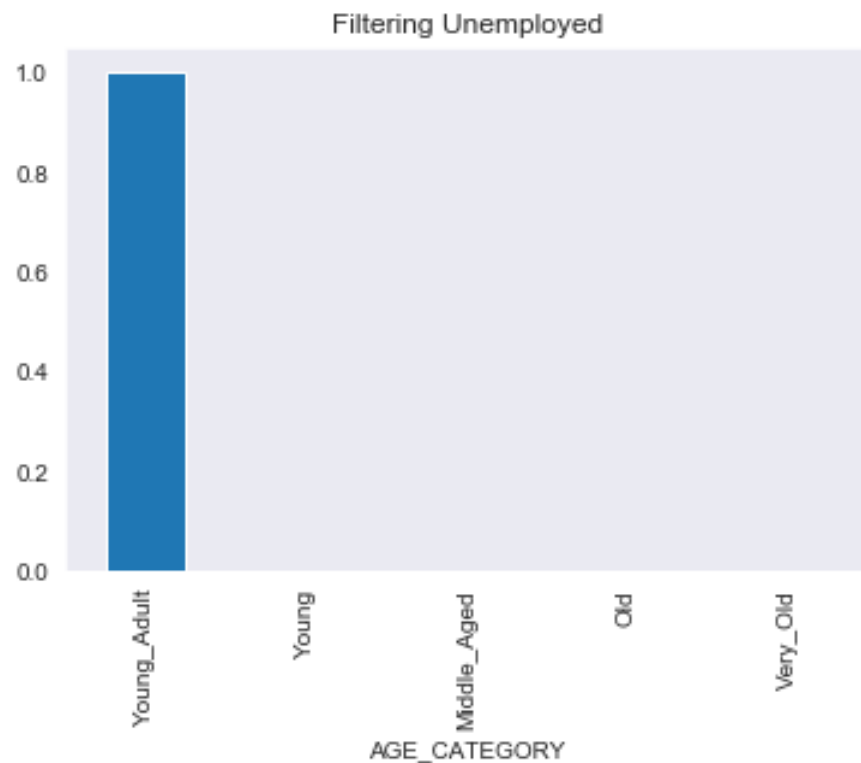


Male Candidates + Refused Application  
AGE\_CATEGORY Analysis

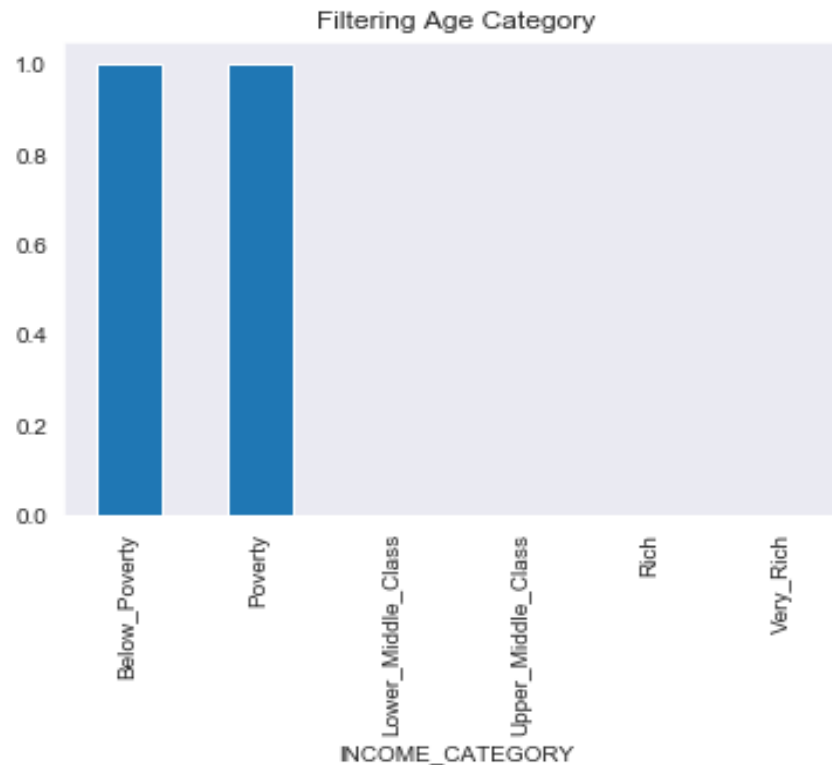




Further these Analysis can be used to filter out rows and fetch the people who have least possibility of paying back the loan



Click  
to  
add  
text



# Combined Insights

- Most people who doesn't own Car or Realty may not be able to payback the loan
- Female Candidates who fall under the category Young and are in Maternity leave may face difficulties in paying back the loan
- Male Candidates who fall under the category Young-Adult and are Unemployed and are in the Below\_Poverty range may face difficulties in paying back the loan