

# EDA REPORT FOR CREDIT CARD DEFAULTERS ANALYSIS

**Muhammad Hanafi Bin Mohd Sani**

---

# BUSINESS UNDERSTANDING & OVERVIEW

---

**Understanding the cause of behind the loan defaults**

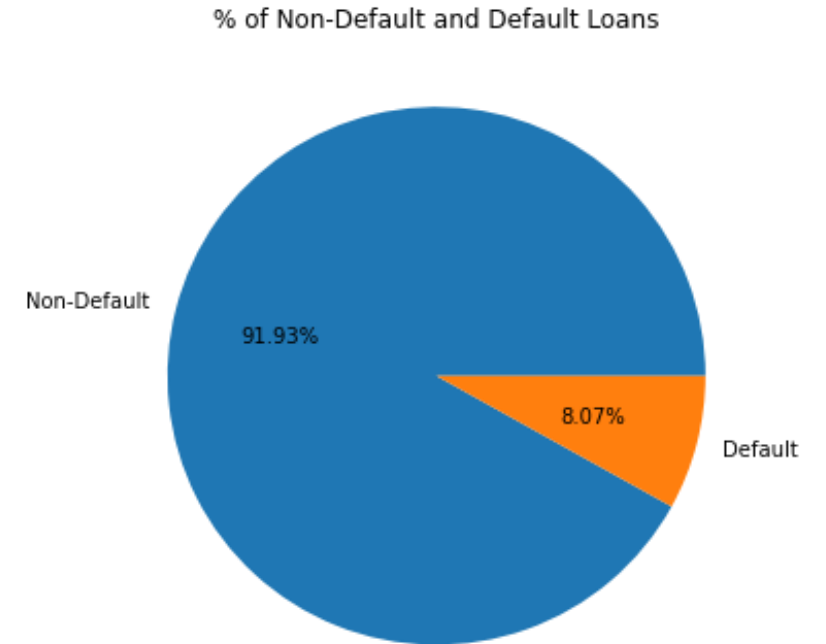
**By performing the EDA, we can pinpoint the cause based on the insight generated from the analysis**

# UNDERSTANDING THE DATA

---

## TARGET VARIABLE:FINDINGS

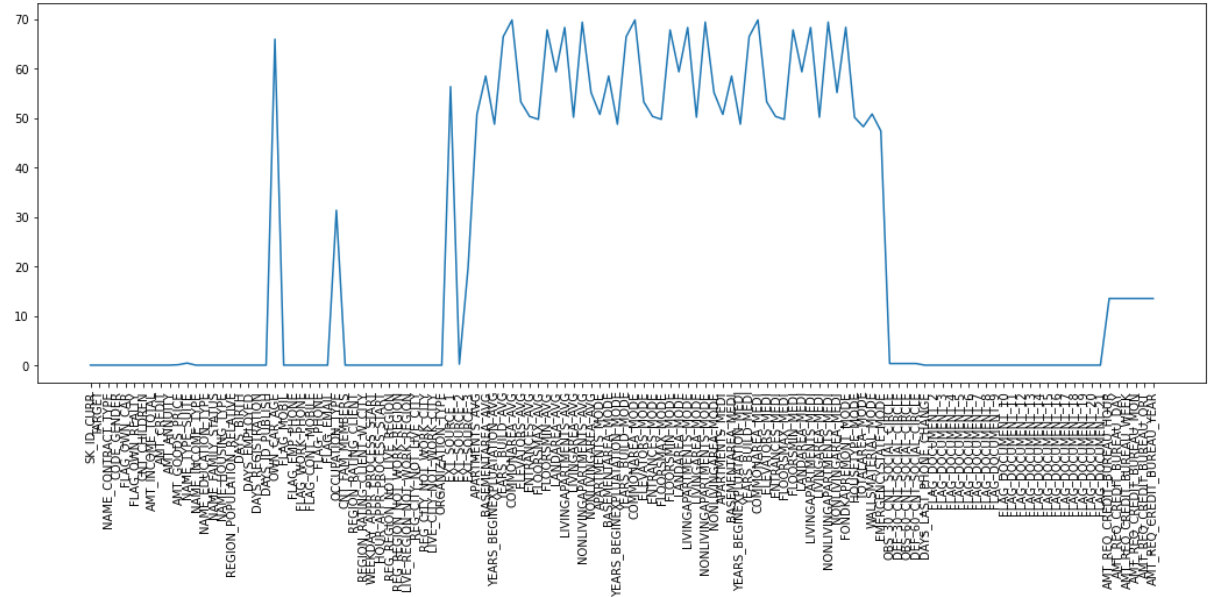
- Target variables are highly unbalanced:  
Only 8.07% for Default and a whopping 91.93% for non-Defaults.
- Meaning, most of the loans are paid on time (non-default)
- Further analysis required



# Missing Data

- 
- The line plot displays the importance of various features for predicting loan status. The y-axis represents importance from 0 to 70. The x-axis lists features. The plot shows that 'SK\_ID\_CURR' is the most important feature, followed by 'AMT\_REQ\_CREDIT\_BUREAU\_YEAR' and 'AMT\_REQ\_CREDIT\_BUREAU\_MON'. Other features like 'YEARS\_BEGINNING\_WHERE\_YOU\_CURRENTLY\_LIVE' and 'YEARS\_BEGINNING\_WHERE\_YOU\_CURRENTLY\_WORK' also show high importance.

# Initial Analysis From The Data (Cont.)



- Resulting columns left after removal:  
'SK\_ID\_CURR', 'TARGET', 'NAME\_CONTRACT\_TYPE',  
'CODE\_GENDER', 'FLAG\_OWN\_REALTY', 'FLAG\_OWN\_CAR', 'CNT\_CHILDREN', 'AMT\_INCOME\_TOTAL',  
'AMT\_CREDIT', 'AMT\_ANNUITY', 'AMT\_GOODS\_PRICE', 'NAME\_INCOME\_TYPE',  
'NAME\_EDUCATION\_TYPE', 'NAME\_FAMILY\_STATUS', 'NAME\_HOUSING\_TYPE', 'DAYS\_BIRTH',  
'DAYS\_EMPLOYED', 'DAYS\_REGISTRATION', 'DAYS\_ID\_PUBLISH', 'FLAG\_MOBIL', 'FLAG\_EMP\_PHONE',  
'FLAG\_CONT\_MOBILE', 'FLAG\_PHONE', 'FLAG\_EMAIL', 'OCCUPATION\_TYPE', 'ORGANIZATION\_TYPE',  
'CNT\_FAM\_MEMBERS', 'REGION\_RATING\_CLIENT\_W\_CITY', 'REGION\_RATING\_CLIENT'  
• FLAG\_WORK\_PHONE is a duplicate of FLAG\_PHONE

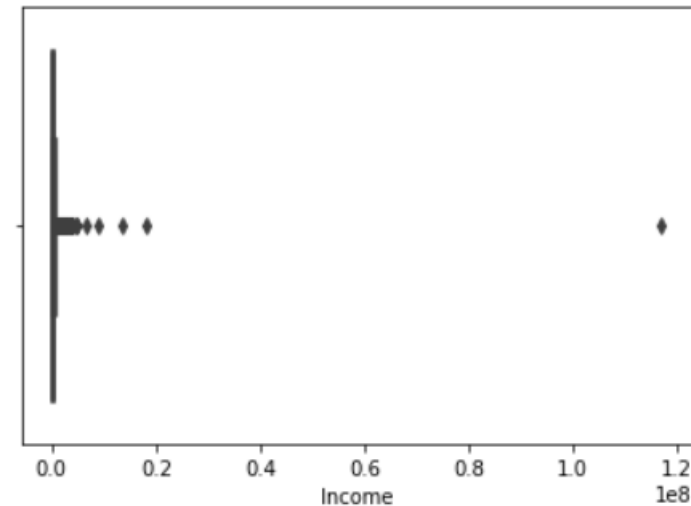
# Incorrect Data

- Some data like ORGANIZATION\_TYPE has “XNA” in it instead of NULL VALUES. To fix this: we will remove the row or convert the XNA into mean, median or mode in the columns depending on the situation.
- Data that should such as DAYS\_EMPLOYED and DAYS\_REGISTRATION are in negative values and should be positive.
- Some of data such as AMT\_INCOME\_TOTAL and have such an extreme value that caused skewedness in analysis. To fix: remove them. Removal of the extreme data does not affected the analysis.
- Data such as ['AMT\_ANNUITY', 'AMT\_APPLICATION', 'AMT\_CREDIT', 'AMT\_GOODS\_PRICE'] are converted into thousands.

```
Business Entity Type 3    67992
XNA                        55374
Self-employed             38412
Other                     16683
Medicine                  11193
Business Entity Type 2    10553
Government                10404
School                    8893
Trade: type 7              7831
Kindergarten              6880
Name: ORGANIZATION_TYPE, dtype: int64
```

```
365243    55374
-200      156
-224      152
-230      151
-199      151
Name: DAYS_EMPLOYED, dtype: int64

-1.0    113
-7.0     98
-6.0     96
-4.0     92
-2.0     92
Name: DAYS_REGISTRATION, dtype: int64
```



# CATEGORICAL ANALYSIS

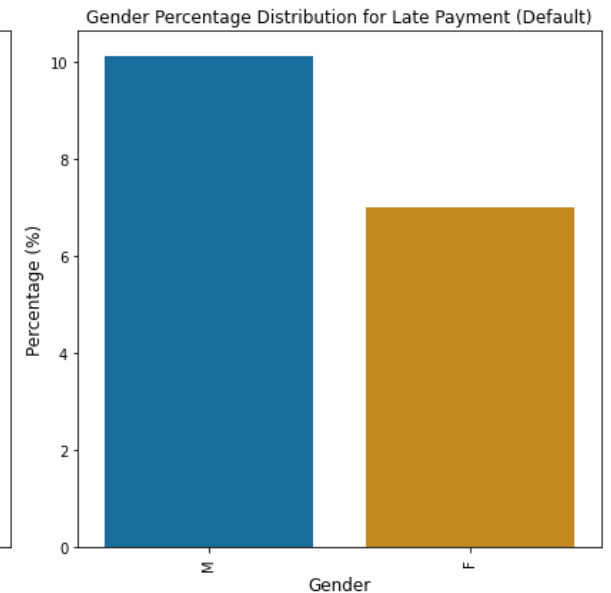
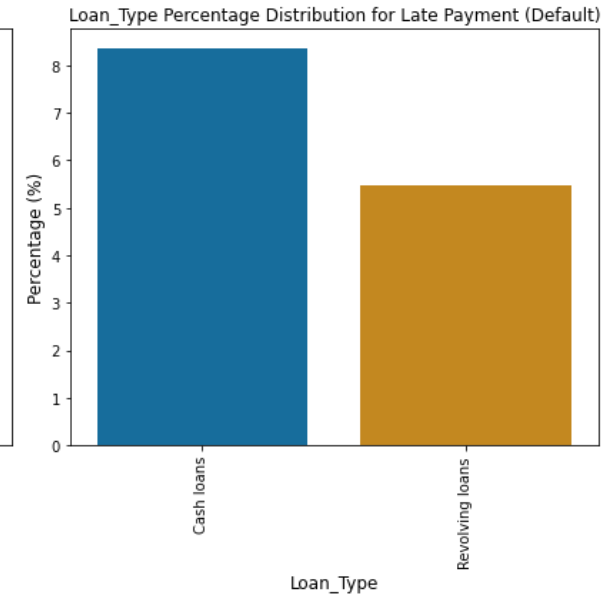
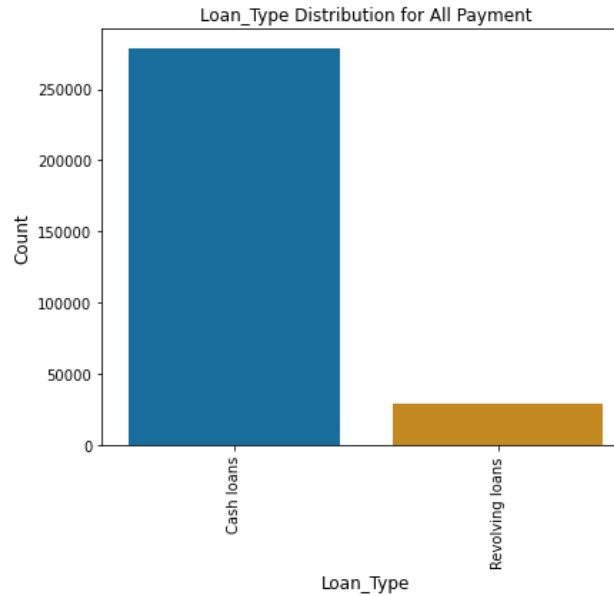
---



# Loan Type, Gender vs. Target Variable

## Findings

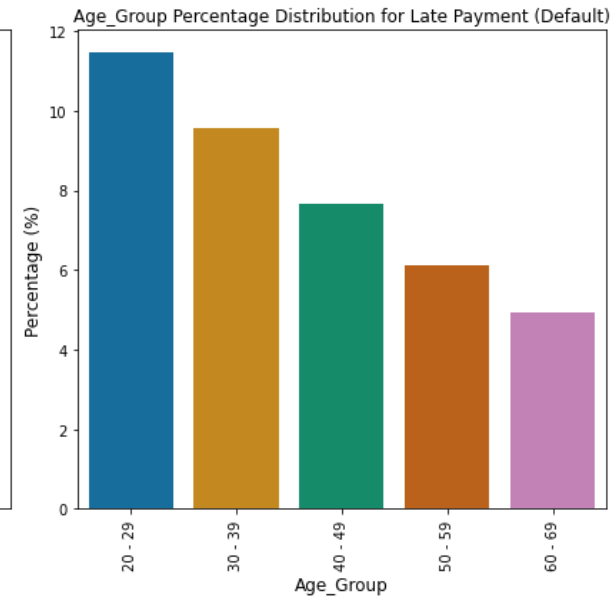
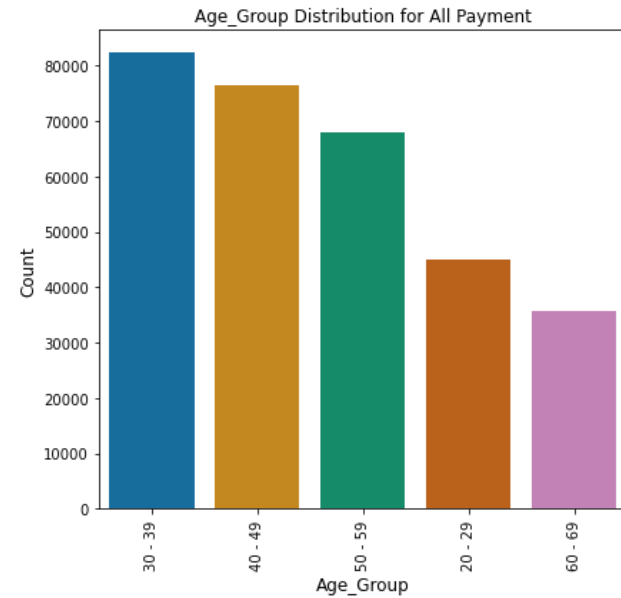
- Cash Loan are the majority for both loans type and contributors to defaults.
- We can see that % of loan non-default of males are low but the defaulted % for Males are higher, this shows that males generally have higher chance of becoming a defaulter than females.



# Age Groups vs. Target Variable

## Findings

- people that are in their 20s and 30s shows higher % in the default graph compared to other age groups. This shows that this group are likely to become a defaulter.



# Income and Loan Groups vs. Target Variable

## Findings

- As you can see from the 2 top pie charts it's clear that major % of loans and its defaulters falls in the Very Low Income Group followed by Moderate.
- In the 2 bottom pie charts, moderate and very low loan group almost shares majority %, while a different case in default, where Low loan group have the majority % followed by Moderate.

Note: The groups are divided as follows:

Very Low: Bottom 25%

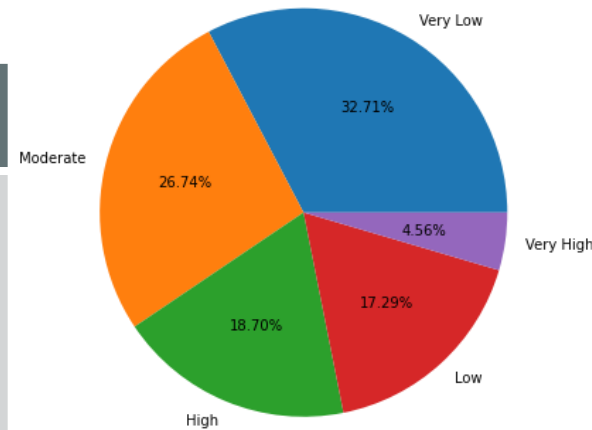
Low: 25th - 50th percentile

Moderate: 50th - 75th percentile

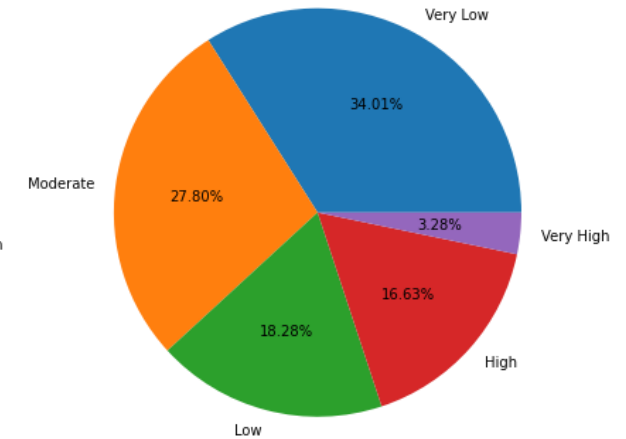
High: 75th - 95th percentile

Very High: Top 5%

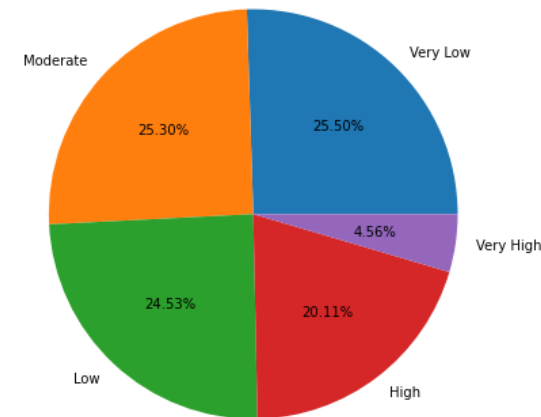
Distribution of All Loans based on Income\_Group



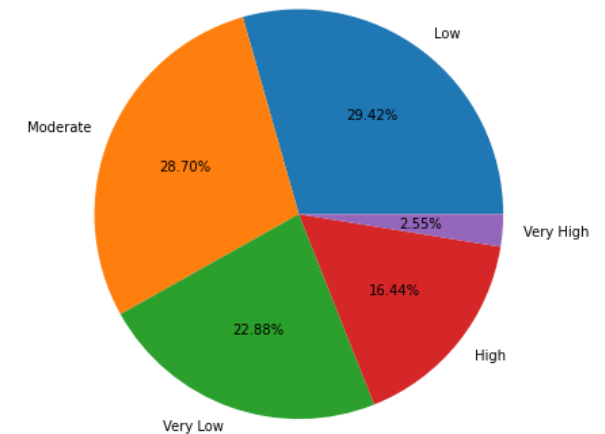
Income\_Group Distribution for Late Payment (Default)



Distribution of All Loans based on Loan\_Group



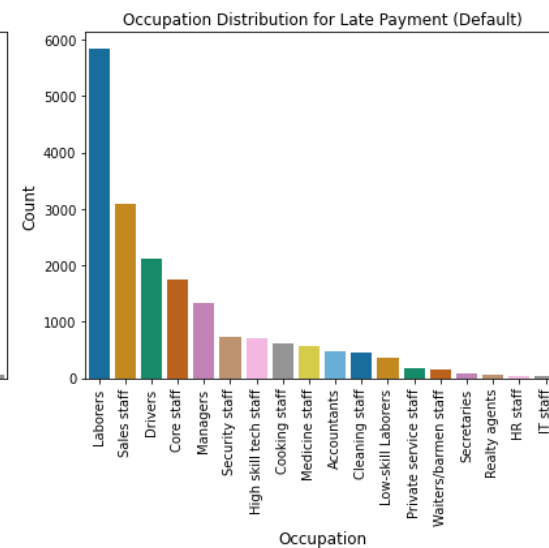
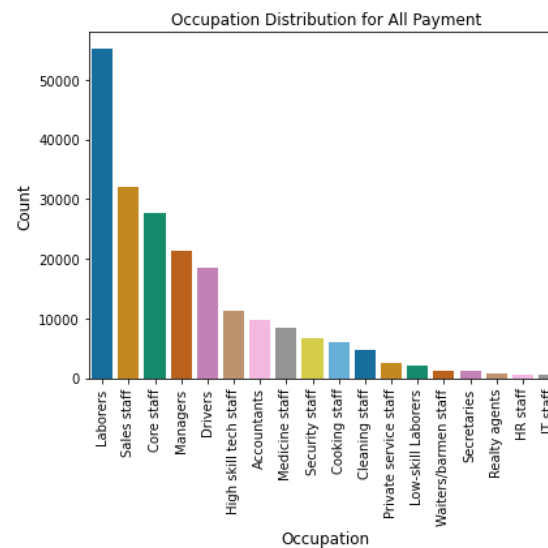
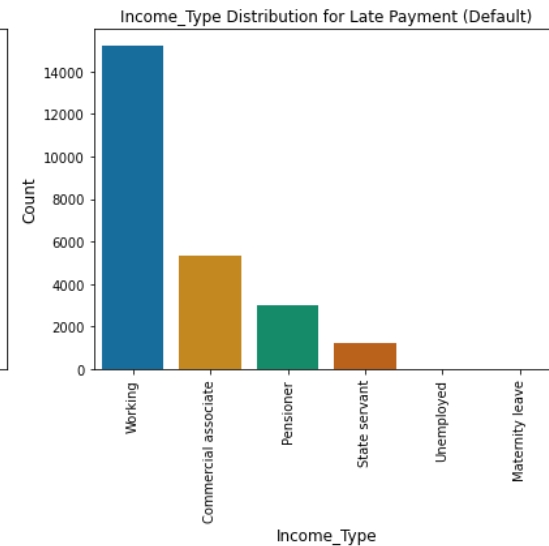
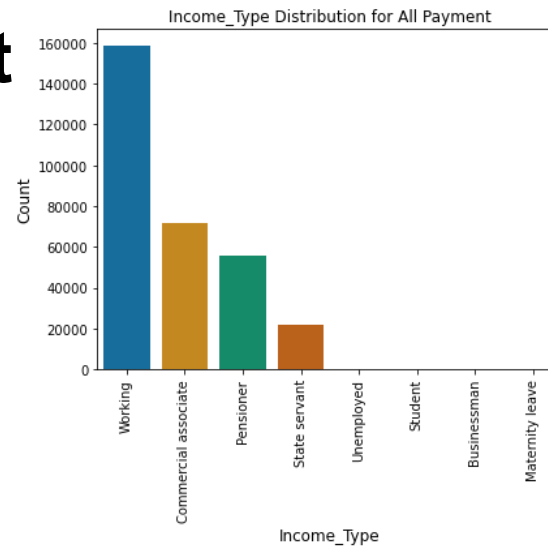
Loan\_Group Distribution for Late Payment (Default)



# Income and Occupation vs. Target Variable

## Findings

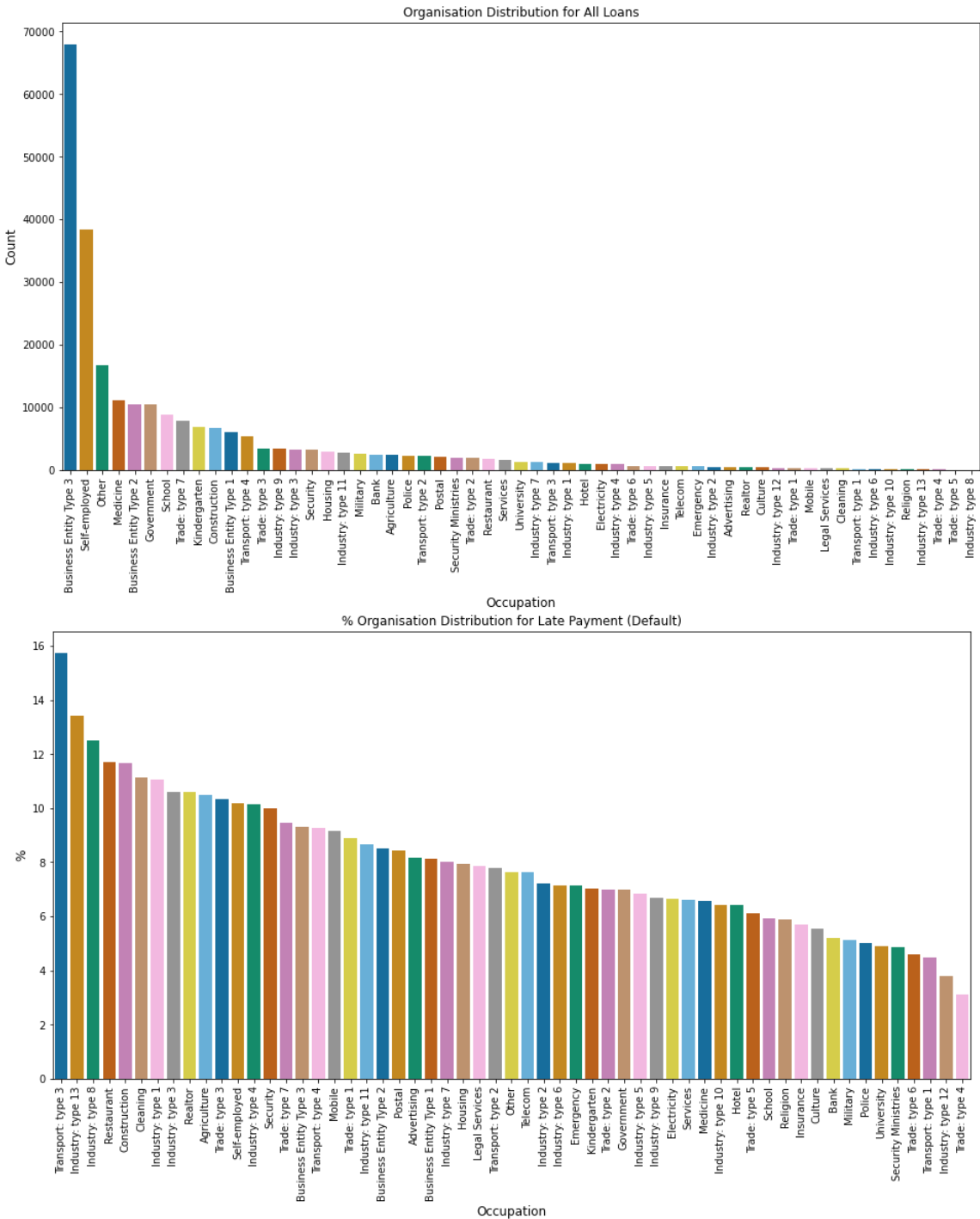
- Top Bar charts: Working people dominates the chart of amount of Loans and defaults: Working people are at the most risk to default
- From bottom bar charts: Laborers are the majority of the loaners and also the most who are defaulted.



# Organisation Type vs. Target Variable

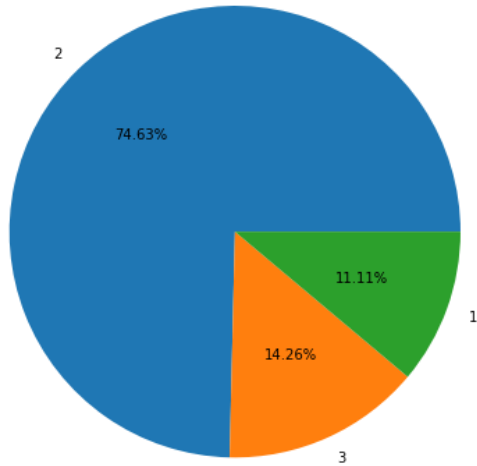
## Findings

- Organisations with the highest % of Defaulters are Transport: Type3 (16%), Industry: Type 13(13.5%) and Industry: Type 8(12.5%).
- Business Entity: Type 3, Trade: type 4 and Industry Type: 12 are the most reliable to be given loan to, as they have a relatively low default compared to loans or low % of loan defaults

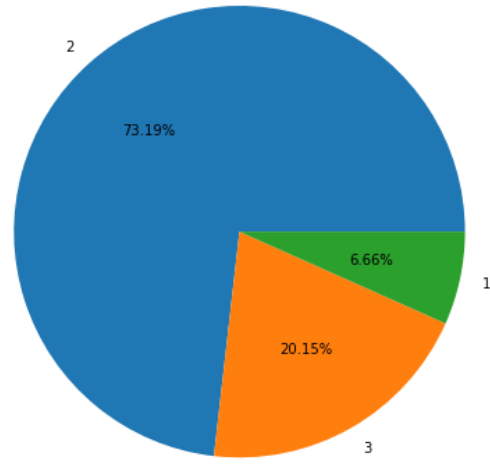


# Area Rating vs. Target Variable

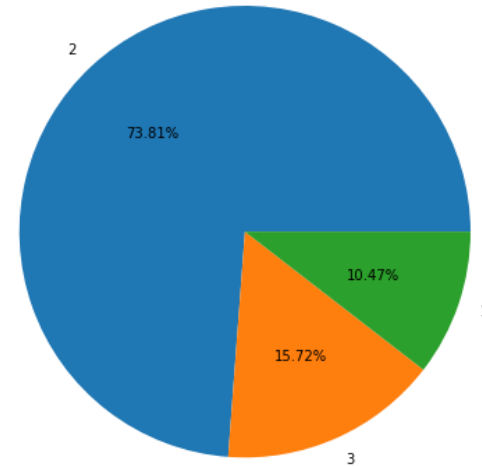
Distribution of All Loans based on City\_Area\_Rating



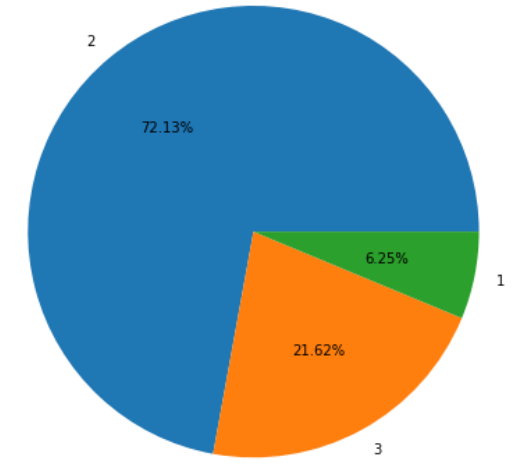
City\_Area\_Rating Distribution for Late Payment (Default)



Distribution of All Loans based on Living\_Area\_Rating



Living\_Area\_Rating Distribution for Late Payment (Default)



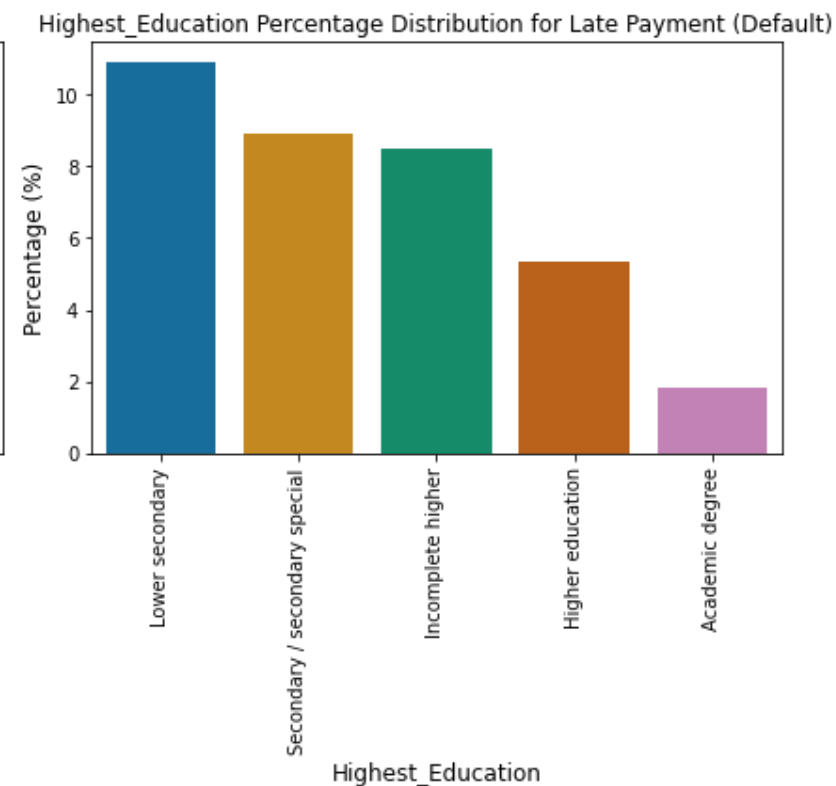
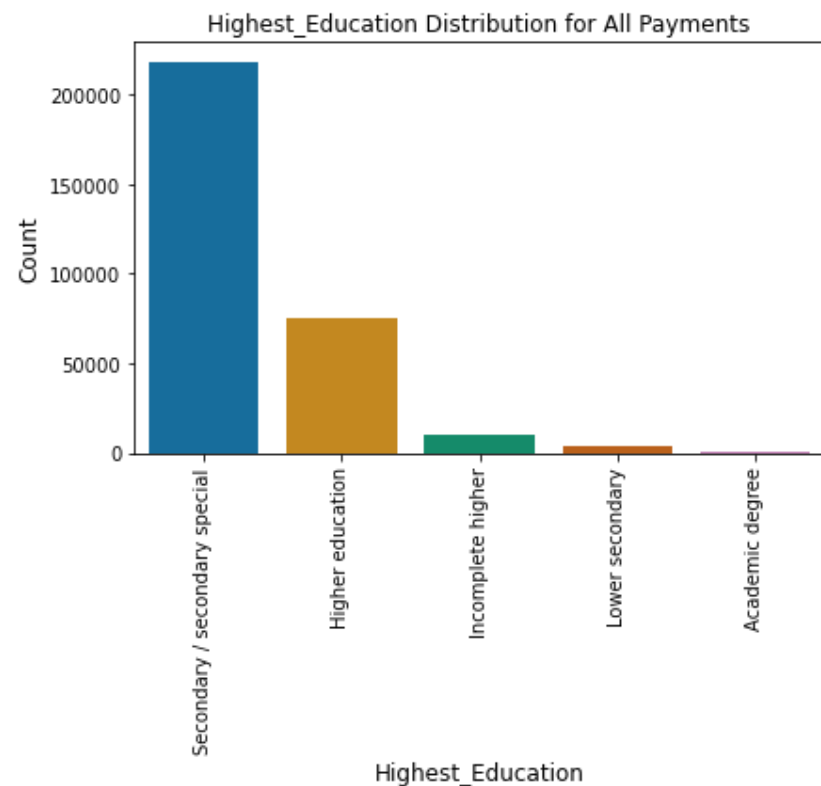
## Findings

- City rating category 2 take the most loan and also have the most default.
- The default % of rating category 3 is higher than In the all loans
- Category 3 City area rating are risky.

## Findings

- Living are rating category 2 take the most loan and also have the most default, the default % of rating category 3 is higher than In the all loans
- Category 3 living area rating are risky.

# Education vs. Target Variables



## Findings

- clients with a secondary education contribute the most to default (10%)

# NUMERICAL ANALYSIS

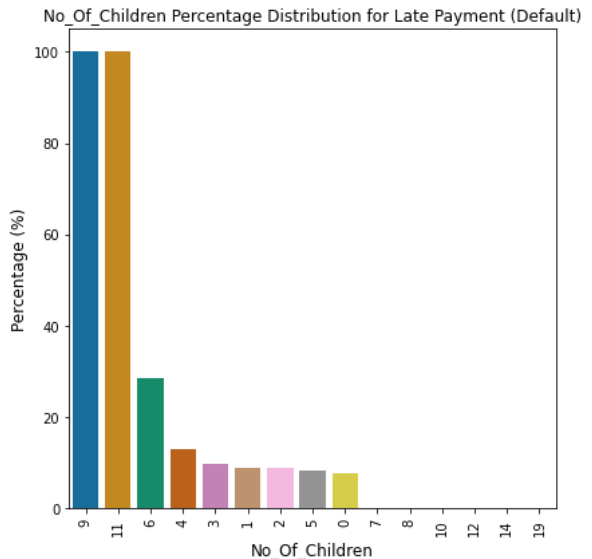
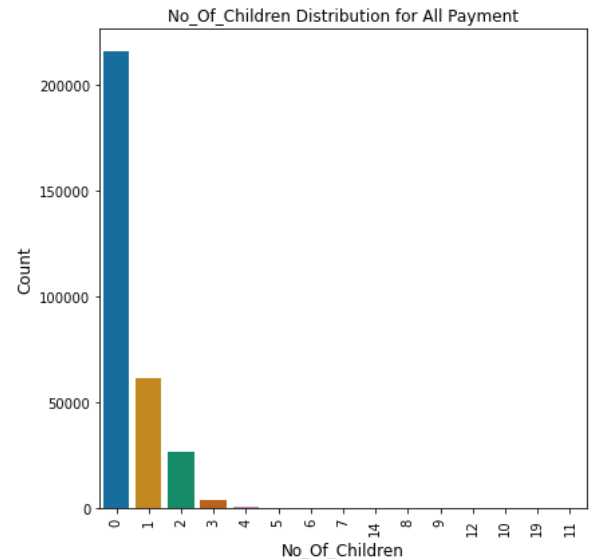
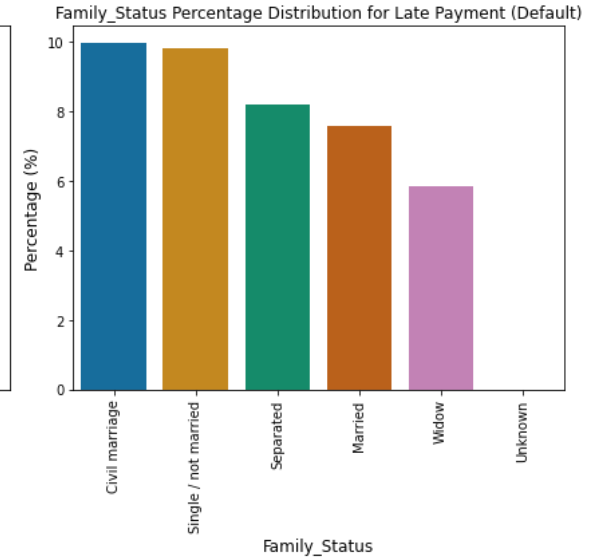
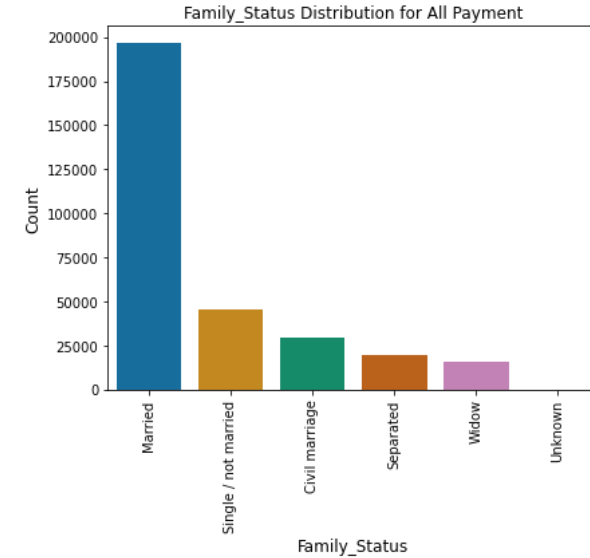
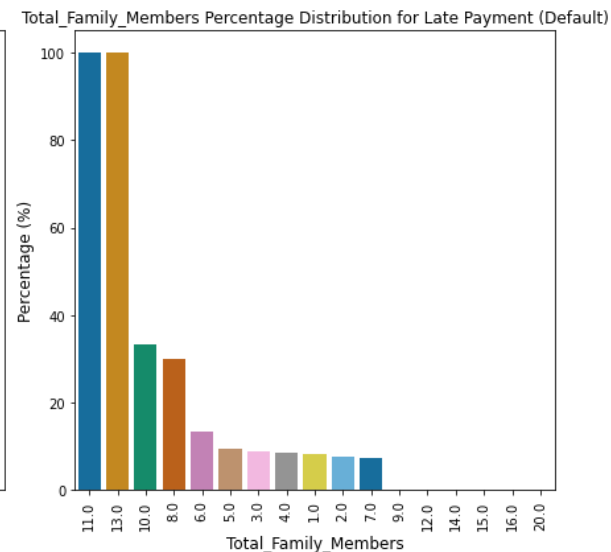
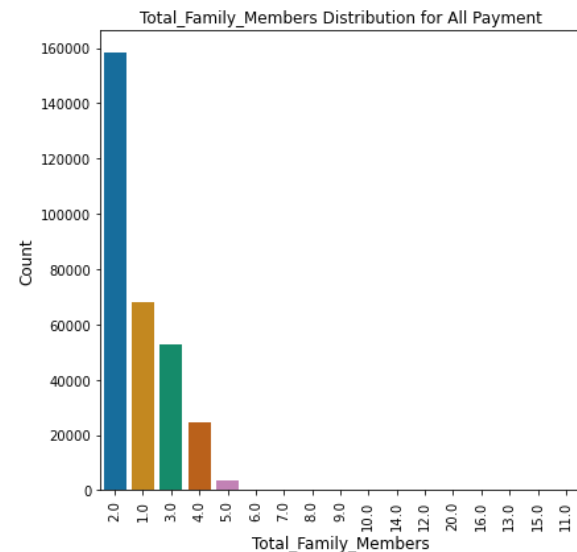
---

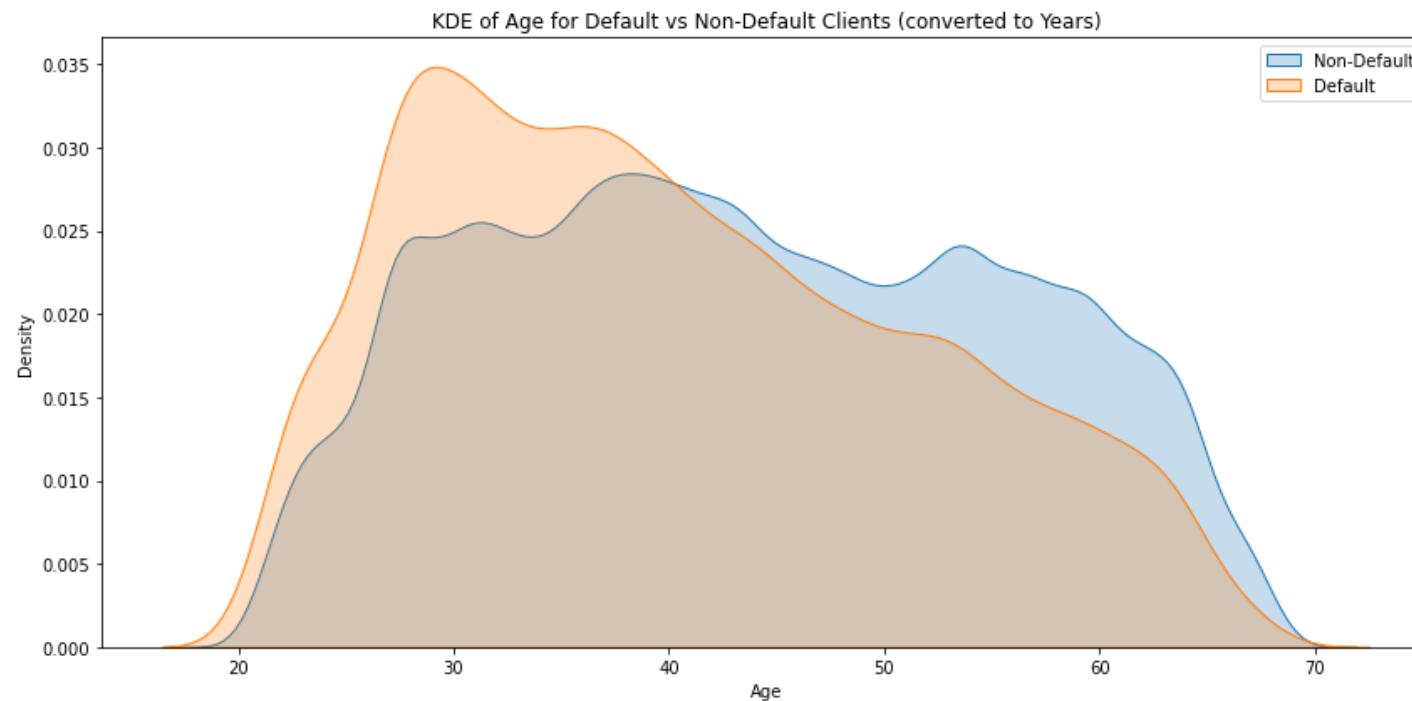


# Family Features

## Findings

- Married clients are take the most loans, Civil marriage have the most defaults (10%), followed by single clients (almost 10%)
- Clients with no or few child will likely to pay their loan on time.
- Clients with higher number of children/total family members are riskier.





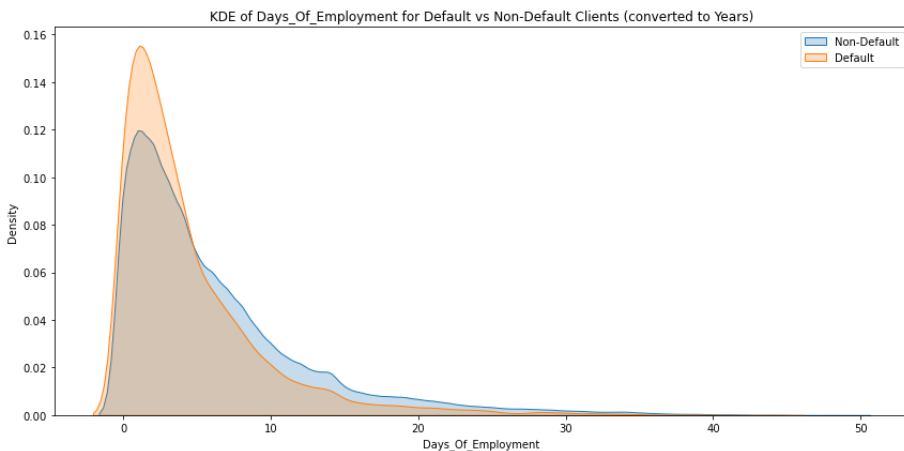
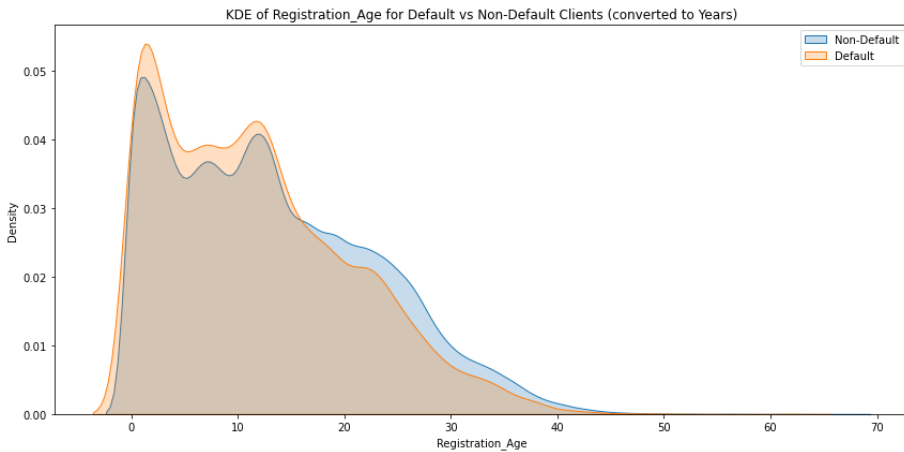
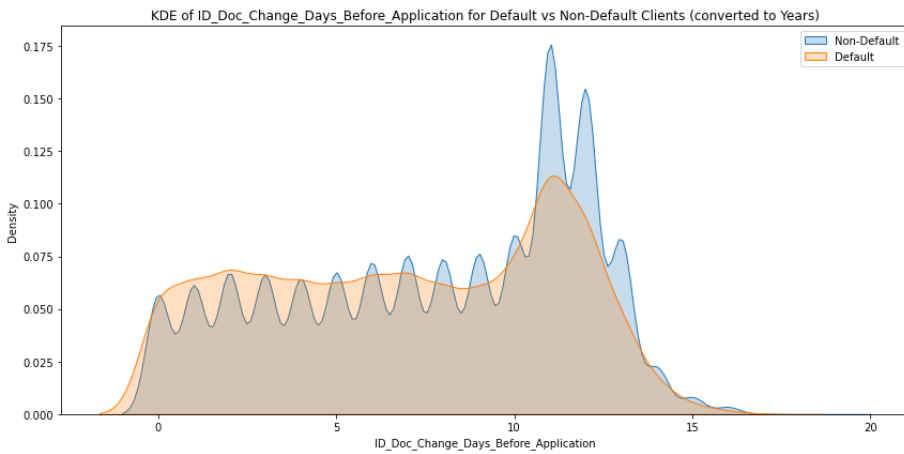
## Findings

### Age distribution:

- As you can see, as the clients are getting older, the clients are often pay they loan on time more often.
- Younger client are less likely to pay on time than older clients.

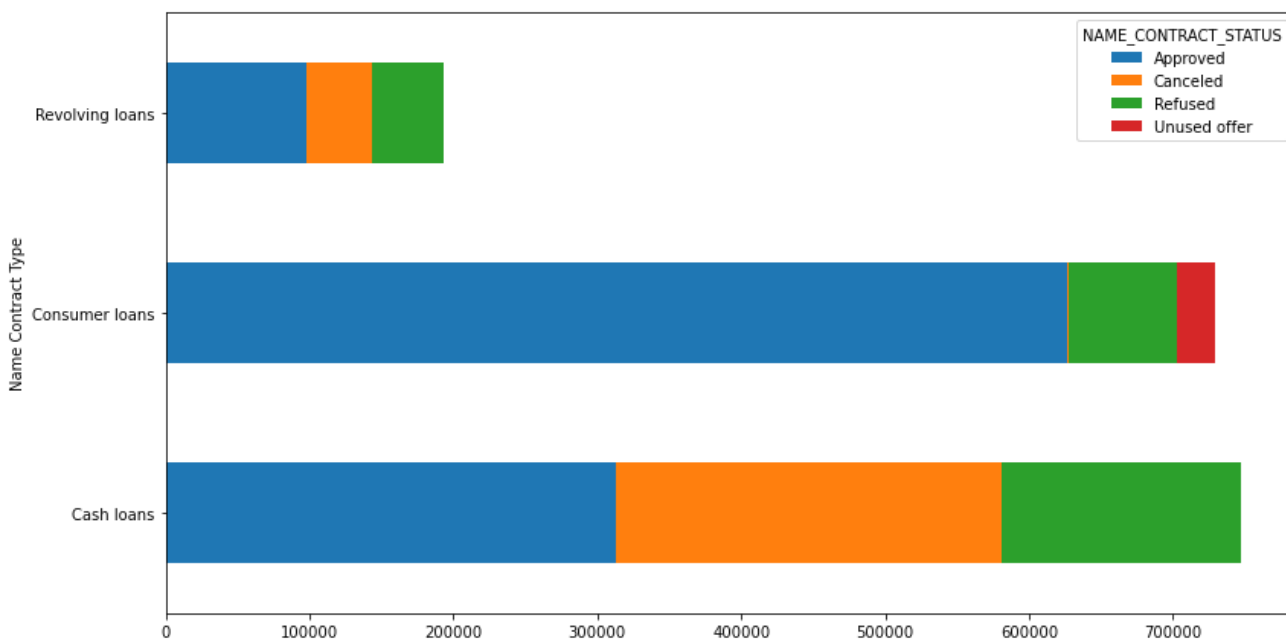
### Age Analysis:

- Consistent with our analysis before, where younger age group have higher amount of defaults than older people.
- Maybe is the difference in financial knowledge, guidance for younger people in financing may be required

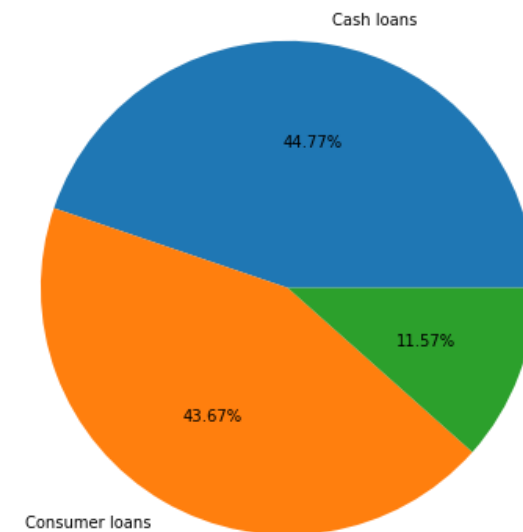


## Findings

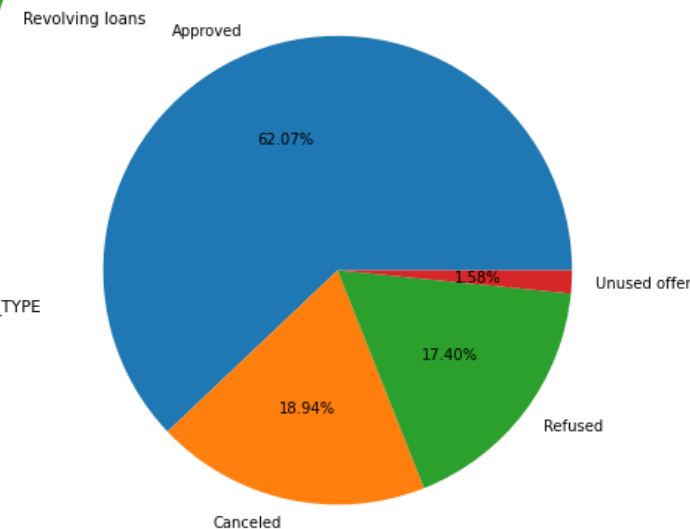
- Clients that change their registration closer to the application are more likely to default
- Clients that change their registration ID closer to the application are less reliable than of those who changed it in advance.
- The more prepared clients, has a higher success in paying back the loan on time.



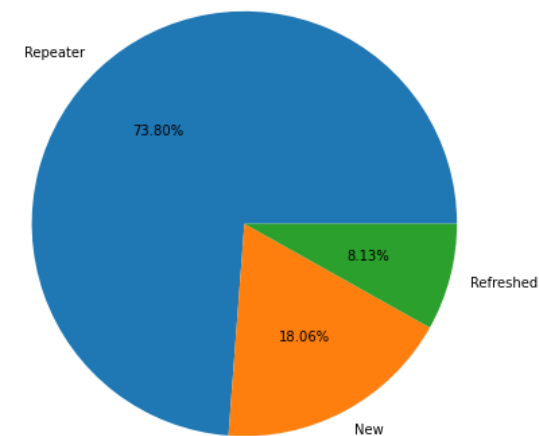
Distribution of All Previous Loans based on NAME\_CONTRACT\_TYPE



Distribution of All Previous Loans based on NAME\_CONTRACT\_STATUS



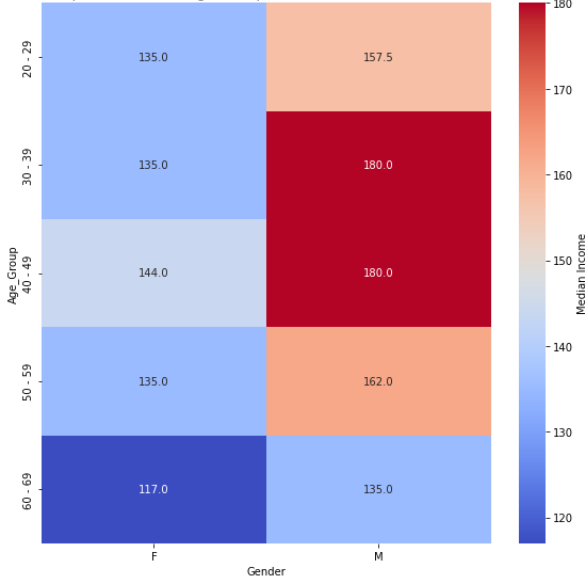
Distribution of All Previous Loans based on NAME\_CLIENT\_TYPE



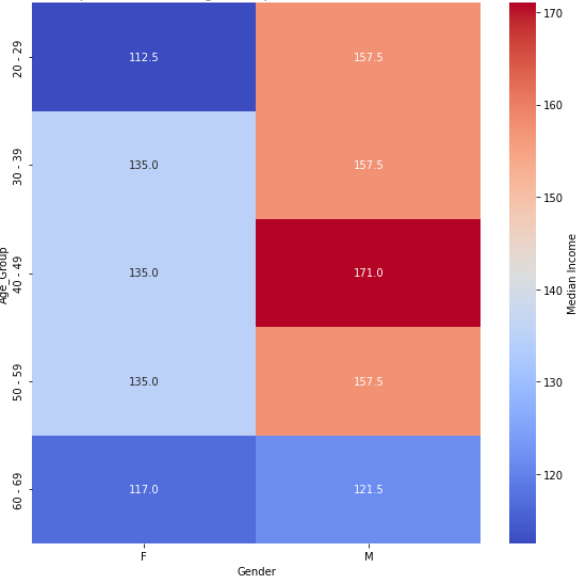
## Findings

- Most of the previous loans are approved.
- Most of the previous loan application are Cash Loans followed by Consumer Loans
- Most of the loans are from Repeaters (those who had been loaned before). Only 18.06% are new.
- Consumer loans are most approved and rarely cancel, they are the most reliable types.

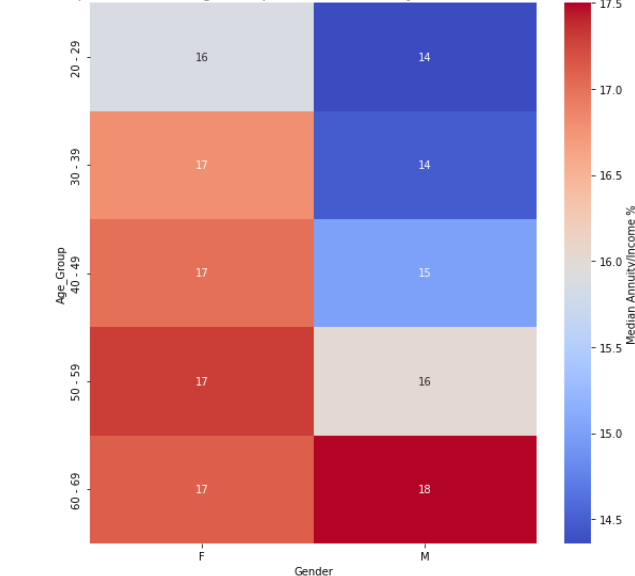
Heatmap of Gender vs Age Group vs Median Income (Non-Defaults)



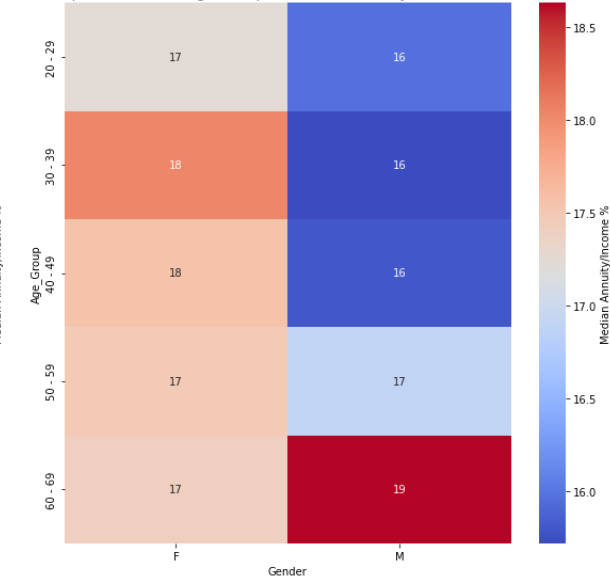
Heatmap of Gender vs Age Group vs Median Income (Defaults)



Heatmap of Gender vs Age Group vs Median Annuity/Income % (Non-Defaults)



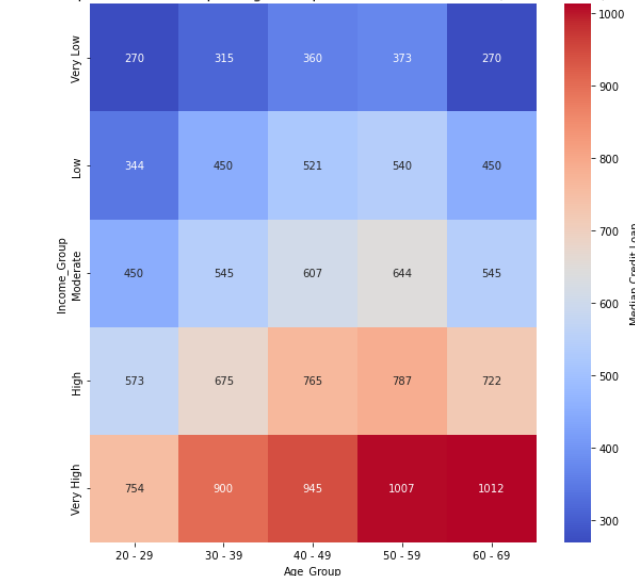
Heatmap of Gender vs Age Group vs Median Annuity/Income % (Defaults)



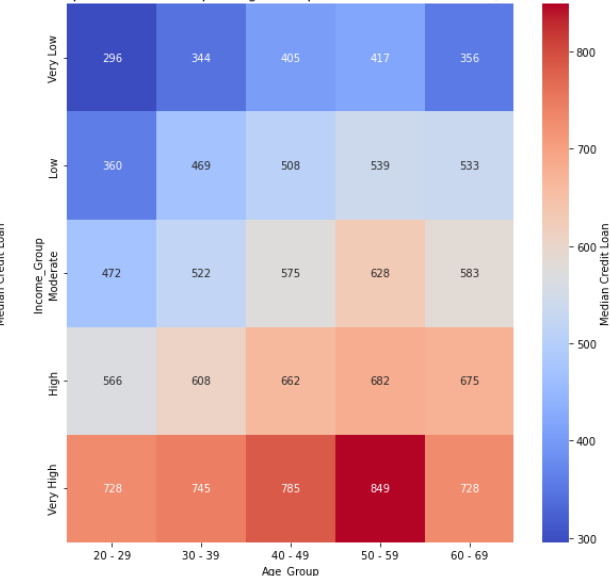
## Findings

- As you can see the Median Income for those who are defaulted are lower than those who are non-defaults.
- The % of Annuity/Income median are higher across Age for defaulters than non-defaulters
- Median Credit Loan of Very low income group of defaulters is high in all age segment as compared to the Non-Default Loans.

Heatmap of Income Group vs Age Group vs Median Credit Loan (Non-Defaults)



Heatmap of Income Group vs Age Group vs Median Credit Loan (Defaults)



# CONCLUSION

---

- Most defaulters are from very low and low income range.
- Younger People are more tend to default.
- Focus more on having client from category 1 & 2 city and living rating area.
- Laborers, Sales Staff, and Drivers are most defaulters.
- Focus less on clients who has Working Income types.
- Attract more repeating clients.
- Approve loan to more prepared clients.
- Clients with more children/family members are likely to default