

# Credit EDA Assignment

1

**Hemavathi AB**  
**ab.hemavathi@gmail.com**



Topic	Slide No
Business Case Study :	
1. Problem Statement 1	03
2. Problem Statement 2	06
Data :	
Application Data : Cleaning , Analysis	08 - 38
Previous Application Data : Cleaning , Analysis	39 – 45
Merging Both Data's : Approach and Solution	46 - 55
Insights: Inferences and Solutions to Problem Statements	56 – 59

# Problem Statement 01

This assignment aims to give an idea of applying EDA in a real business scenario. In this assignment, apart from applying the techniques that is learnt in the EDA module, should also develop a basic understanding of risk analytics in banking and financial services and understand how data is used to minimise the risk of losing money while lending to customers.

# Business Objective:

4

This case study aims to identify patterns which indicate if a client has difficulty paying their instalments which may be used for taking actions such as denying the loan, reducing the amount of loan, lending (to risky applicants) at a higher interest rate, etc. This will ensure that the consumers capable of repaying the loan are not rejected. Identification of such applicants using EDA is the aim of this case study



# Data Understanding

5 The data given below contains the information about the loan application at the time of applying for the loan. It contains two types of scenarios:

1. The client with payment difficulties: he/she had late payment more than X days on at least one of the first Y instalments of the loan in our sample,
2. All other cases: when the payment is paid on time.

When a client applies for a loan, there are four types of decisions that could be taken by the client/company):

1. Approved: The Company has approved loan Application
2. Cancelled: The client cancelled the application sometime during approval. Either the client changed her/his mind about the loan or in some cases due to a higher risk of the client, he received worse pricing which he did not want.
3. Refused: The company had rejected the loan (because the client does not meet their requirements etc.).
4. Unused offer: Loan has been cancelled by the client but at different stage the process.

# Problem Statement 02

## Expectation from Learners:

1. Present the overall approach of the analysis in a presentation. Mention the problem statement and the analysis approach briefly.
2. Identify the missing data and use appropriate method to deal with it. (Remove columns/or replace it with an appropriate value)
3. Identify if there are outliers in the dataset. Also, mention why do you think it is an outlier. Again, remember that for this exercise, it is not necessary to remove any data points.
4. Identify if there is data imbalance in the data. Find the ratio of data imbalance.
5. Analysis for the 'Target variable' in the dataset one can plot in terms of percentage or absolute value to show mix of univariate and bivariate analysis.

Explain the results of univariate, segmented univariate, bivariate analysis, etc. in business terms.

6 Find the top 10 correlation for the Client with payment difficulties and all other cases (Target variable). Note that you have to find the top correlation by segmenting the data frame w.r.t to the target variable and then find the top correlation for each of the segmented data and find if any insight is there.

# Problem Statement 02

## Expectation from Learners:

8. Include visualisations and summarise the most important results in the presentation. Insights should explain why the variable is important for differentiating the clients with payment difficulties with all other cases.
9. You need to submit one/two Ipython notebook which clearly explains the thought process behind your analysis (either in comments of markdown text), code and relevant plots.
9. The presentation file needs to be in PDF format and should contain the points discussed above with the necessary visualisations. Also, all the visualisations and plots must be done in Python(should be present in the Ipython notebook), though they may be recreated in Tableau for better aesthetics in the PPT file.

# Application Data: Cleaning

1. 22 columns were reduced to 73 as 49 columns had null values which were more than 32%.
2. Dropped 26 columns further as they were client documents and contact information. Hence the number of active columns became 47.
3. OCCUPATION\_TYPE column had 31.3% null values marked them as "Unknown"
4. Replaced null values with median statistics for the following columns : EXT\_SOURCE\_3, EXT\_SOURCE\_2, AMT\_GOODS\_PRICE, AMT\_ANNUITY

8

Replaced null values with mode statistics for the following columns :

AMT\_REQ\_CREDIT\_BUREAU\_YEAR, AMT\_REQ\_CREDIT\_BUREAU\_QRT,  
AMT\_REQ\_CREDIT\_BUREAU\_MON, AMT\_REQ\_CREDIT\_BUREAU\_WEEK,  
AMT\_REQ\_CREDIT\_BUREAU\_DAY, AMT\_REQ\_CREDIT\_BUREAU\_HOUR, NAME\_TYPE\_SUITE,  
OBS\_30\_CNT\_SOCIAL\_CIRCLE, DEF\_30\_CNT\_SOCIAL\_CIRCLE, OBS\_60\_CNT\_SOCIAL\_CIRCLE,  
DEF\_60\_CNT\_SOCIAL\_CIRCLE, CNT\_FAM\_MEMBERS, DAYS\_LAST\_PHONE\_CHANGE



# Application Data: Anomalies

1. Changed negative values to absolute for the following columns :  
DAYS\_LAST\_PHONE\_CHANGE, DAYS\_BIRTH, DAYS\_REGISTRATION, DAYS\_EMPLOYED, DAYS\_ID\_PUBLISH
2. CODE\_GENDER column had XNA replaced them with mode statistics "F"
3. ORGANIZATION\_TYPE column had XNA replaced them with "Unknow"
4. Following columns had outliers so did binning : AMT\_INCOME\_TOTAL, AMT\_CREDIT, AMT\_GOODS\_PRICE, AMT\_ANNUITY
5. Since DAYS\_BIRTH has absolute values to predict the Age divided it by 365 . Then did binning to this column too.
6. Added 6 extra columns: AGE, AGE\_RANGE, AMT\_INCOME\_TOTAL\_RANGE, AMT\_CREDIT\_RANGE, AMT\_GOODS\_PRICE\_RANGE, AMT\_ANNUITY\_RANGE
7. Final total number of rows and columns is : 307511 rows and 53 columns

# Application Data: Analysis

TARGET is the dependent variable as it acts as first column to be checked if its 0 or 1 . If 0 then client with payment difficulties and if 0 then client with no payment difficulties.

Diving the application data (df\_app) into two data frames df0 (TARGET = 1 , client with payment difficulties )

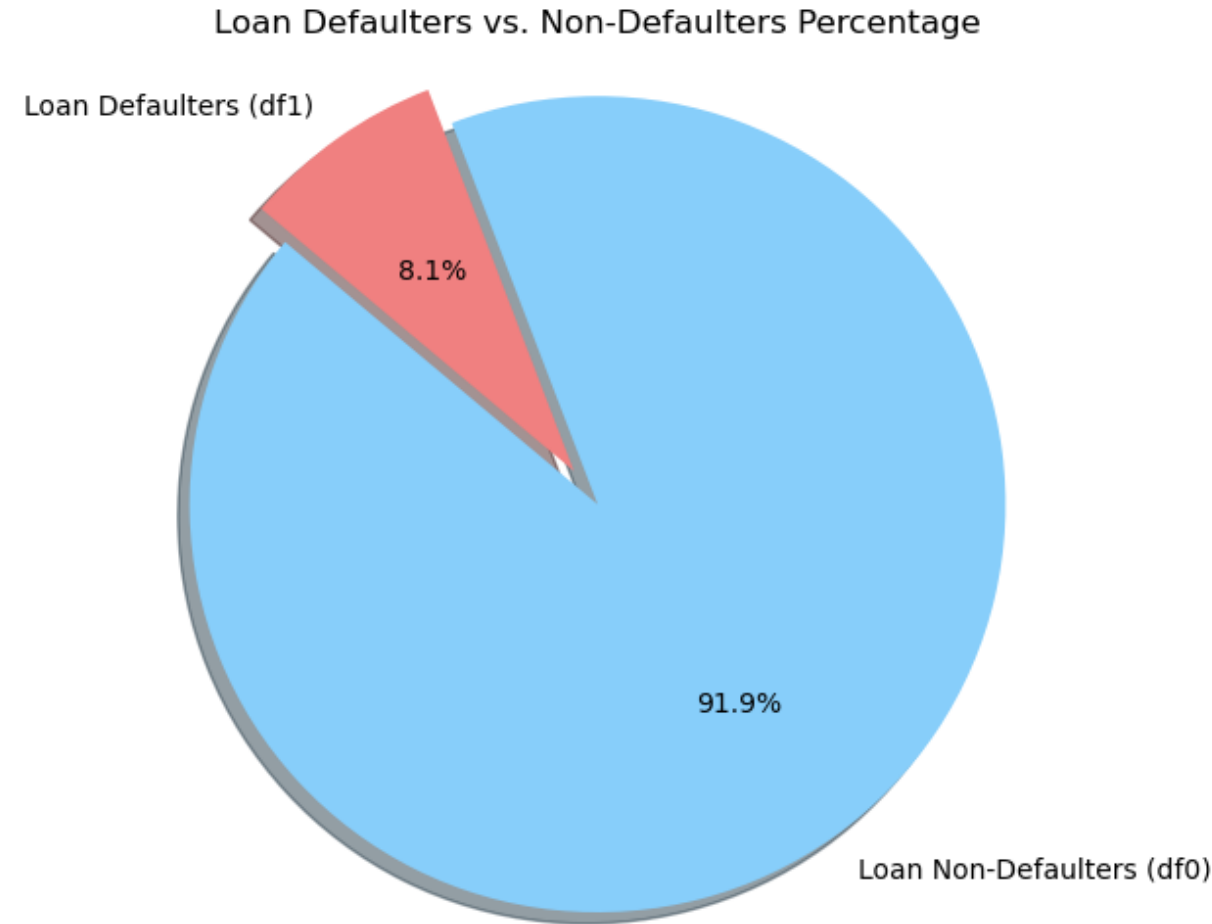
10

df1 (TARGET = 0 , all others)

Total rows with TARGET = 0 is 282,686  
(Non-Defaulters)

Total rows with TARGET = 1 is 24,825(Defaulters)

Hence, we see from the graph that there are **8.1% population who are Defaulters.**



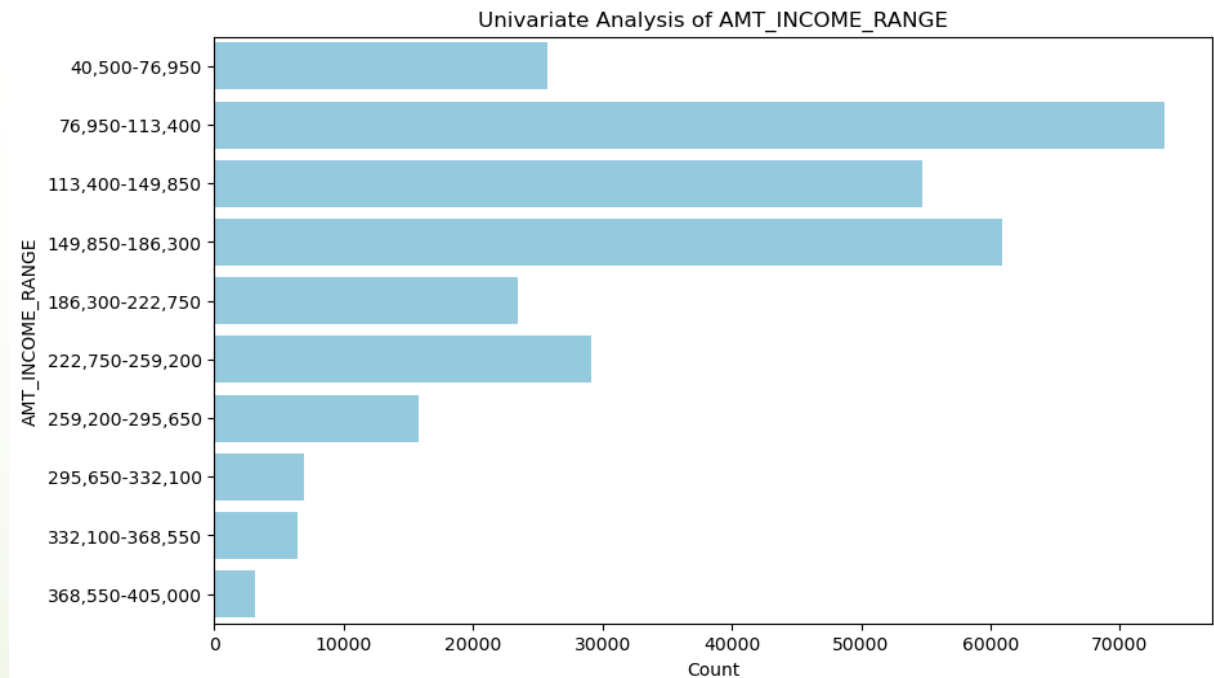
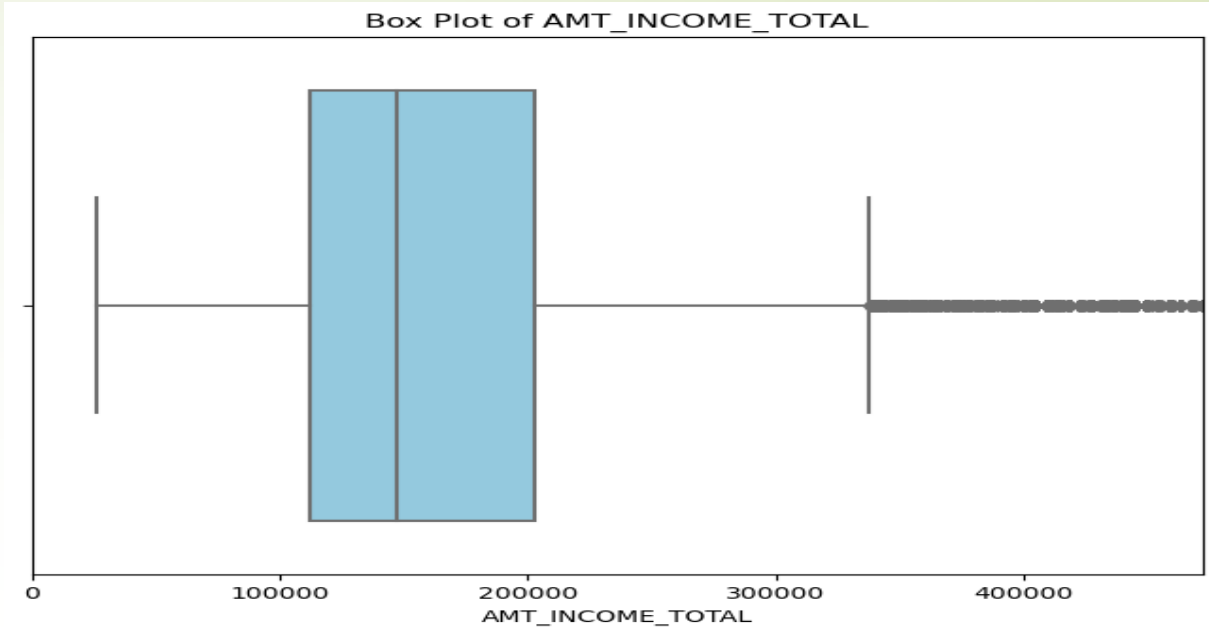
# Univariate Analysis

11

## AMT\_INCOME:

As the raw data was skewed towards 75<sup>th</sup> percentile and has outliers, binning gave it more meaningful information.

1. We see a greater number of clients were in the income range 76 to 113K
2. Second largest were in the income range 149 to 186K followed by 113 to 149K



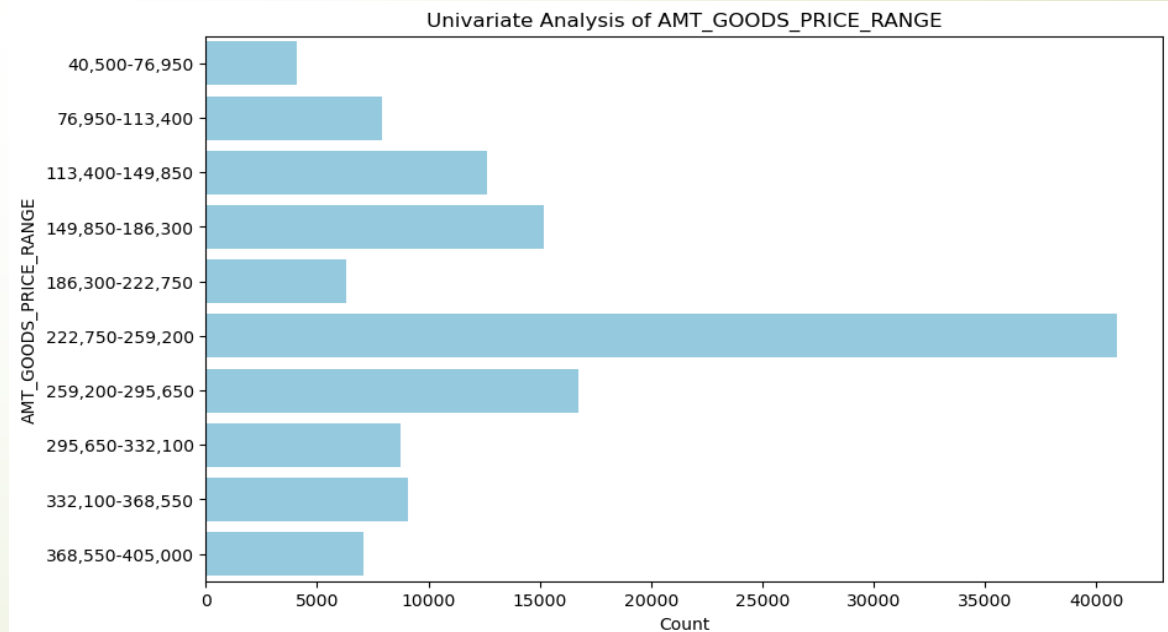
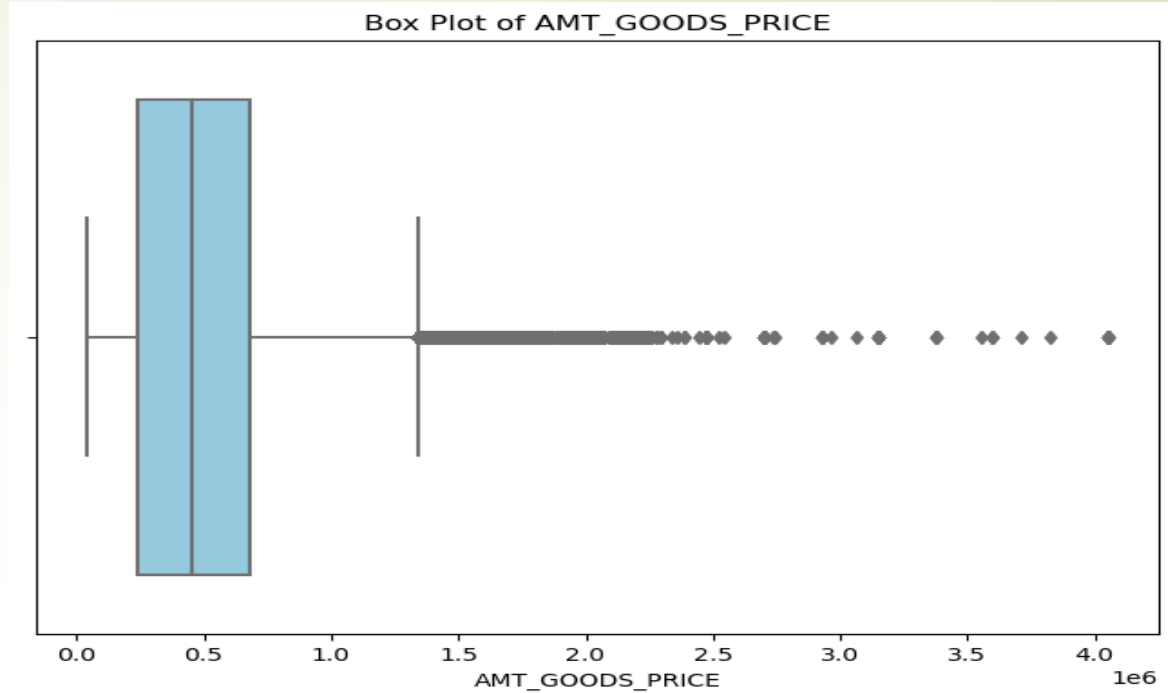
# Univariate Analysis

12

## GOODS\_PRICE:

As the raw data was skewed towards 75<sup>th</sup> percentile and has outliers, binning gave it more meaningful information .

1. We see a greater number of clients took loan for the goods which were priced in the range of 222 to 259K
2. Second largest were in the range of 259 to 295K and 149 to 186K
3. Third was in the range of 113 to 149K





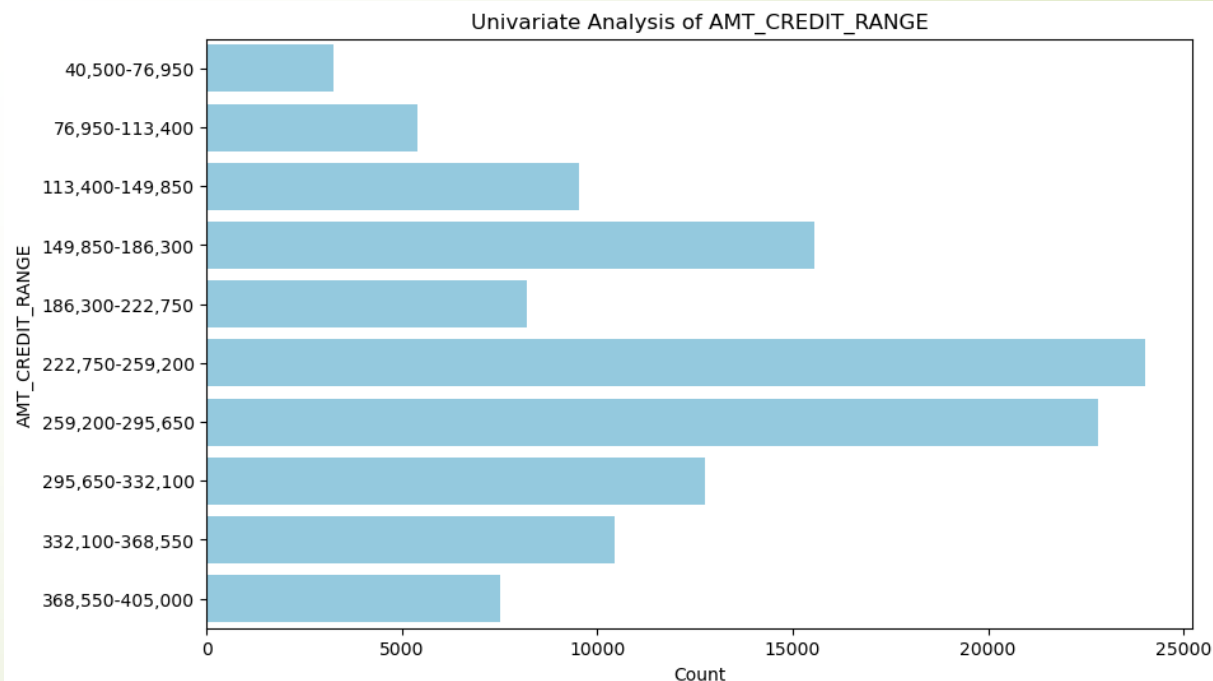
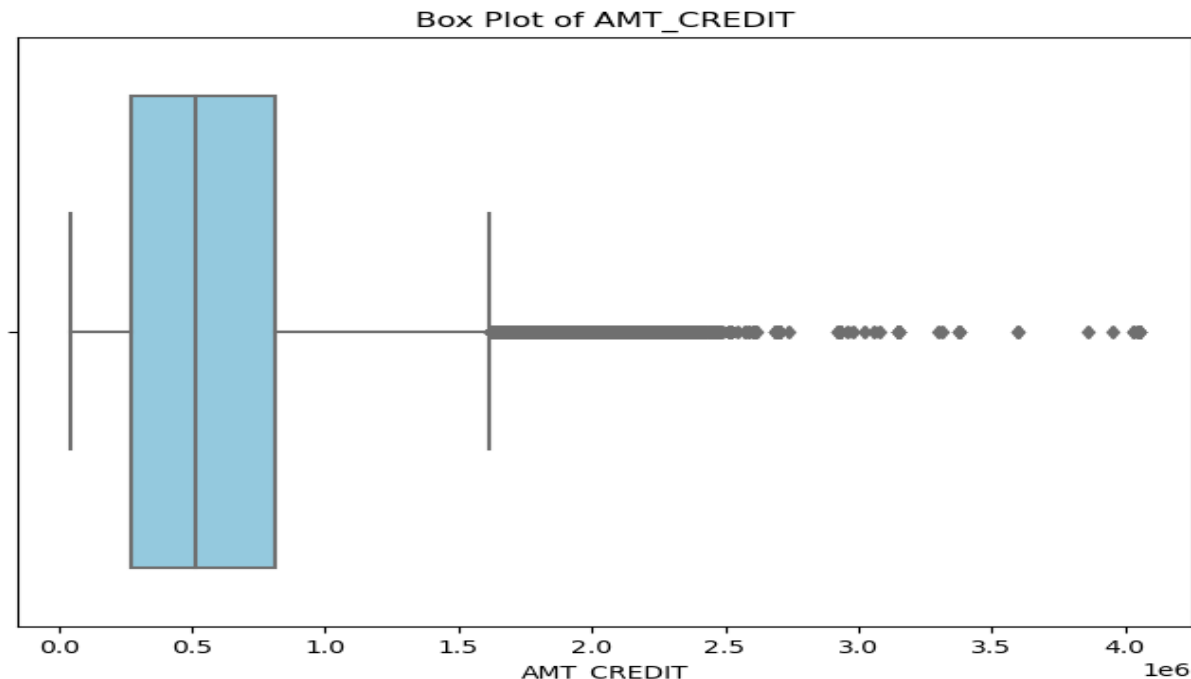
# Univariate Analysis

13

## AMT\_CREDIT:

As the raw data was skewed towards 75<sup>th</sup> percentile and has outliers, binning gave it more meaningful information.

1. We see a greater number of clients took credit in the range of 222 to 259K
2. Second largest were in the range 259 to 295K followed by 140 to 186K



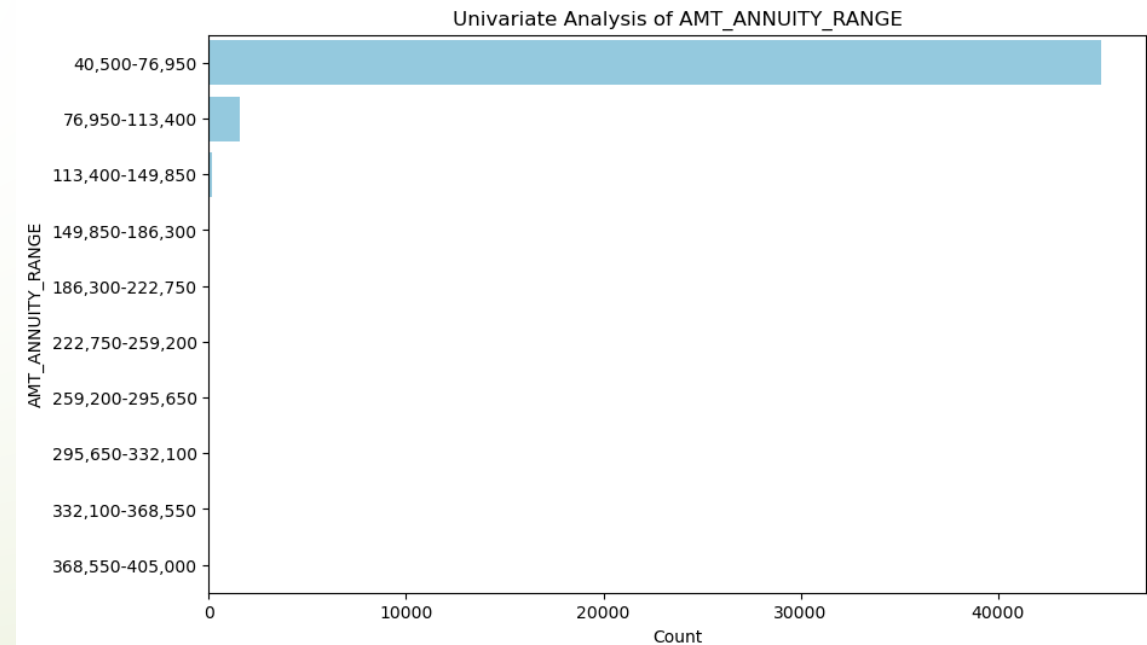
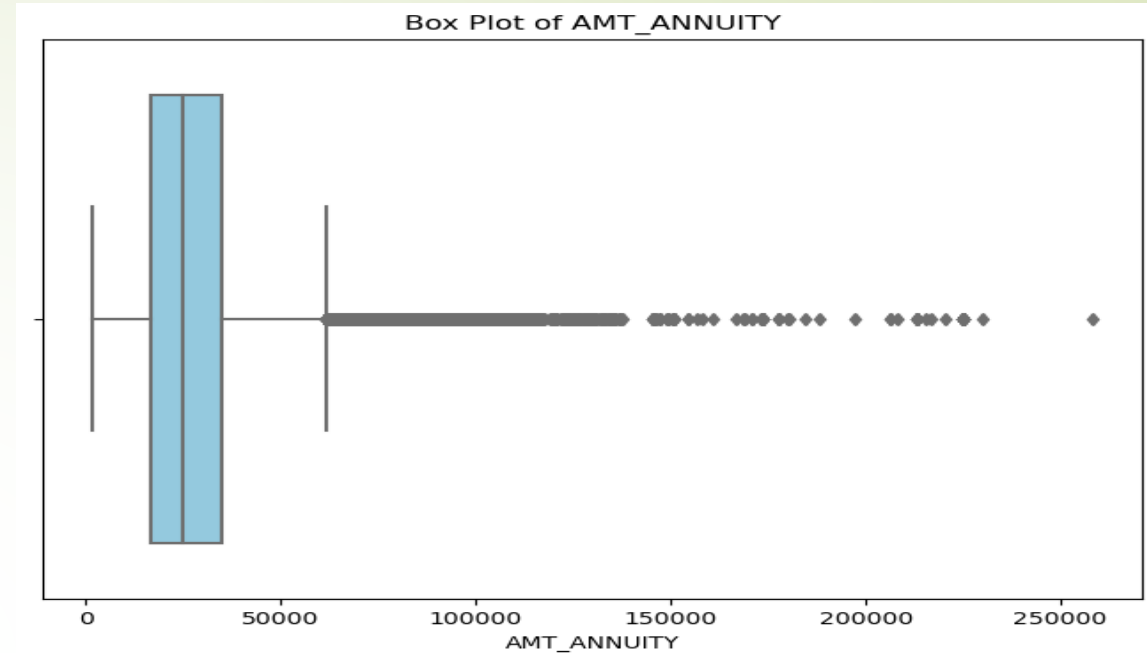
# Univariate Analysis

14

## AMT\_ANNUITY:

As the raw data was skewed towards 75<sup>th</sup> percentile and has outliers, binning gave it more meaningful information .

1. We see a greater number of clients were in the loan annuity range 40 to 76K followed by 76 to 113K



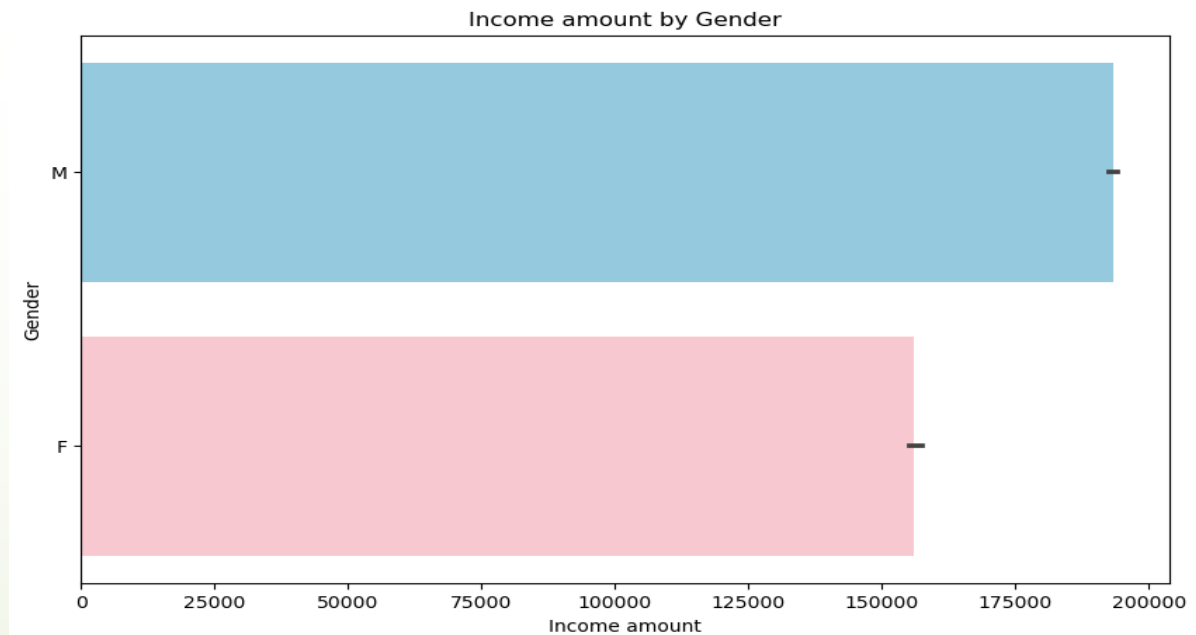
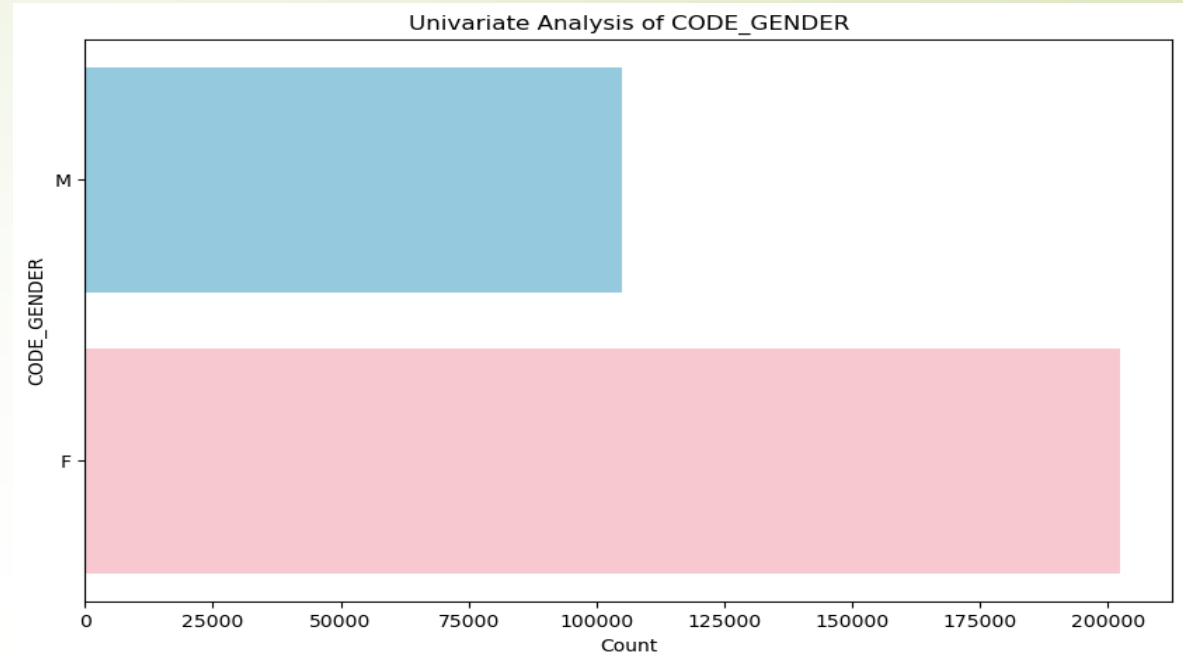
# Univariate Analysis

15

## CODE\_GENDER:

We see more number of female population in the group then male who have taken loan.

But when compared with the income that both gender take , male population earns significantly higher than females.

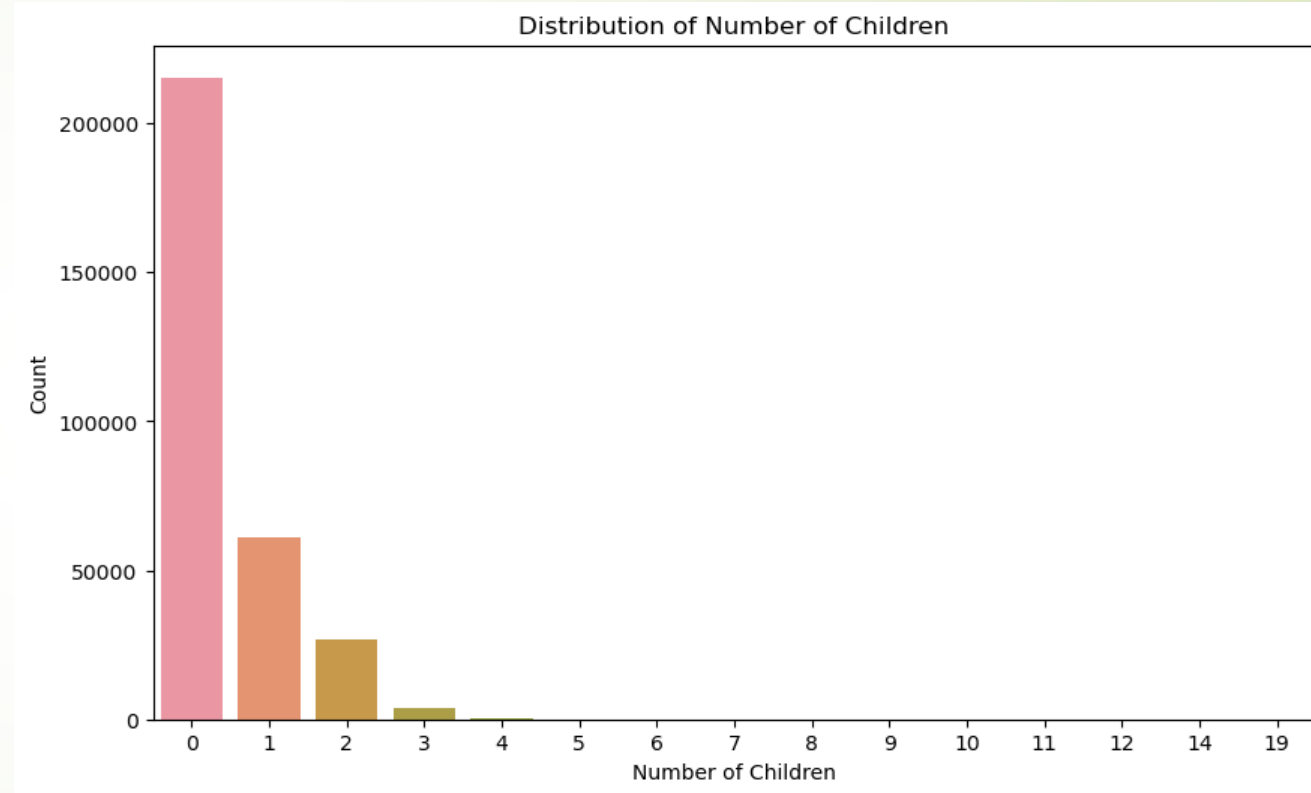


# Univariate Analysis

16

## CNT\_CHILDREN:

We see significantly high number of clients who do not have children.



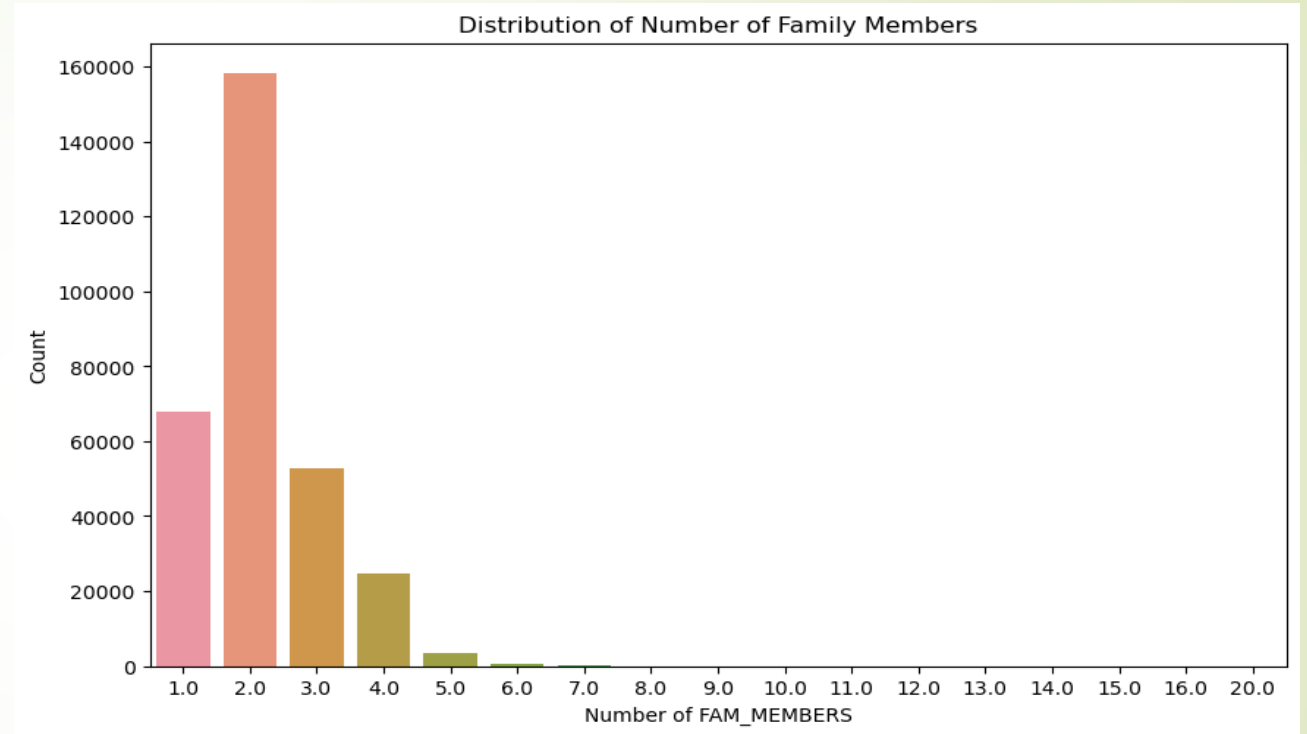


# Univariate Analysis

17

## CNT\_FAM\_MEMBERS:

We see a significantly large number of clients who have 2 family members

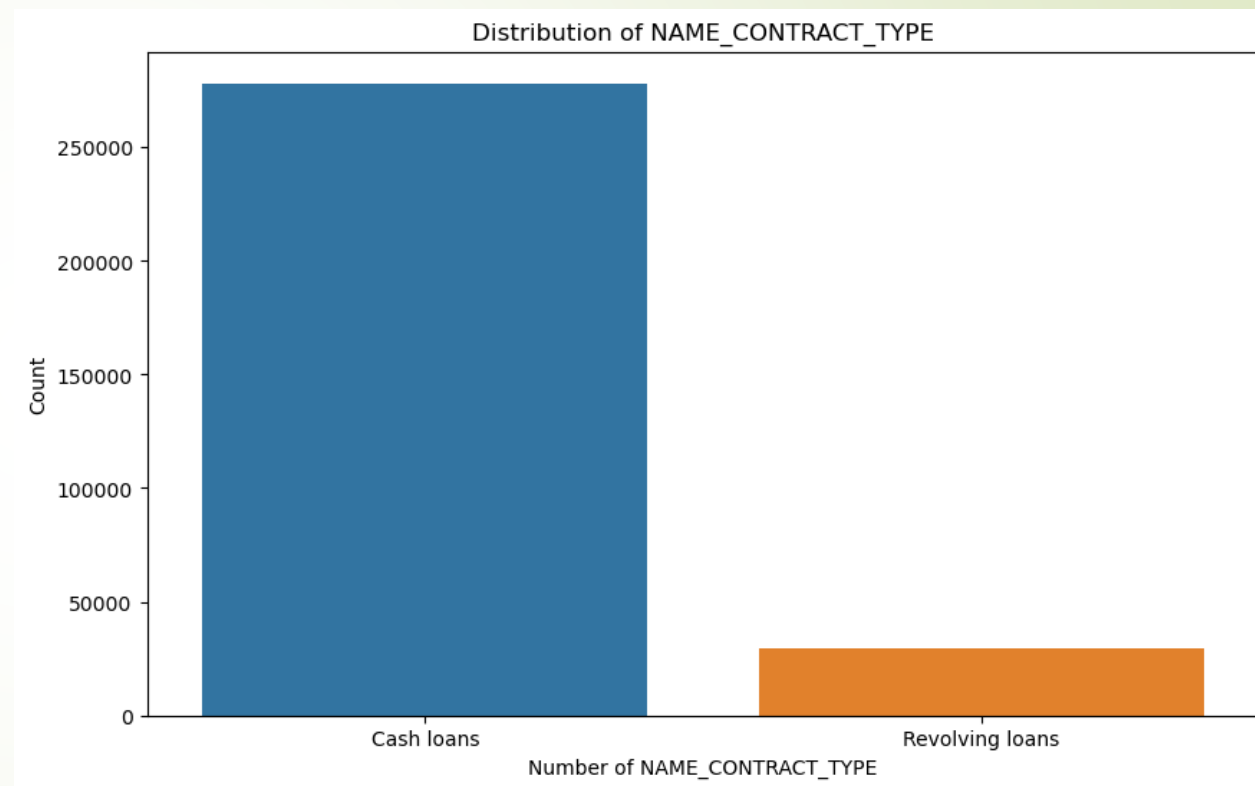


# Univariate Analysis

18

## NAME\_CONTRACT\_TYPE:

We see a significantly large number of clients who have availed Cash Loans

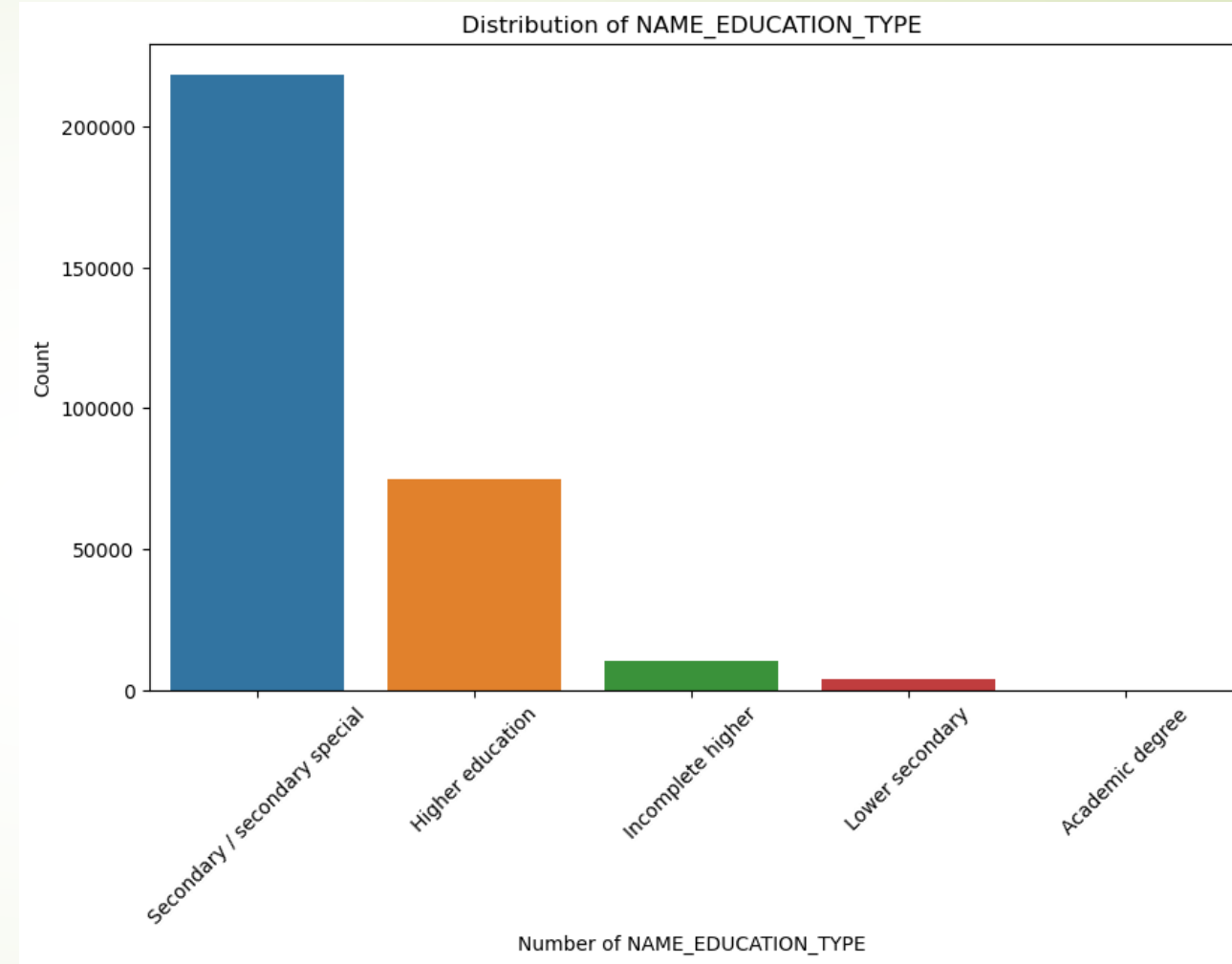


# Univariate Analysis

19

## NAME\_EDUCATION\_TYPE:

We see a significantly large number of clients who are Secondary/secondary special education group who have taken loan

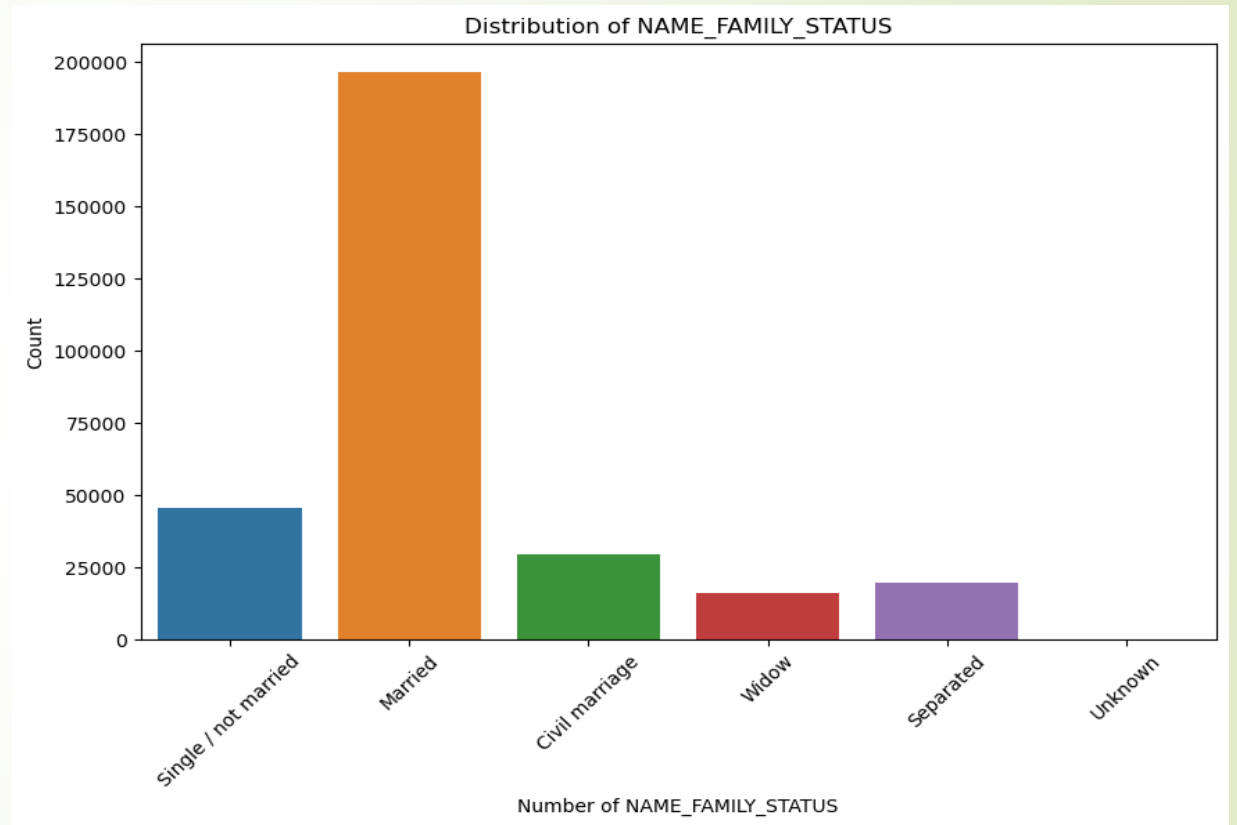


# Univariate Analysis

20

## NAME\_FAMILY\_STATUS

We see a significantly large number of clients who are from Married group have taken loan



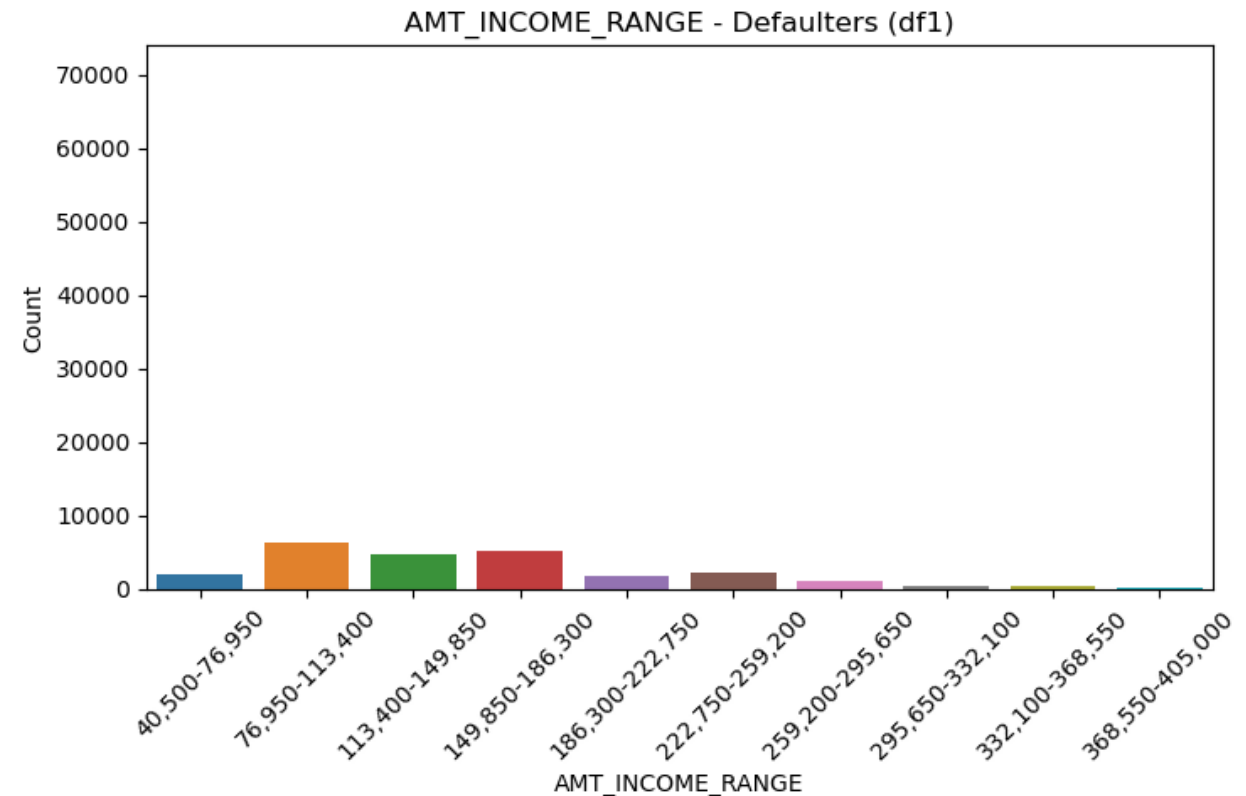
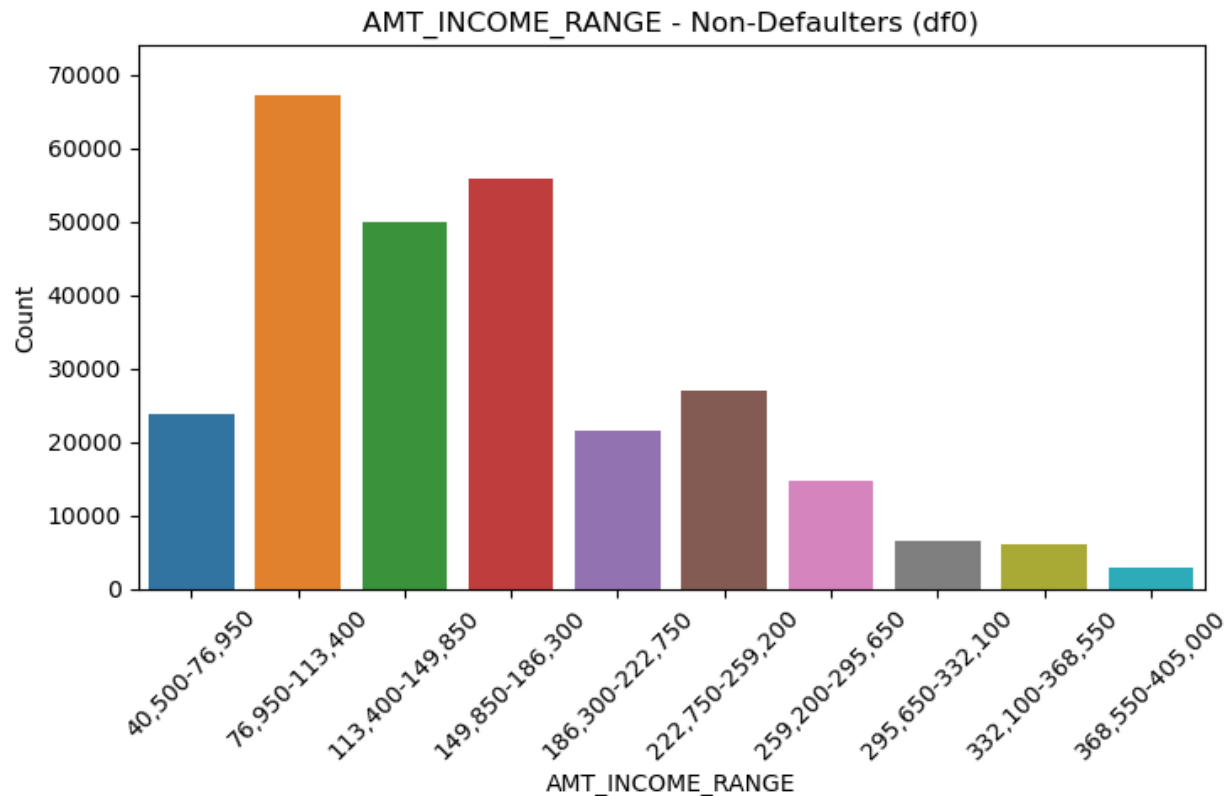


# Bivariate Analysis

21

To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- Both groups(defaulters & non-defaulters) have highest **INCOME** in the range of 76 to 113K



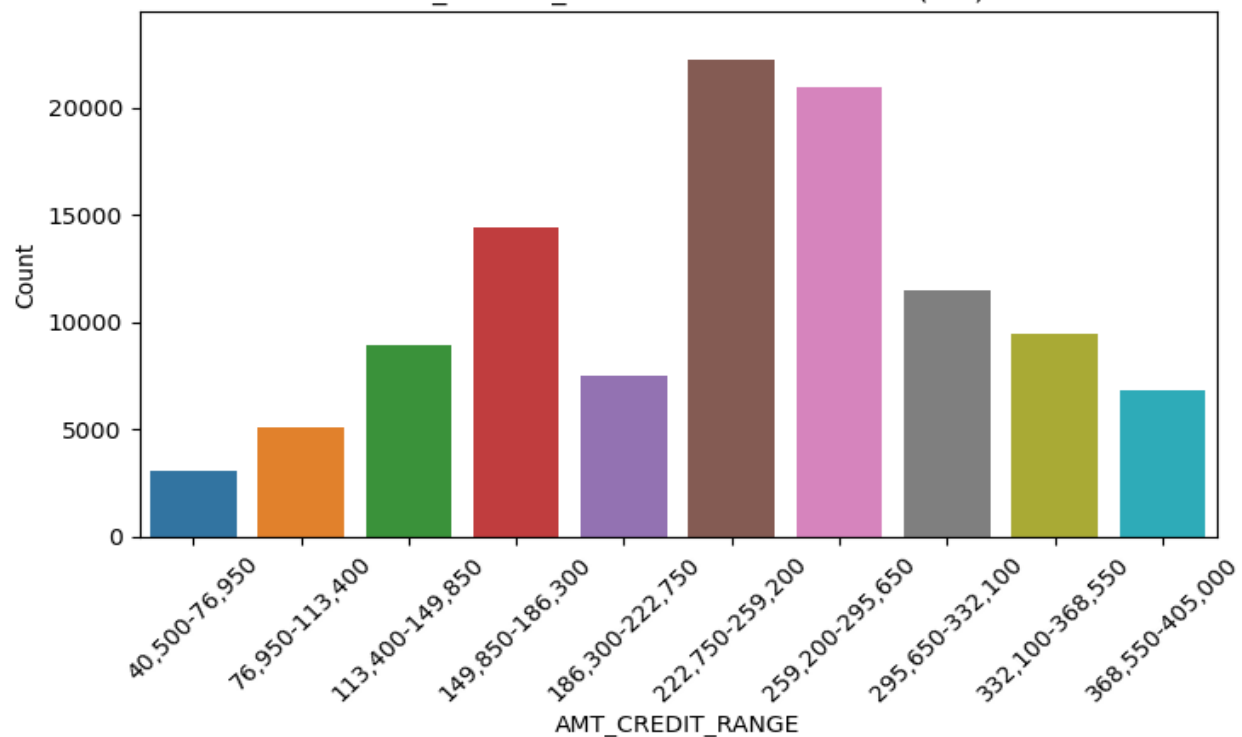
# Bivariate Analysis

22

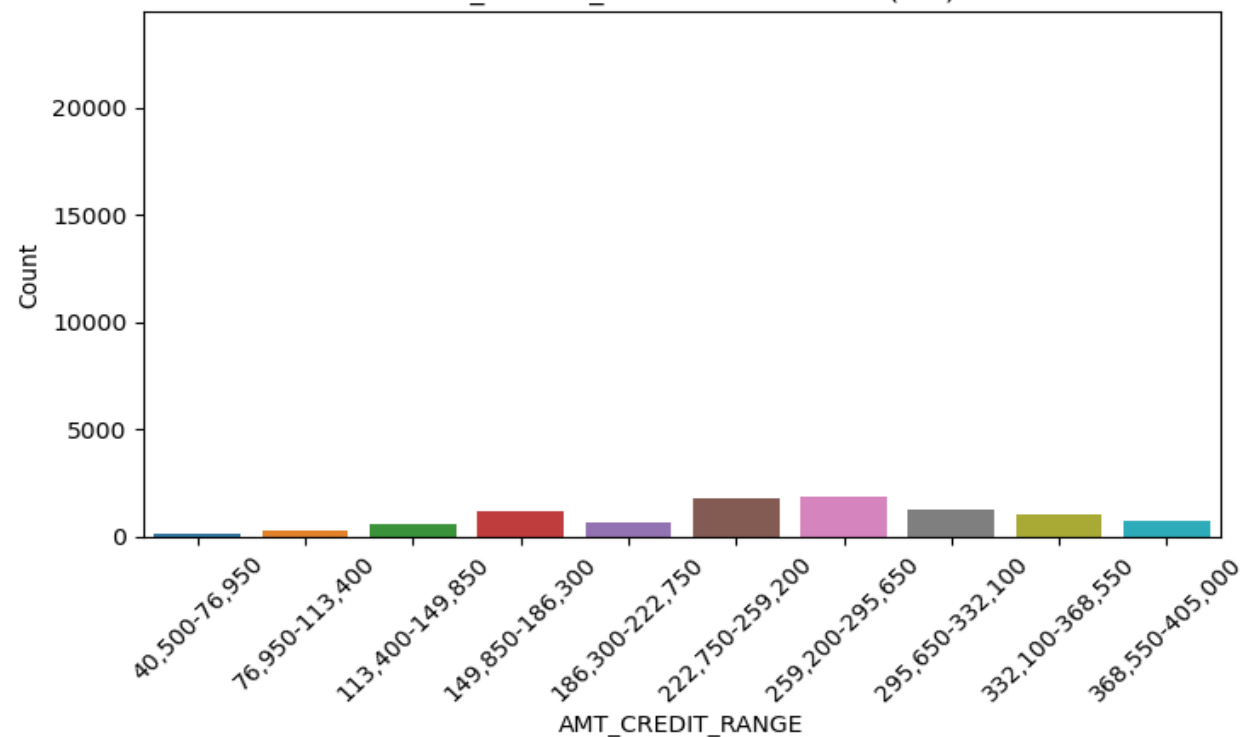
To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- More number of non-defaulter clients have taken higher **CREDIT\_AMT** then defaulters.

AMT\_CREDIT\_RANGE - Non-Defaulters (df0)



AMT\_CREDIT\_RANGE - Defaulters (df1)

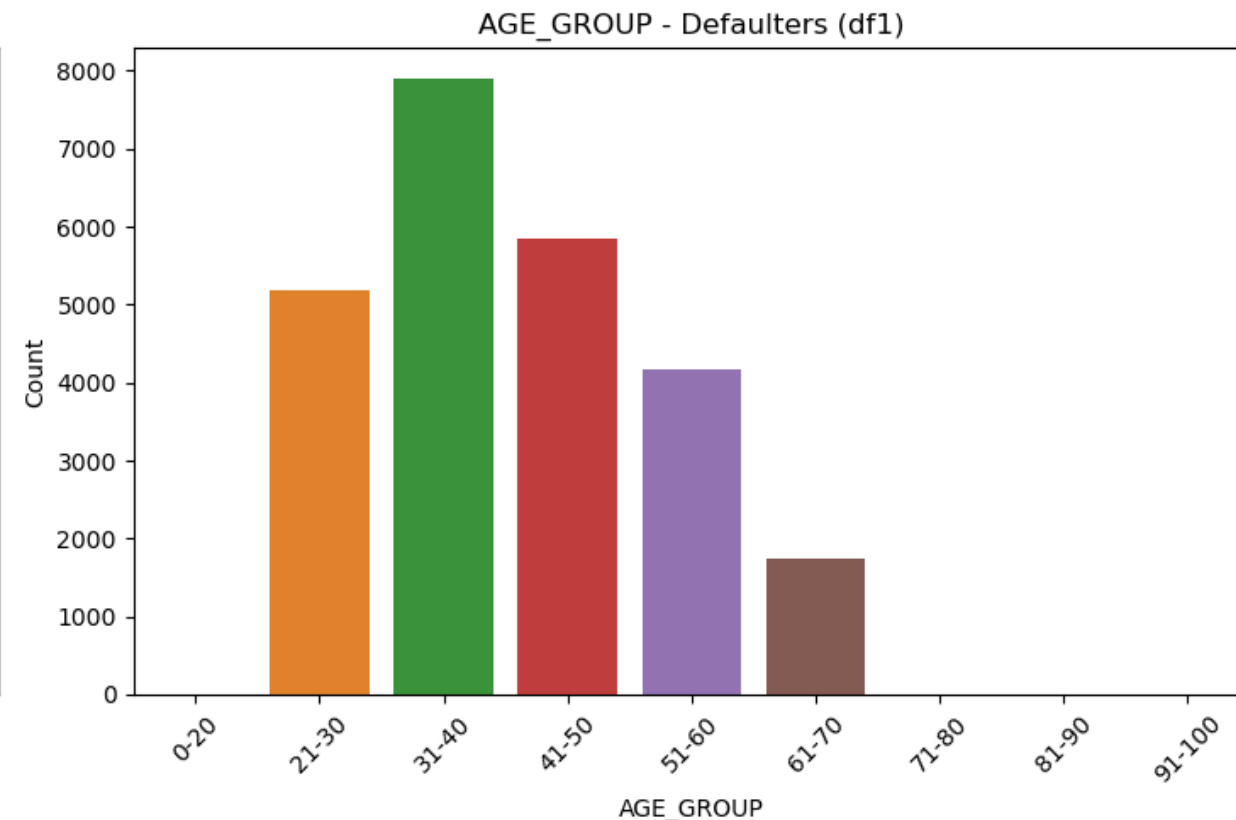
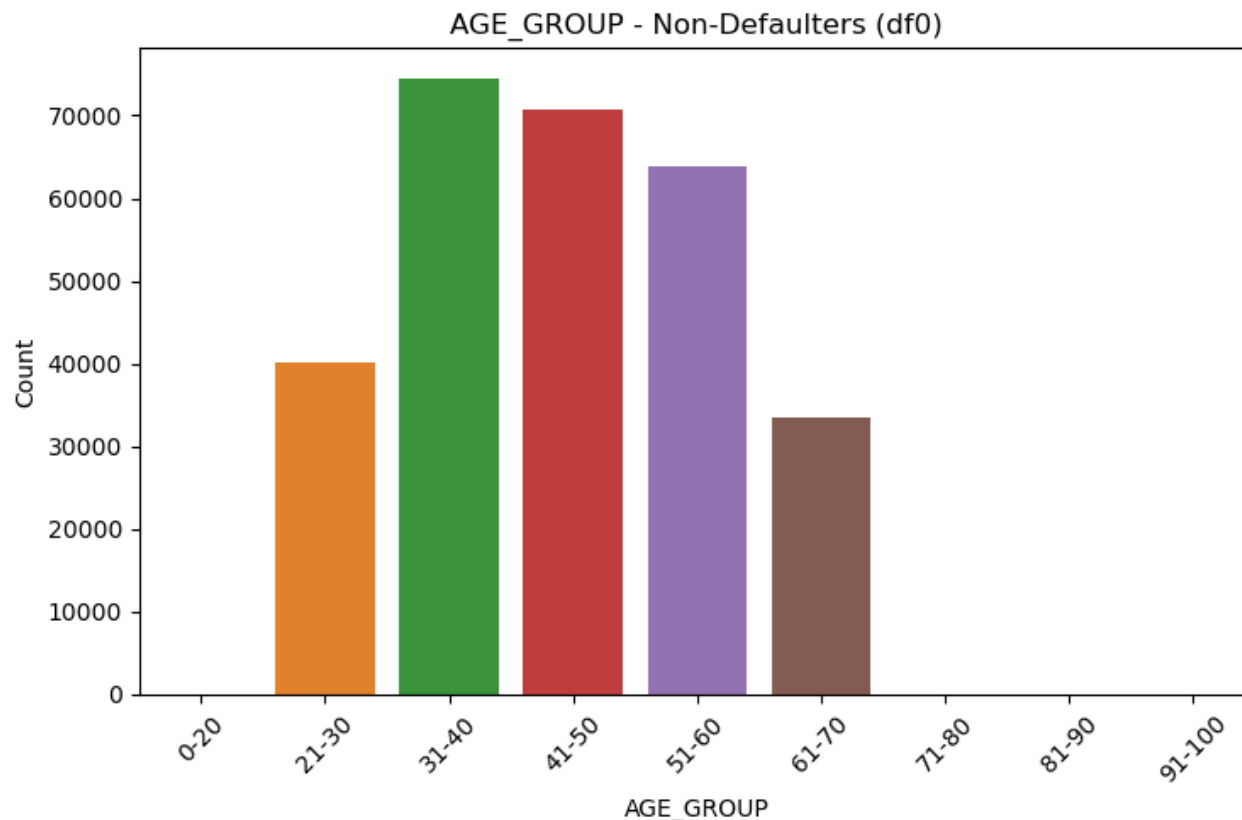


# Bivariate Analysis

23

To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- **31- to 40-year-olds** were higher defaulters then other age groups in “Defaulters(df1)

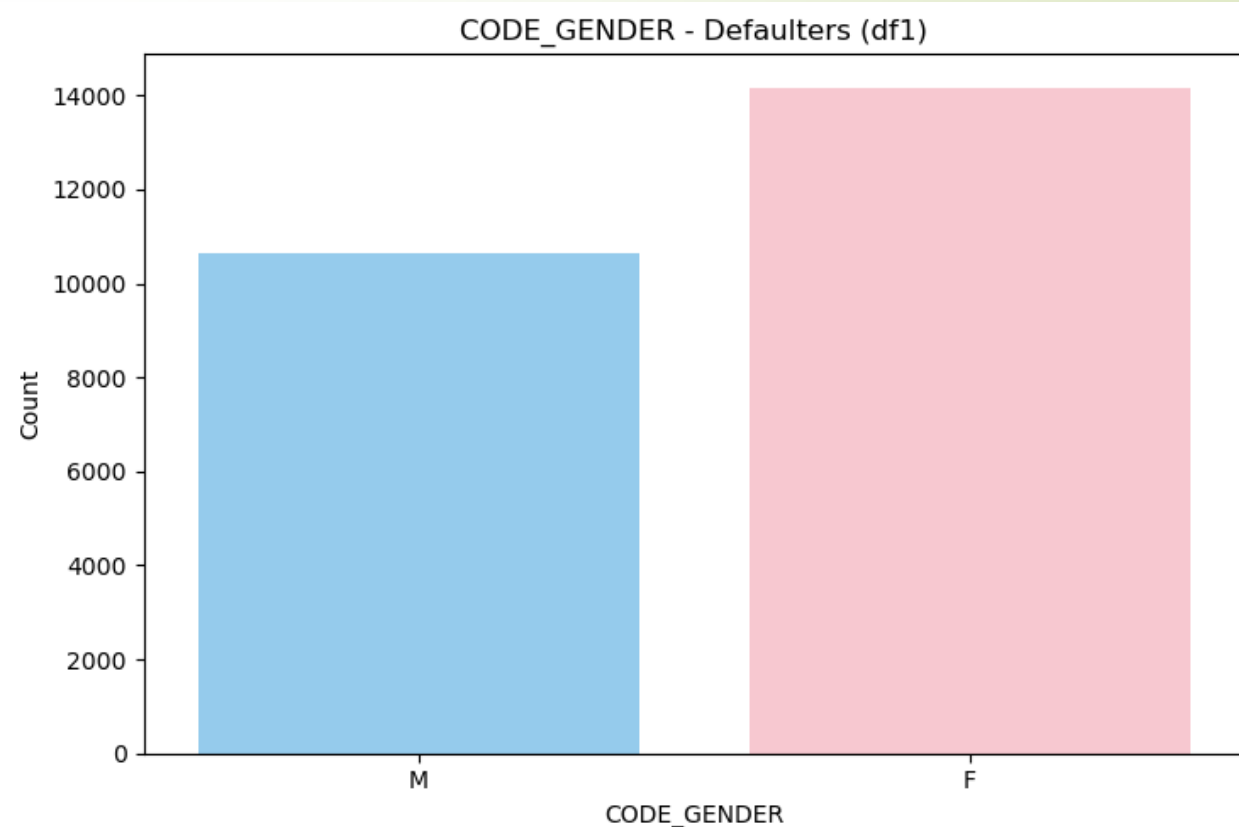
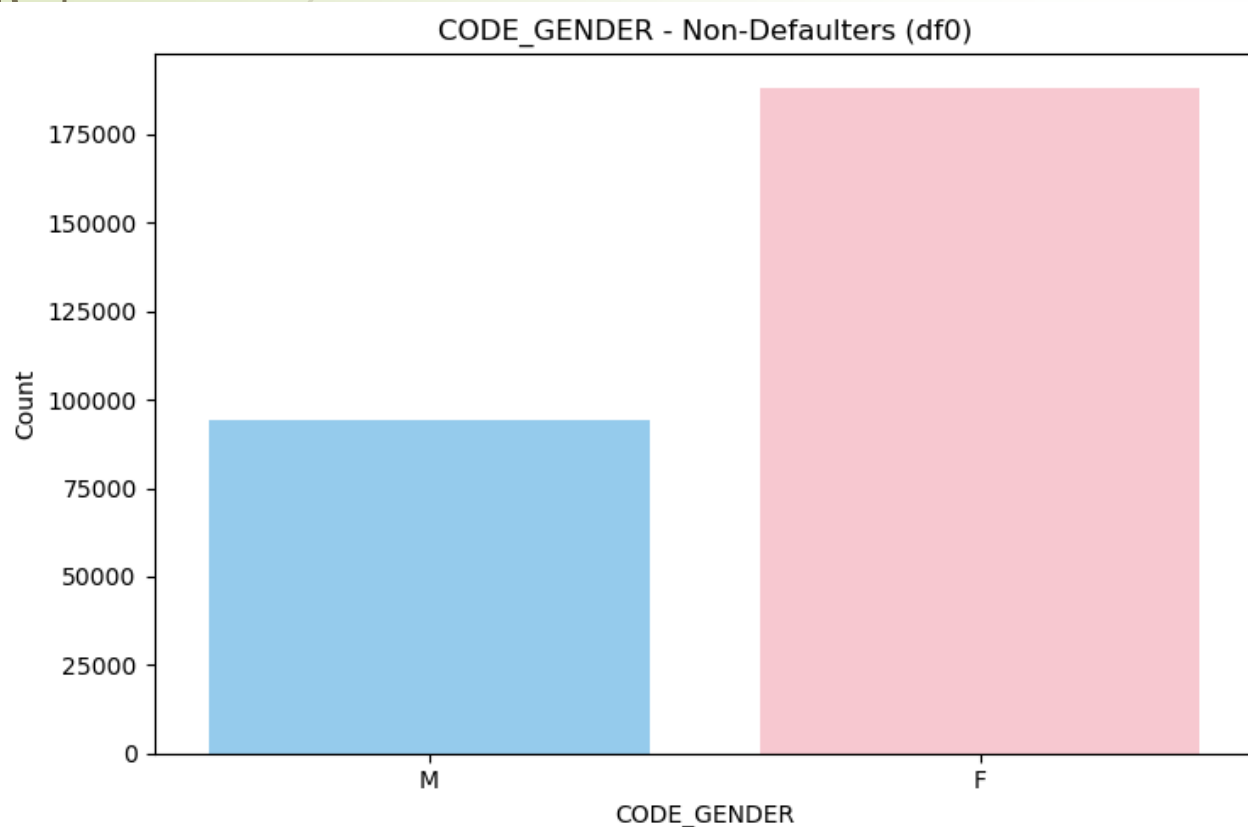


# Bivariate Analysis

24

To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- In both the TARGET groups, **Female** population was highest.





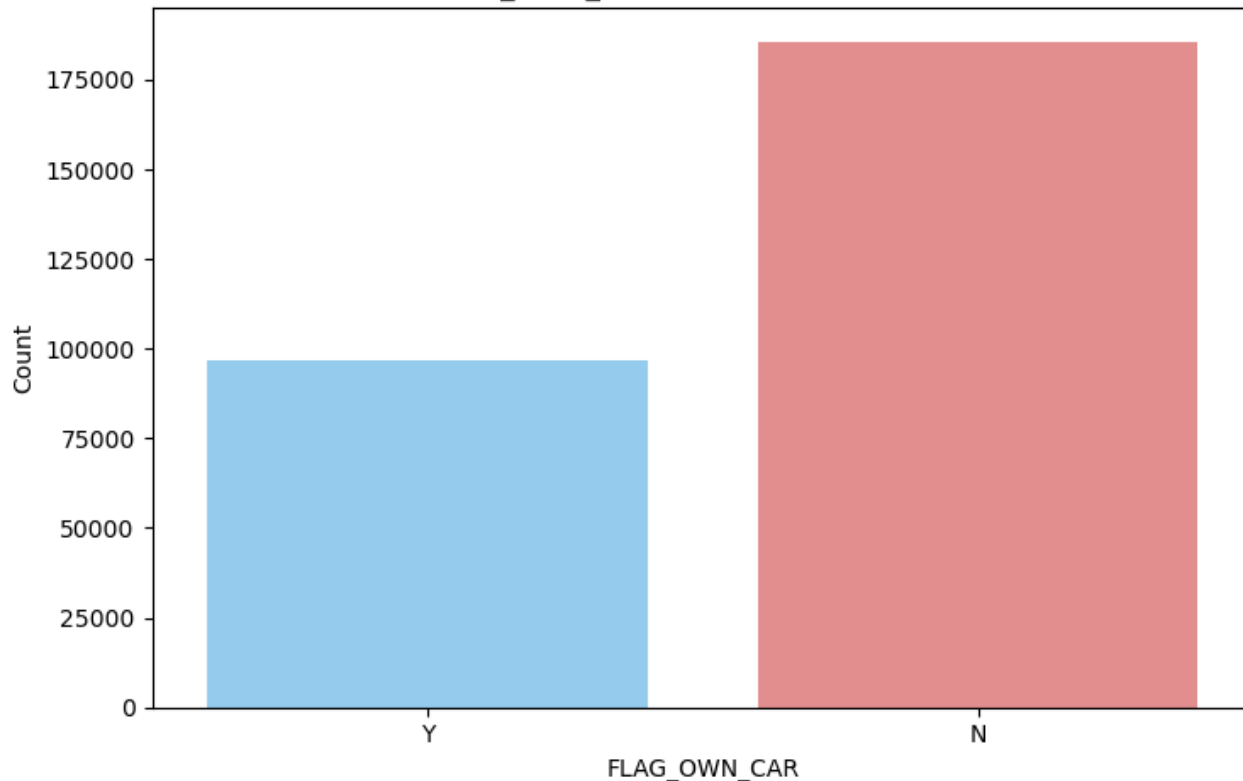
# Bivariate Analysis

25

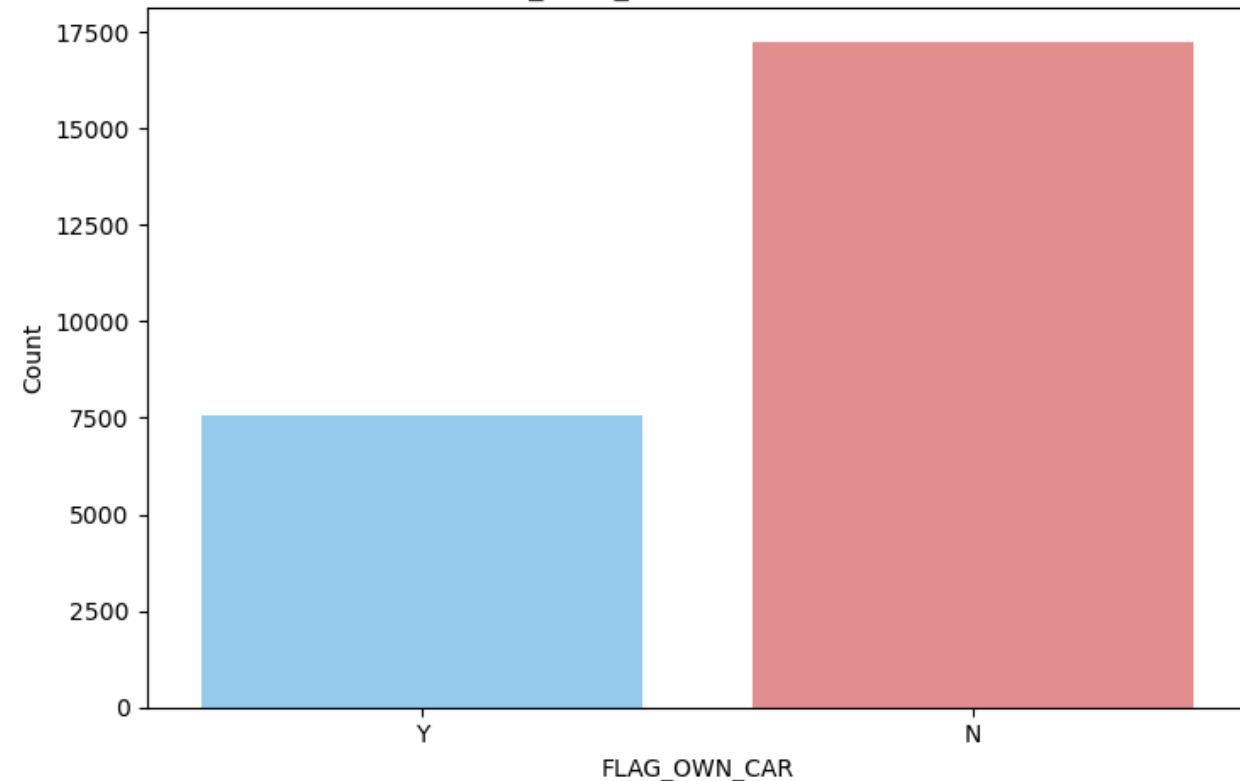
To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- In both the TARGET groups, most of them did not own a **car** while applying for the loan.

FLAG\_OWN\_CAR - Non-Defaulters (df0)



FLAG\_OWN\_CAR - Defaulters (df1)

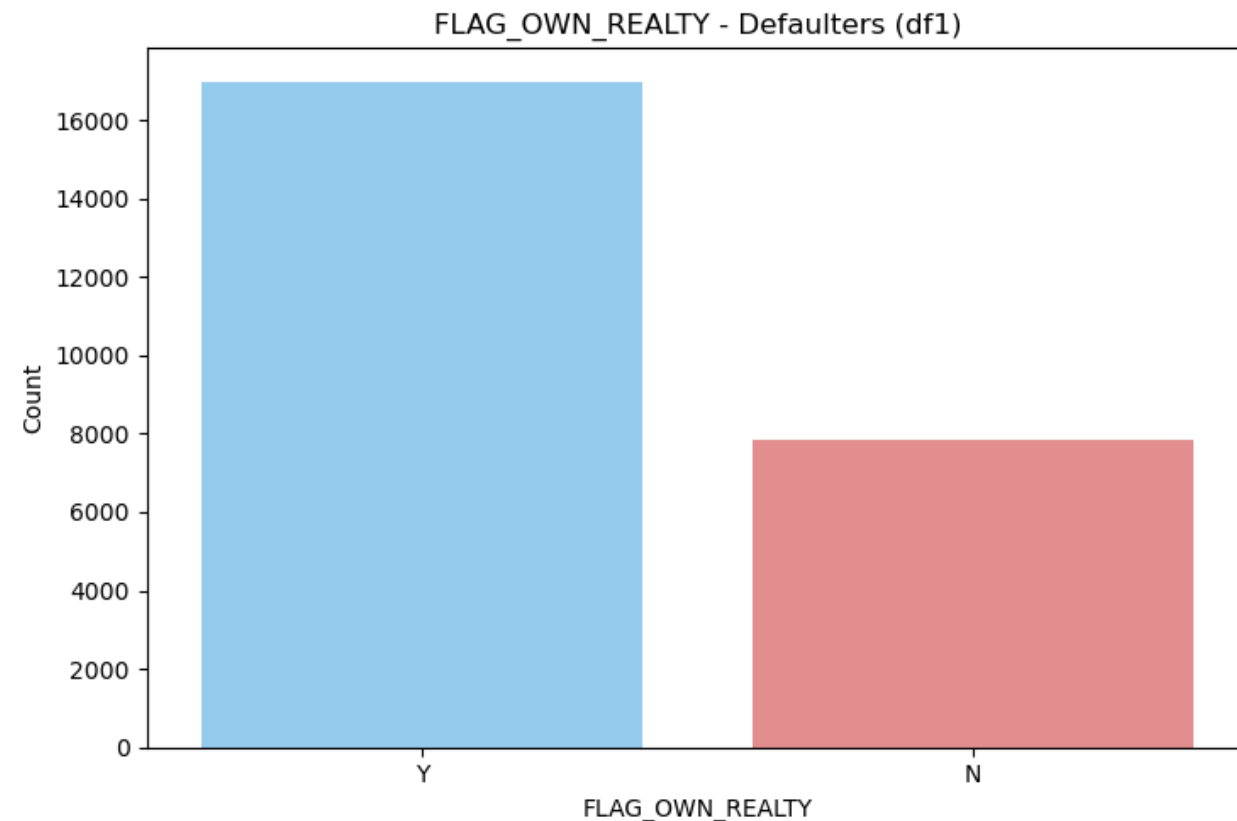
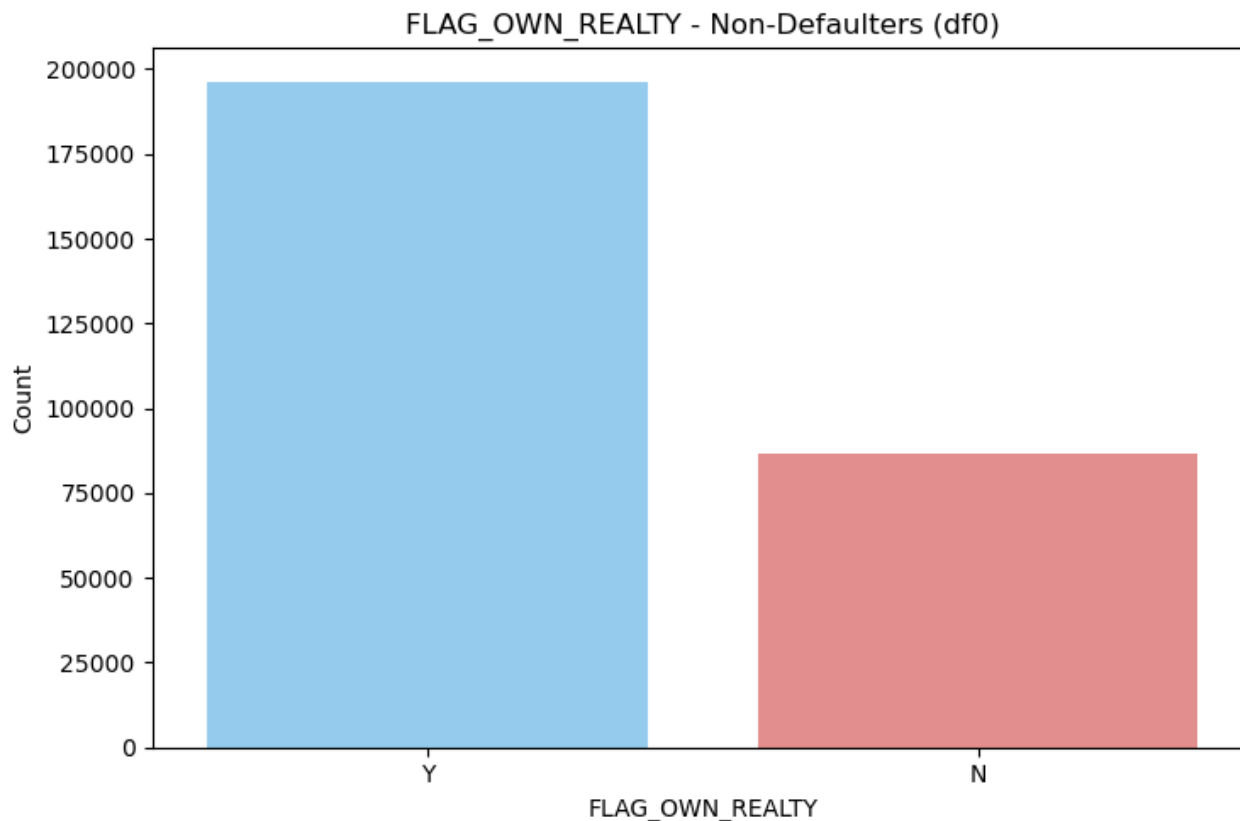


# Bivariate Analysis

26

To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- In both the TARGET groups, most of them owned a **house or flat** while applying for the loan.



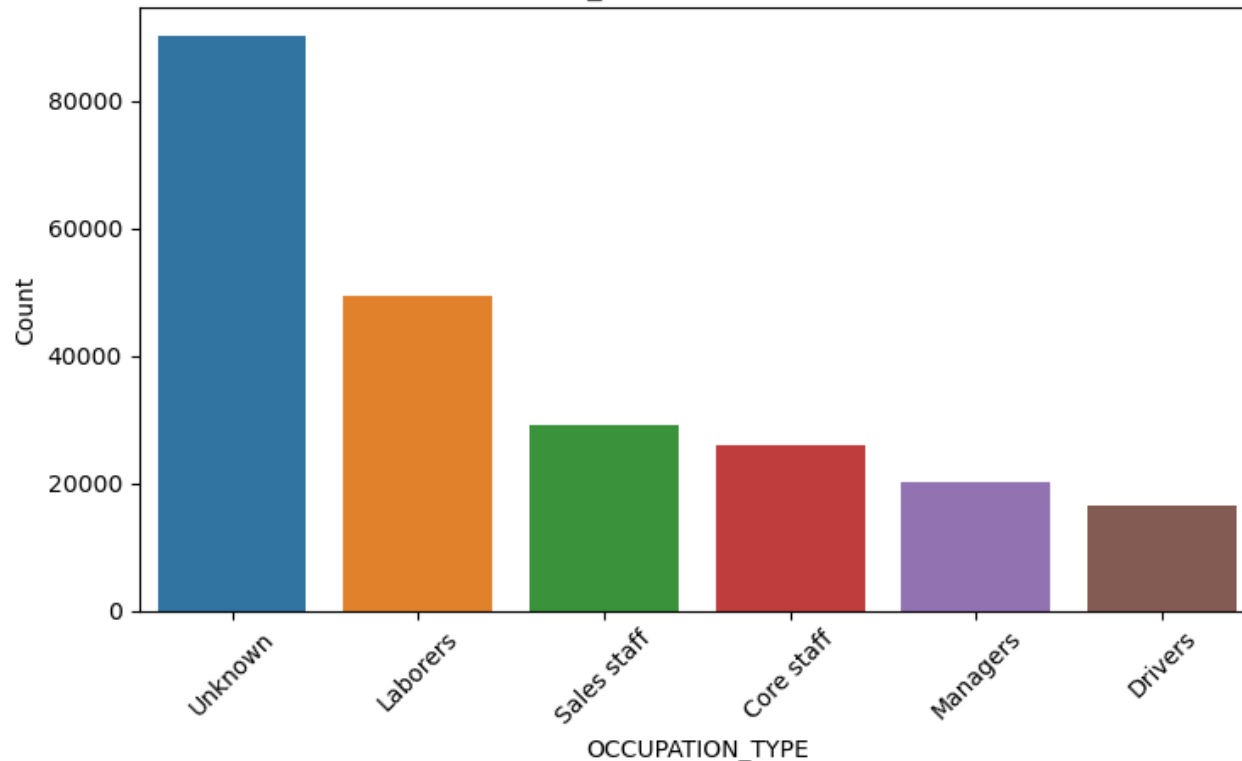
# Bivariate Analysis

27

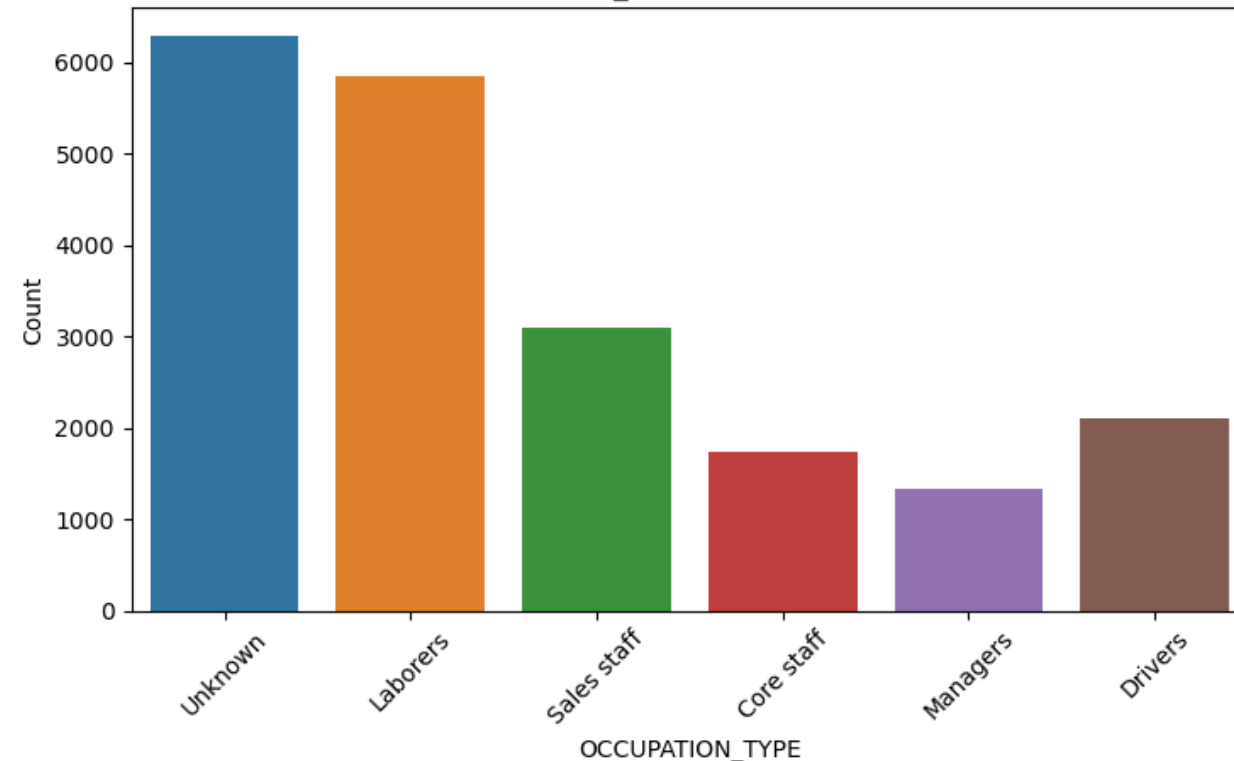
To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- In both the TARGET groups , “**OCCUPATION\_TYPE**” column saw **Laborers** who were the one's who took more loan

OCCUPATION\_TYPE - Non-Defaulters (df0)



OCCUPATION\_TYPE - Defaulters (df1)



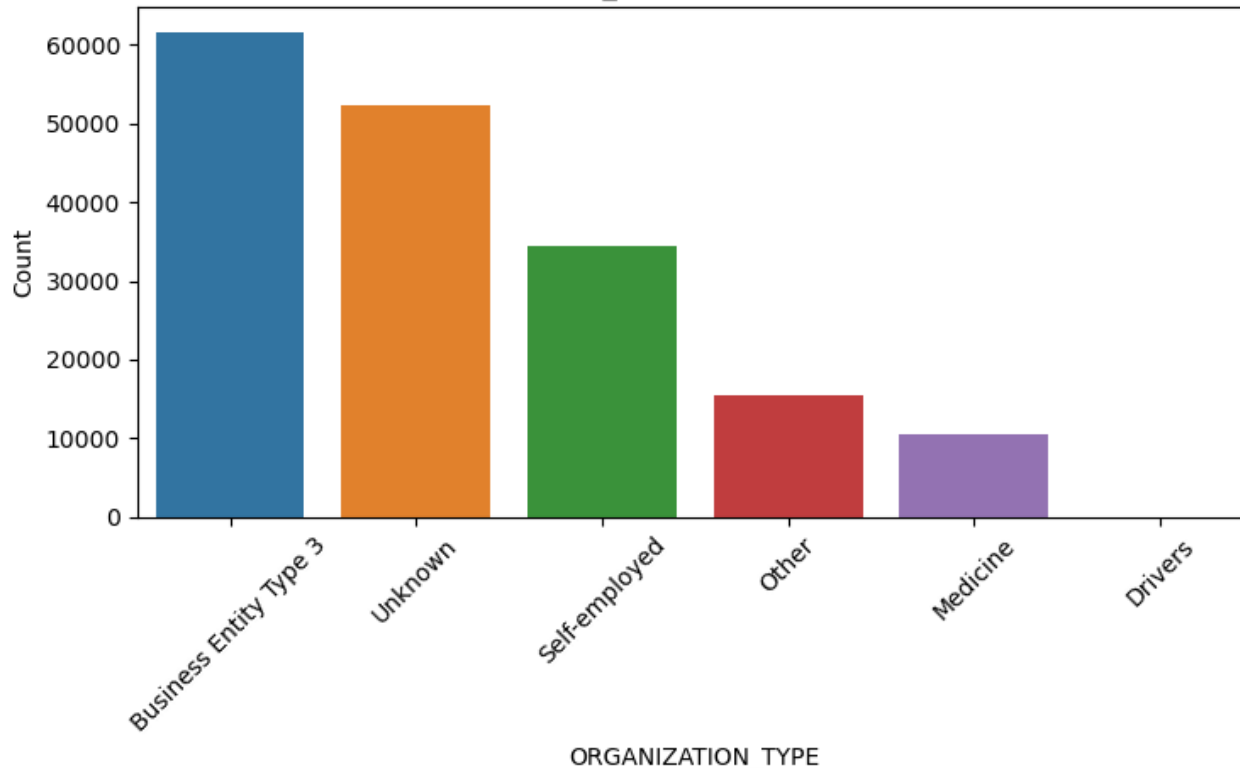
# Bivariate Analysis

28

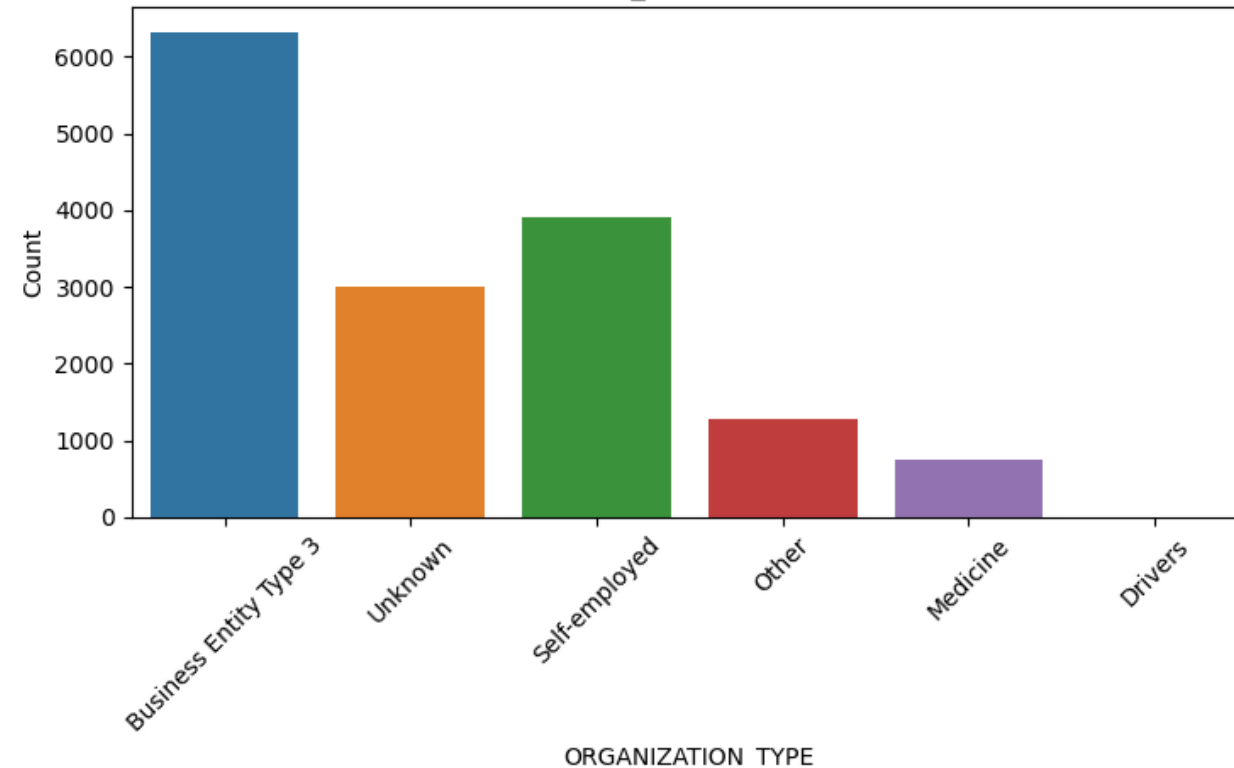
## To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- In both the TARGET groups, “**ORGANIZATION\_TYPE**” column saw “**Business Entity Type 3**” and “**Self-Employed**” working population availed more loans than others. There were more number of rows which did not show which type of work they did, and they were grouped as “**Unknown**”

ORGANIZATION\_TYPE - Non-Defaulters (df0)



ORGANIZATION\_TYPE - Defaulters (df1)

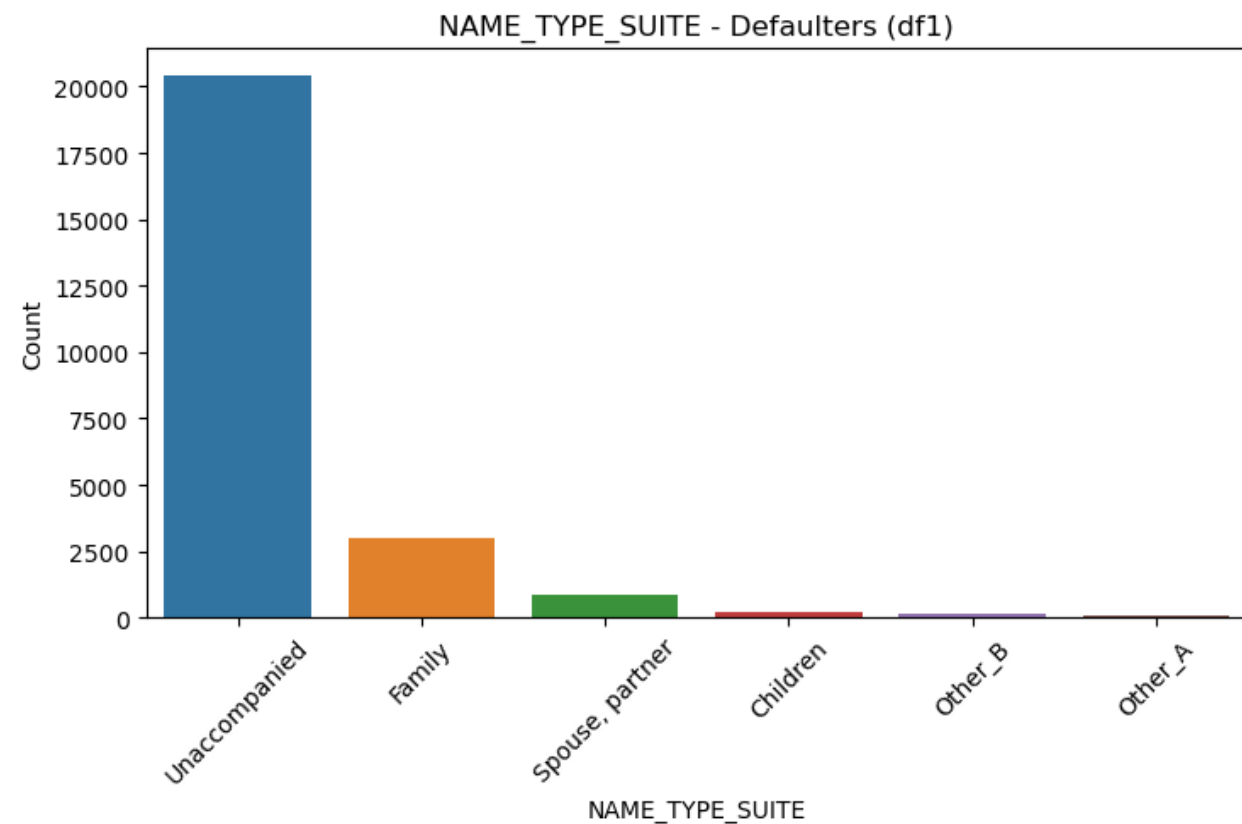
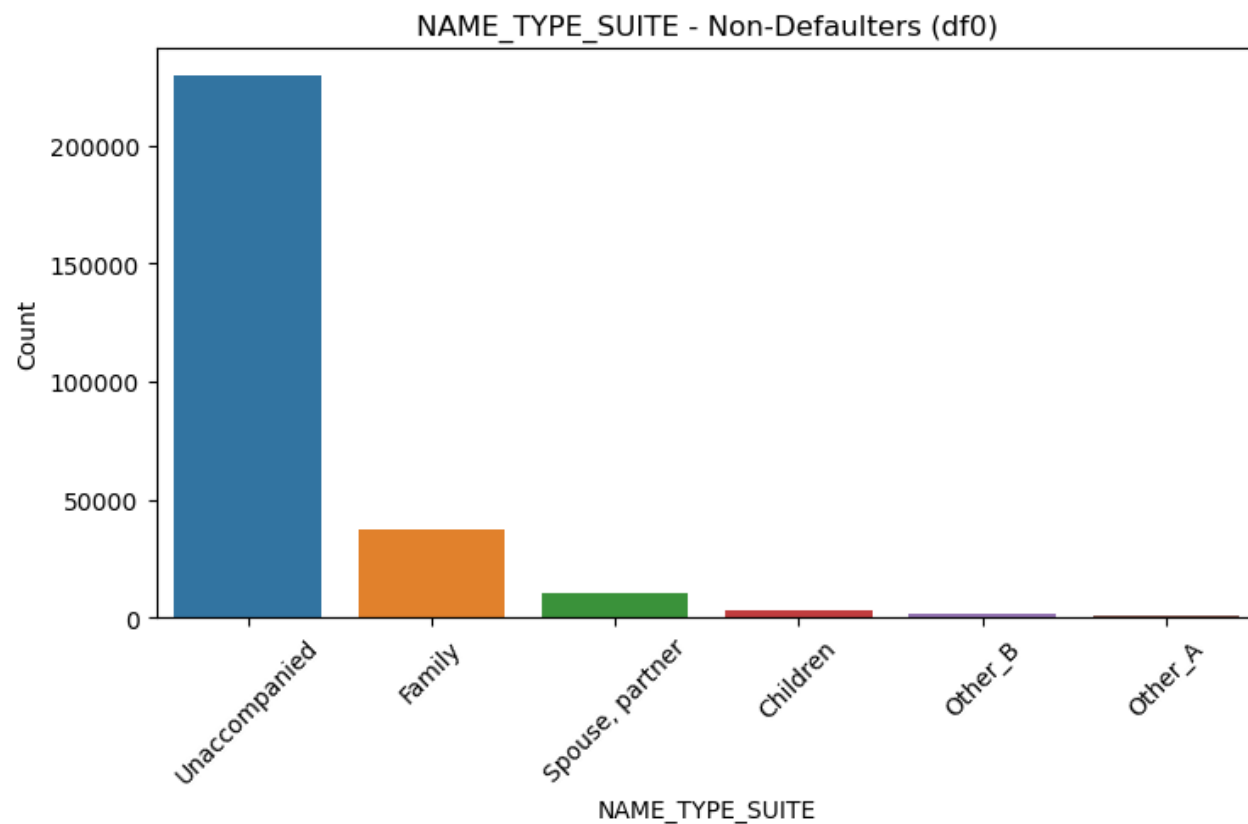


# Bivariate Analysis

29

## To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- In both the TARGET groups, “**NAME\_TYPE\_SUITE**” column saw “**Unaccompanied**” category significantly higher. That is most of the population were not accompanied by anybody.



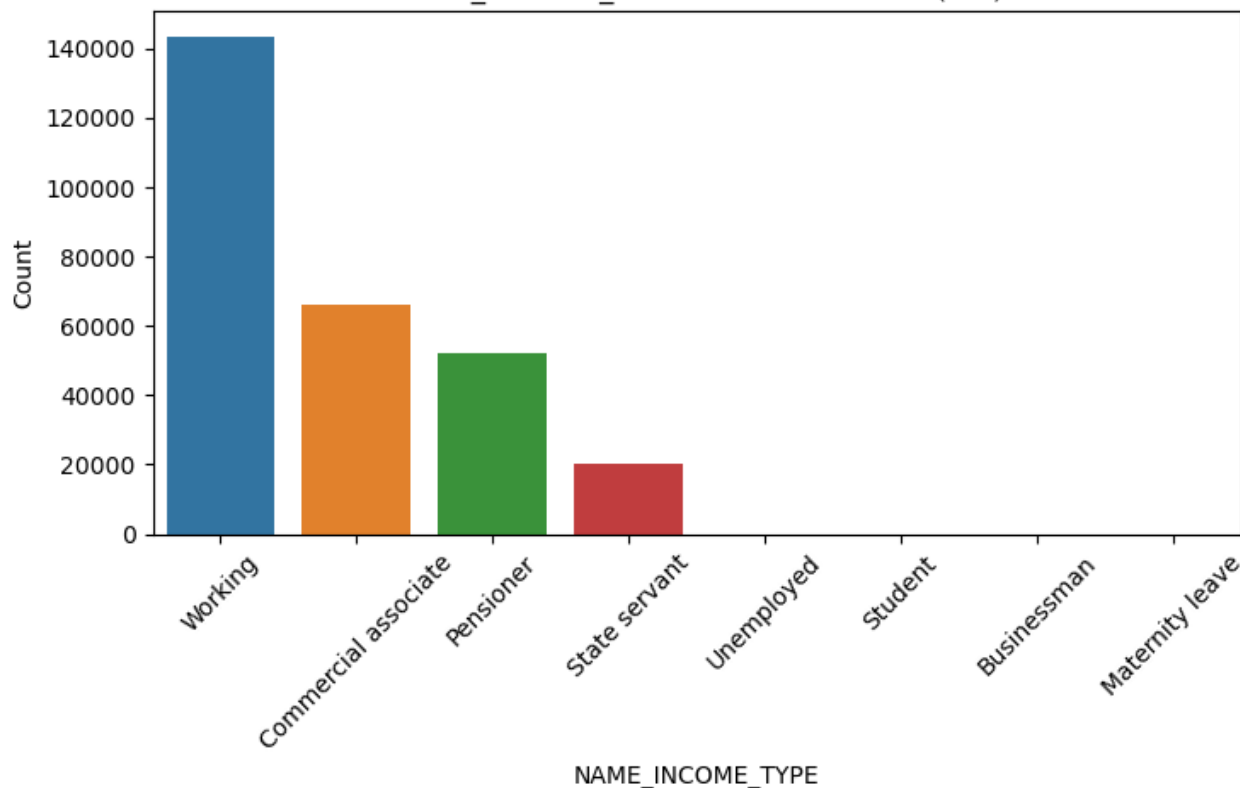
# Bivariate Analysis

30

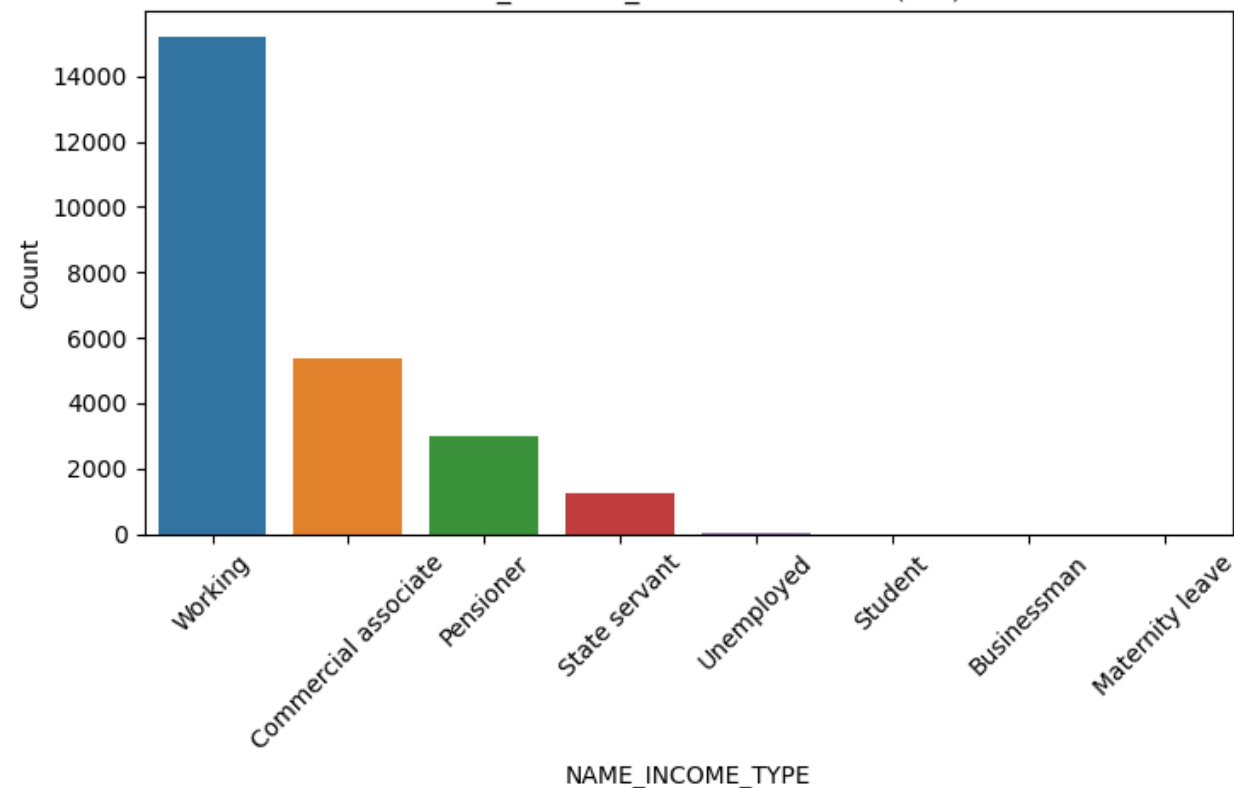
To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- In both the TARGET groups, **working** category saw significantly large number of population followed by **Commercial associate** and **Pensioner** in the **NAME\_INCOME\_TYPE** column.

NAME\_INCOME\_TYPE - Non-Defaulters (df0)



NAME\_INCOME\_TYPE - Defaulters (df1)



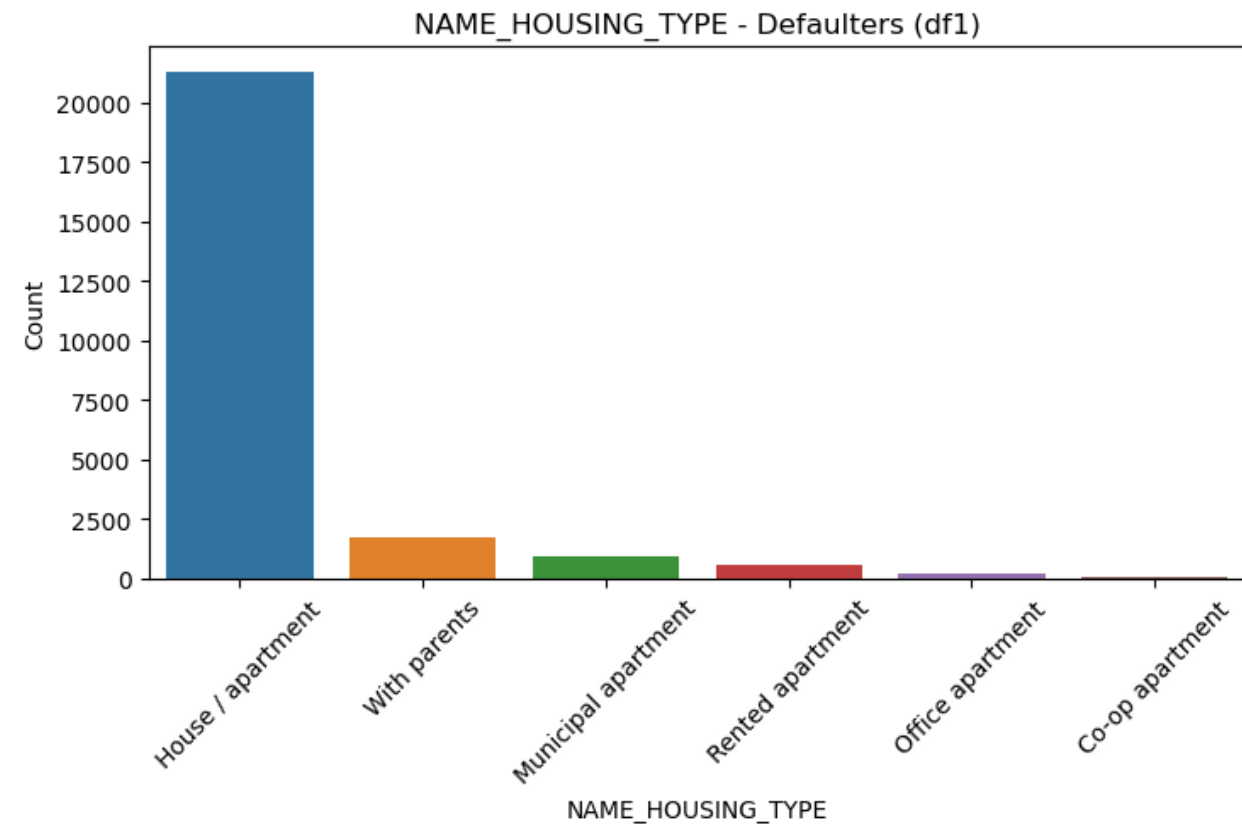
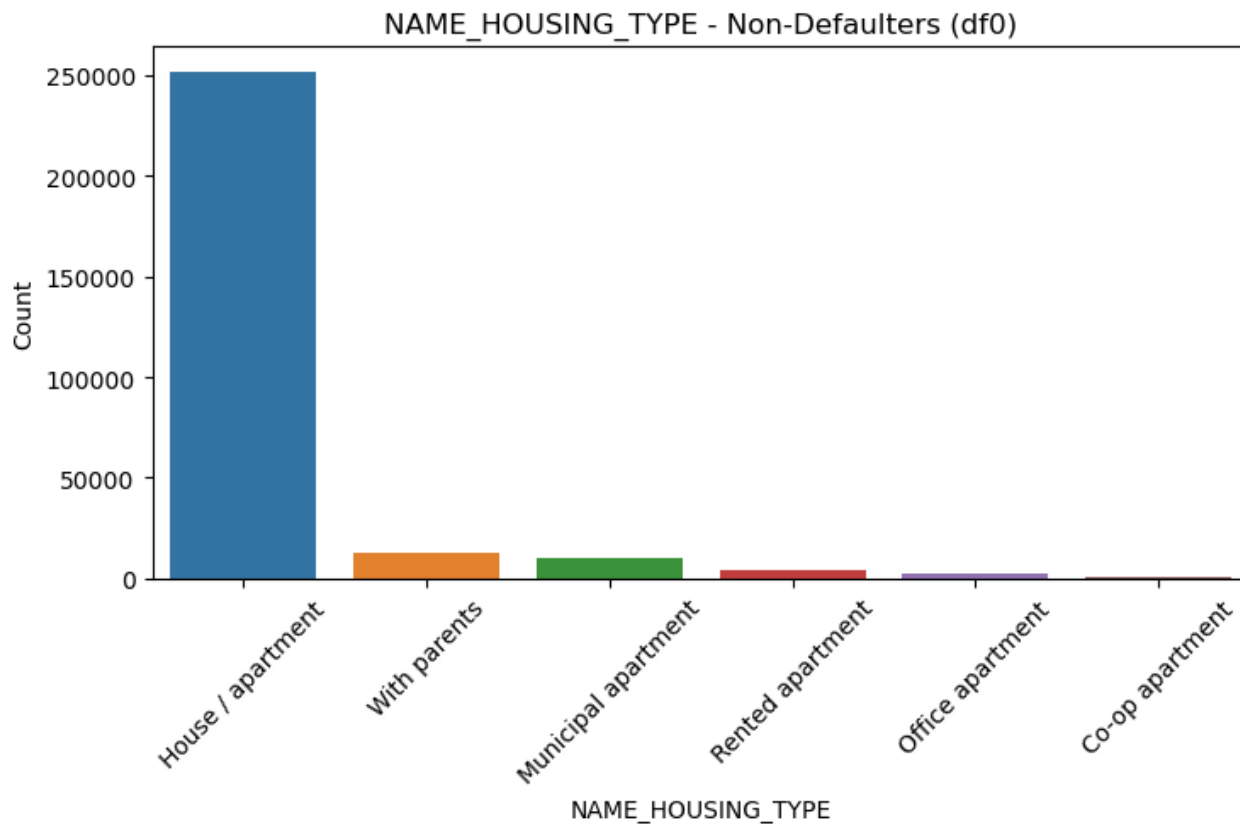


# Bivariate Analysis

31

To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- In both the TARGET groups, population living in **House/apartment** category was significantly large.



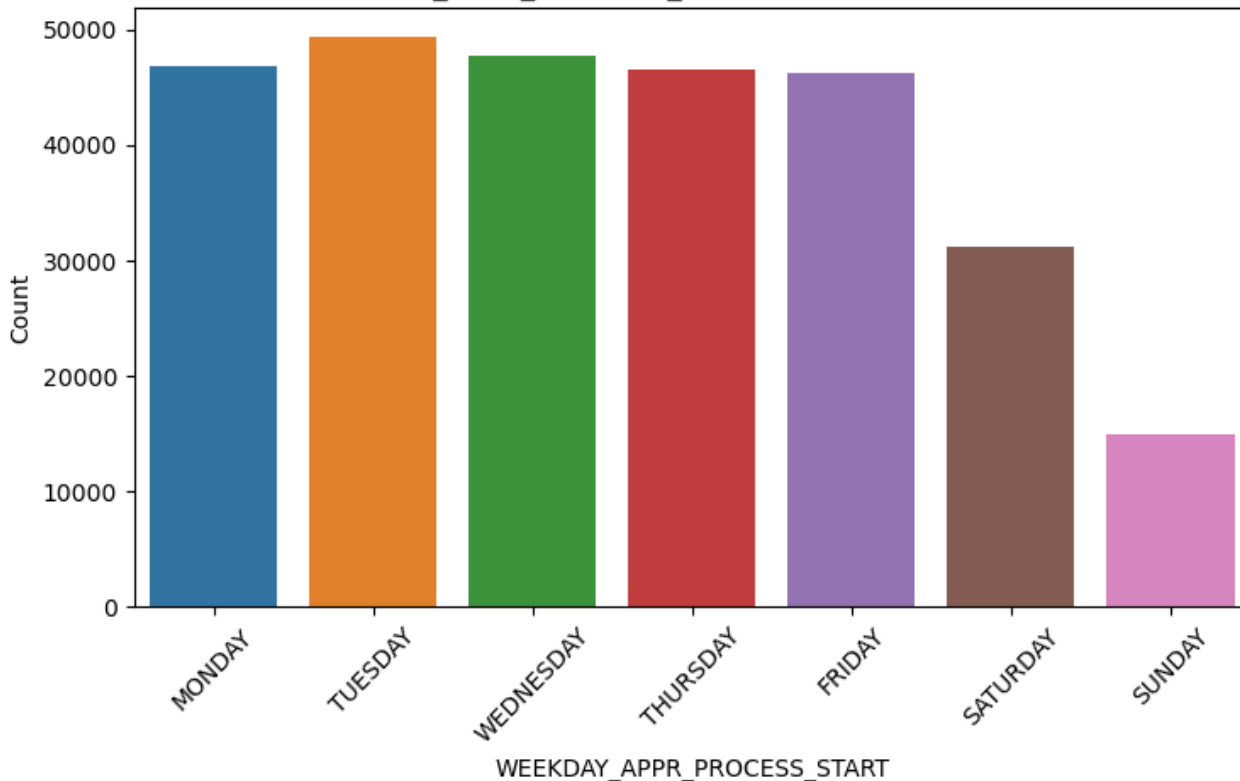
# Bivariate Analysis

32

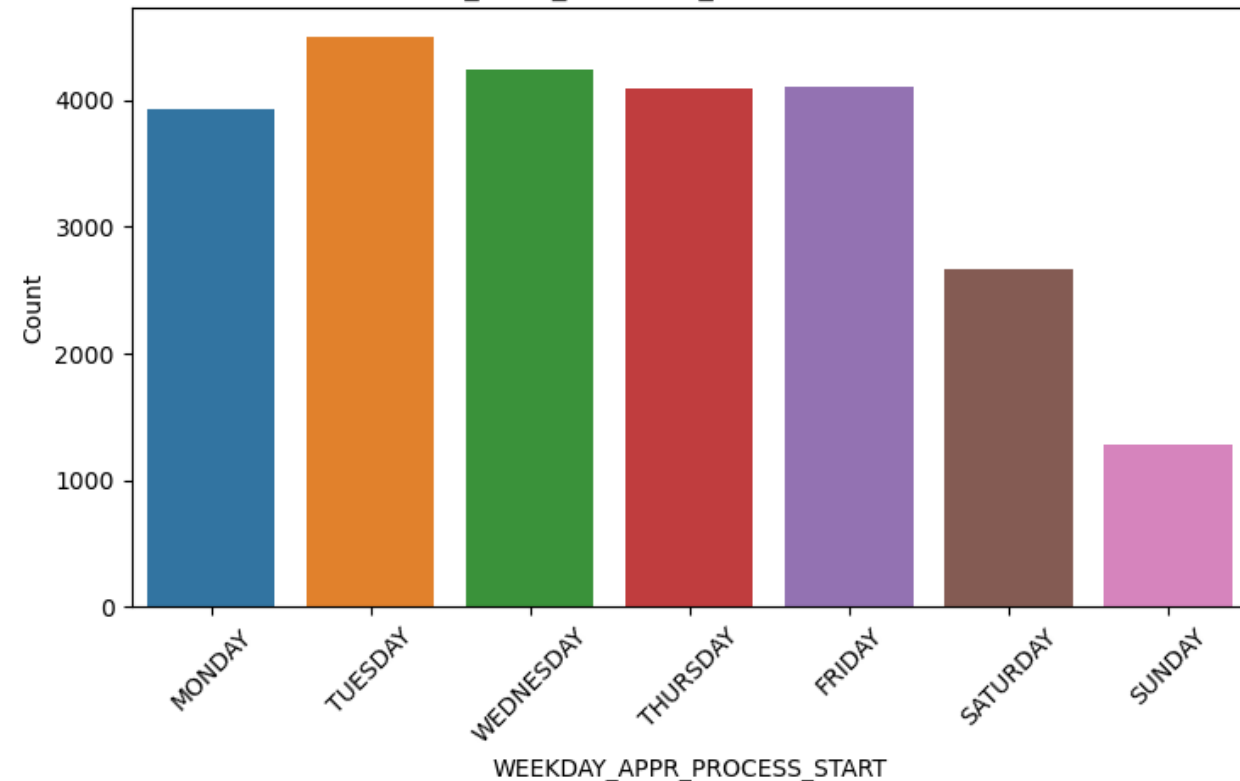
To cross check at columns from TARGET perspective(defaulters and non-defaulters)

- **Tuesdays** saw a increase in “**WEEKDAY\_APPR\_PROCESS\_START**” then other days.

WEEKDAY\_APPR\_PROCESS\_START - Non-Defaulters (df0)



WEEKDAY\_APPR\_PROCESS\_START - Defaulters (df1)



# Multivariate Analysis

33

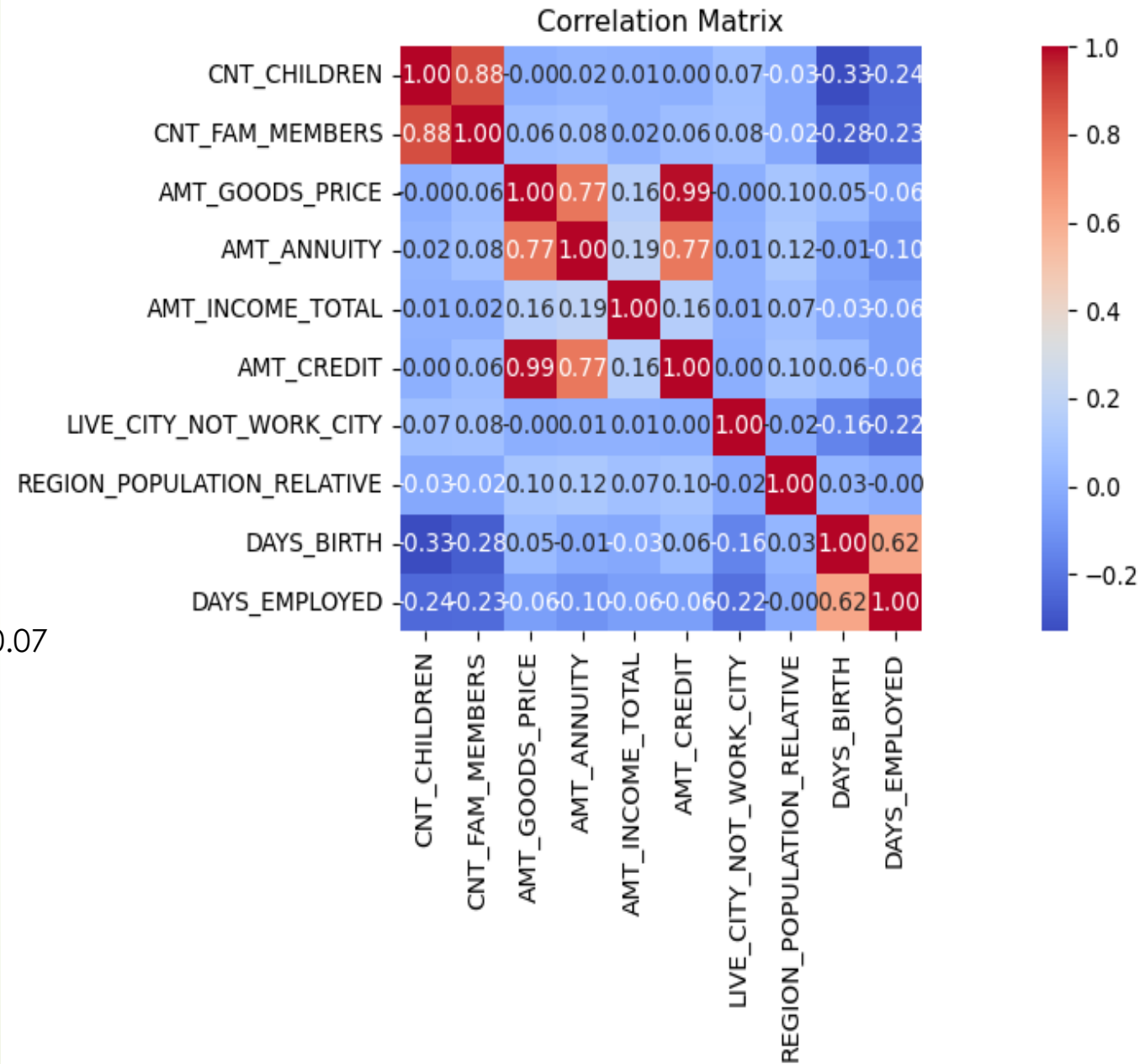
## On the Overall Application Data

There is high correlation between

- $\text{AMT\_CREDIT} \text{ \& } \text{AMT\_GOODS\_PRICE} = 0.99$
- $\text{AMT\_GOODS\_PRICE} \text{ \& } \text{AMT\_ANNUITY} = 0.77$
- $\text{DAYS\_BIRTH} \text{ \& } \text{DAYS\_EMPLOYED} = 0.62$
- $\text{CNT\_CHILDREN} \text{ \& } \text{CNT\_FAM\_MEMBERS} = 0.88$

There is low Correlation between

- $\text{LIVE\_CITY\_NOT\_WORK\_CITY} \text{ \& } \text{AMT\_INCOME} = 0.008$
- $\text{AMT\_CREDIT} \text{ \& } \text{CNT\_FAM\_MEMBERS} = 0.06$
- $\text{REGION\_POPULATION\_RELATIVE} \text{ \& } \text{AMT\_INNCOME} = 0.07$

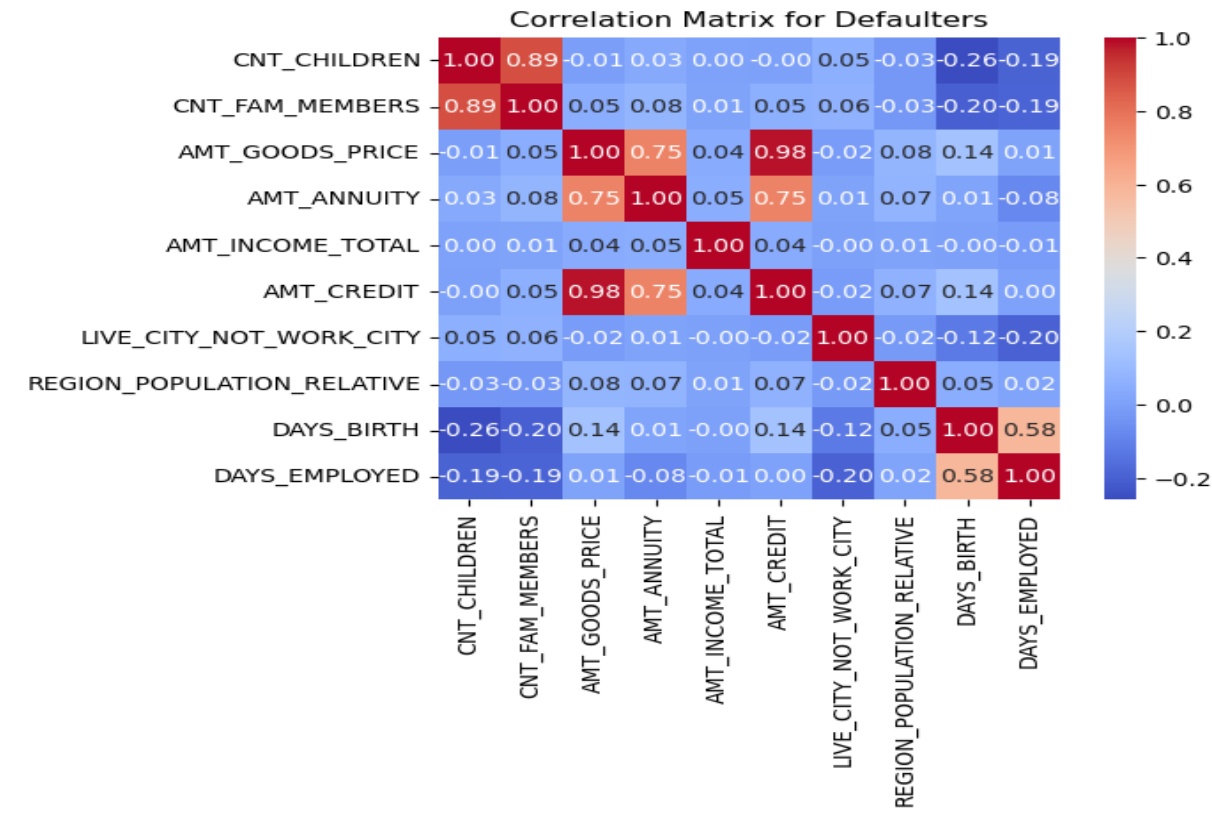
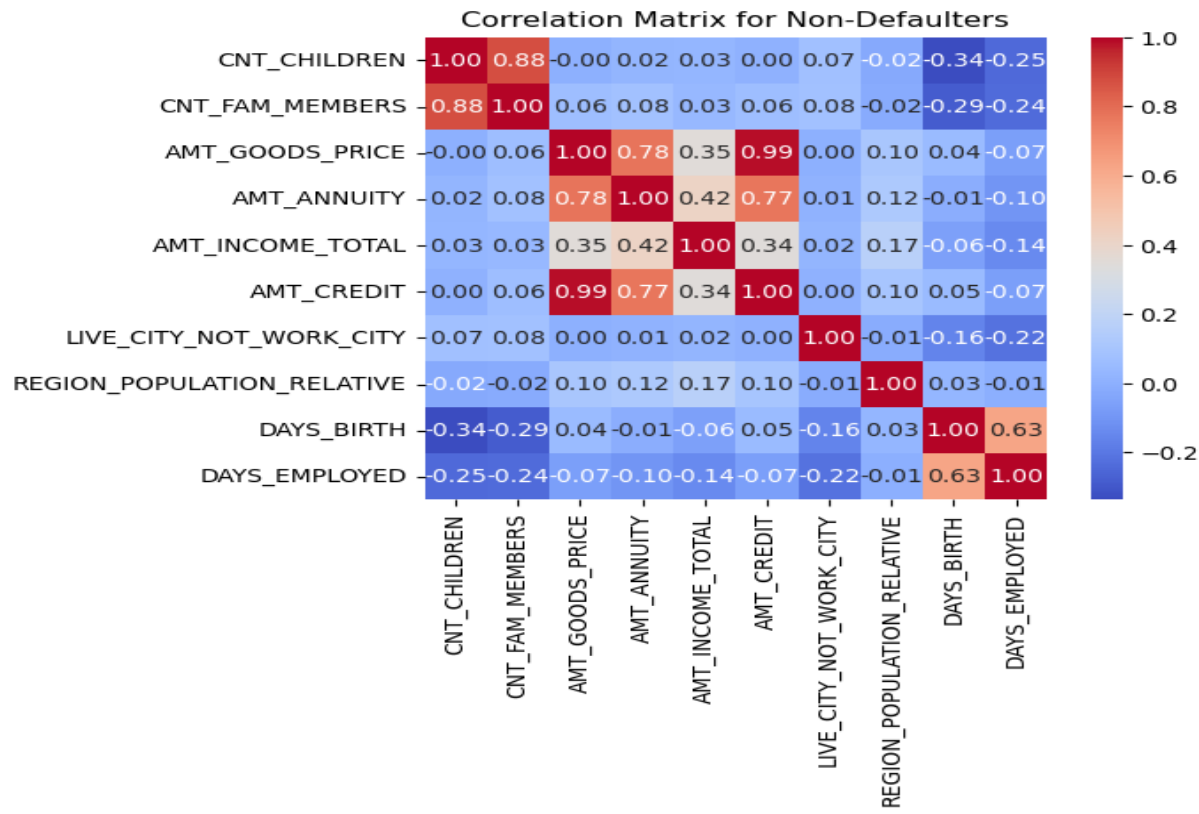


# Multivariate Analysis

34

To cross check using TARGET columns for defaulters and non-defaulters

Among , **non-defaulters** we see **positive good correlation** between **AMT\_INCOME\_TOTAL** & **AMT\_GOODS\_PRICE** , **AMT\_INCOME\_TOTAL** & **AMT\_ANNUITY**, **AMT\_CREDIT** & **AMOUNT\_INCOME\_TOTAL** which is **not seen** in defaulters.



# Multivariate Analysis

35

Correlation matrix for Non-Defaulters:

	CNT_CHILDREN	CNT_FAM_MEMBERS	AMT_GOODS_PRICE	\
CNT_CHILDREN	1.000000	0.878571	-0.000559	
CNT_FAM_MEMBERS	0.878571	1.000000	0.062763	
AMT_GOODS_PRICE	-0.000559	0.062763	1.000000	
AMT_ANNUITY	0.020909	0.075789	0.776421	
AMT_INCOME_TOTAL	0.027397	0.034254	0.349426	
AMT_CREDIT	0.003081	0.064536	0.987022	
LIVE_CITY_NOT_WORK_CITY	0.070988	0.078843	0.001285	
REGION_POPULATION_RELATIVE	-0.024363	-0.023425	0.103826	
DAYS_BIRTH	-0.336966	-0.285823	0.044650	
DAYS_EMPLOYED	-0.245174	-0.238300	-0.068527	

	AMT_ANNUITY	AMT_INCOME_TOTAL	AMT_CREDIT	\
CNT_CHILDREN	0.020909	0.027397	0.003081	
CNT_FAM_MEMBERS	0.075789	0.034254	0.064536	
AMT_GOODS_PRICE	0.776421	0.349426	0.987022	
AMT_ANNUITY	1.000000	0.418948	0.771297	
AMT_INCOME_TOTAL	0.418948	1.000000	0.342799	
AMT_CREDIT	0.771297	0.342799	1.000000	
LIVE_CITY_NOT_WORK_CITY	0.010577	0.020684	0.002506	
REGION_POPULATION_RELATIVE	0.120977	0.167851	0.100604	
DAYS_BIRTH	-0.012260	-0.062609	0.047378	
DAYS_EMPLOYED	-0.104975	-0.140392	-0.070104	

Correlation matrix for Defaulters:

	CNT_CHILDREN	CNT_FAM_MEMBERS	AMT_GOODS_PRICE	\
CNT_CHILDREN	1.000000	0.885484	-0.008111	
CNT_FAM_MEMBERS	0.885484	1.000000	0.047367	
AMT_GOODS_PRICE	-0.008111	0.047367	1.000000	
AMT_ANNUITY	0.031257	0.075711	0.752295	
AMT_INCOME_TOTAL	0.004796	0.006654	0.037591	
AMT_CREDIT	-0.001675	0.051224	0.982783	
LIVE_CITY_NOT_WORK_CITY	0.053515	0.061316	-0.016703	
REGION_POPULATION_RELATIVE	-0.031975	-0.030163	0.076053	
DAYS_BIRTH	-0.259109	-0.203267	0.135738	
DAYS_EMPLOYED	-0.192864	-0.186515	0.006648	

	AMT_ANNUITY	AMT_INCOME_TOTAL	AMT_CREDIT	\
CNT_CHILDREN	0.031257	0.004796	-0.001675	
CNT_FAM_MEMBERS	0.075711	0.006654	0.051224	
AMT_GOODS_PRICE	0.752295	0.037591	0.982783	
AMT_ANNUITY	1.000000	0.046421	0.752195	
AMT_INCOME_TOTAL	0.046421	1.000000	0.038131	
AMT_CREDIT	0.752195	0.038131	1.000000	
LIVE_CITY_NOT_WORK_CITY	0.009902	-0.001353	-0.016509	
REGION_POPULATION_RELATIVE	0.071690	0.009135	0.069161	
DAYS_BIRTH	0.014303	-0.003096	0.135316	
DAYS_EMPLOYED	-0.081207	-0.014977	0.001930	



# Multivariate Analysis

36

## Non-Defaulters

	LIVE_CITY_NOT_WORK_CITY \
CNT_CHILDREN	0.070988
CNT_FAM_MEMBERS	0.078843
AMT_GOODS_PRICE	0.001285
AMT_ANNUITY	0.010577
AMT_INCOME_TOTAL	0.020684
AMT_CREDIT	0.002506
LIVE_CITY_NOT_WORK_CITY	1.000000
REGION_POPULATION_RELATIVE	-0.013502
DAYS_BIRTH	-0.160072
DAYS_EMPLOYED	-0.221387

	REGION_POPULATION_RELATIVE	DAYS_BIRTH \
CNT_CHILDREN	-0.024363	-0.336966
CNT_FAM_MEMBERS	-0.023425	-0.285823
AMT_GOODS_PRICE	0.103826	0.044650
AMT_ANNUITY	0.120977	-0.012260
AMT_INCOME_TOTAL	0.167851	-0.062609
AMT_CREDIT	0.100604	0.047378
LIVE_CITY_NOT_WORK_CITY	-0.013502	-0.160072
REGION_POPULATION_RELATIVE	1.000000	0.025244
DAYS_BIRTH	0.025244	1.000000
DAYS_EMPLOYED	-0.007198	0.626114

## Defaulters

	LIVE_CITY_NOT_WORK_CITY \
CNT_CHILDREN	0.053515
CNT_FAM_MEMBERS	0.061316
AMT_GOODS_PRICE	-0.016703
AMT_ANNUITY	0.009902
AMT_INCOME_TOTAL	-0.001353
AMT_CREDIT	-0.016509
LIVE_CITY_NOT_WORK_CITY	1.000000
REGION_POPULATION_RELATIVE	-0.020428
DAYS_BIRTH	-0.123623
DAYS_EMPLOYED	-0.198484

	REGION_POPULATION_RELATIVE	DAYS_BIRTH \
CNT_CHILDREN	-0.031975	-0.259109
CNT_FAM_MEMBERS	-0.030163	-0.203267
AMT_GOODS_PRICE	0.076053	0.135738
AMT_ANNUITY	0.071690	0.014303
AMT_INCOME_TOTAL	0.009135	-0.003096
AMT_CREDIT	0.069161	0.135316
LIVE_CITY_NOT_WORK_CITY	-0.020428	-0.123623
REGION_POPULATION_RELATIVE	1.000000	0.048190
DAYS_BIRTH	0.048190	1.000000
DAYS_EMPLOYED	0.015532	0.582185



# Multivariate Analysis

37

## Non-Defaulters

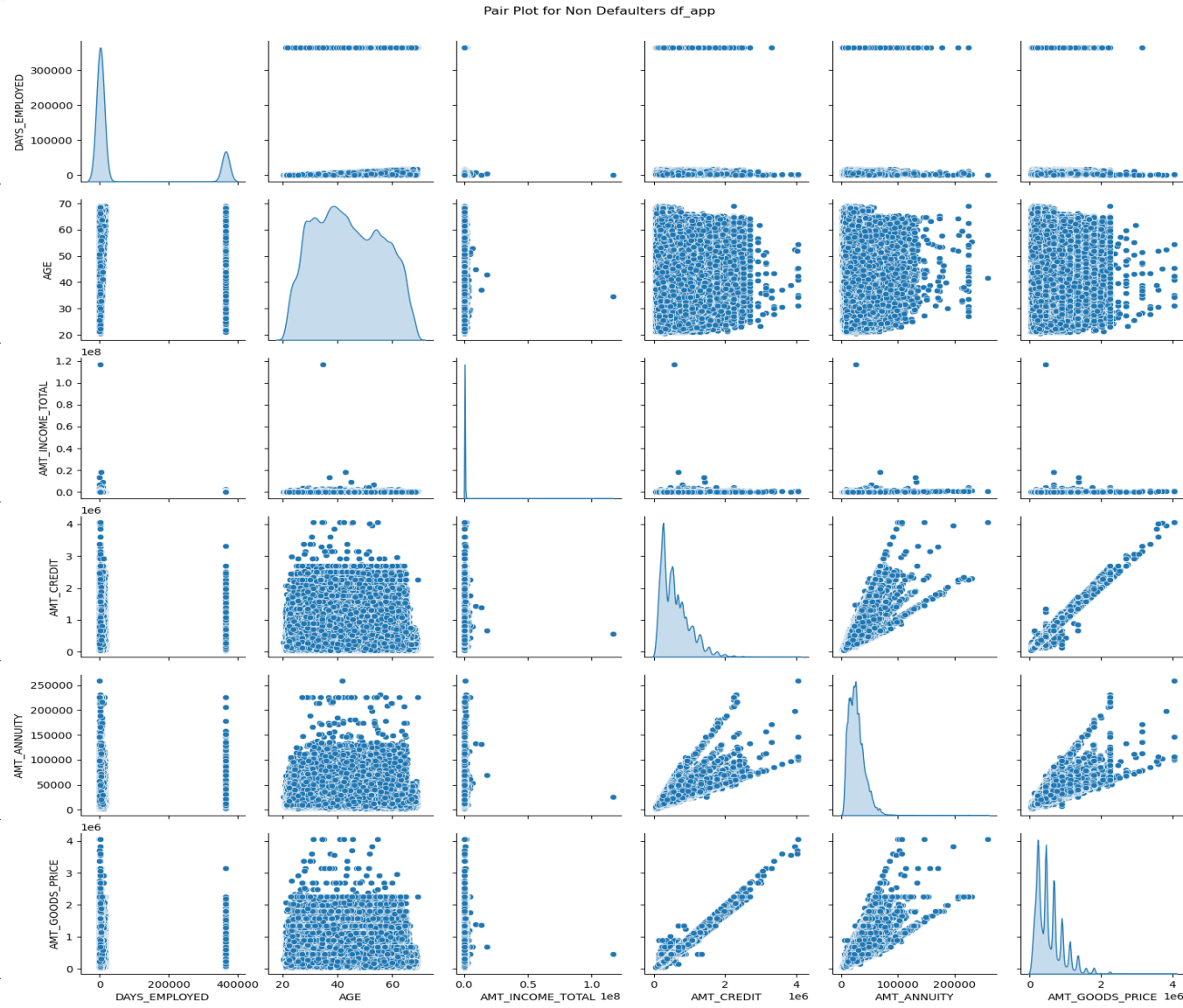
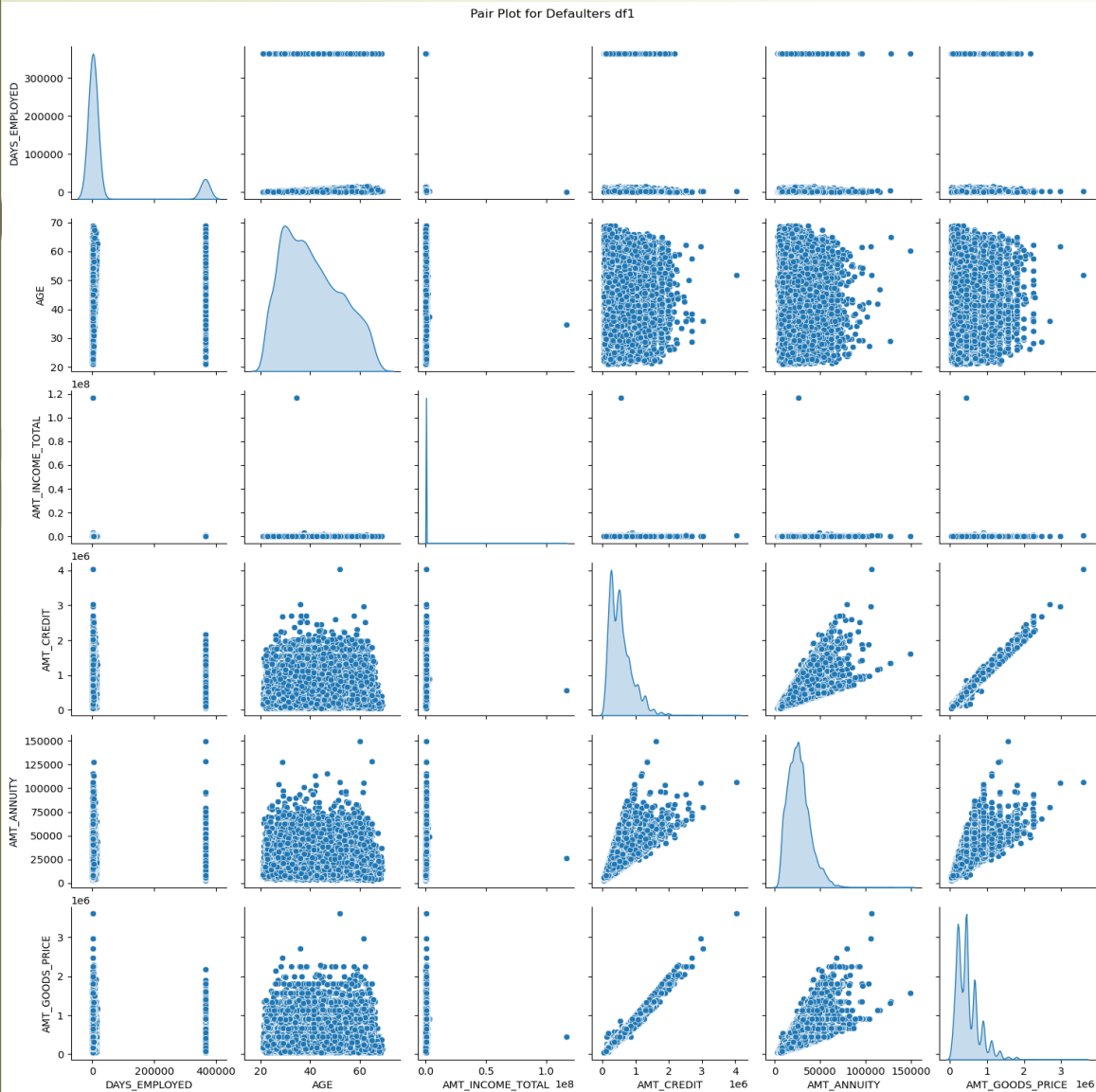
	DAYS_EMPLOYED
CNT_CHILDREN	-0.245174
CNT_FAM_MEMBERS	-0.238300
AMT_GOODS_PRICE	-0.068527
AMT_ANNUITY	-0.104975
AMT_INCOME_TOTAL	-0.140392
AMT_CREDIT	-0.070104
LIVE_CITY_NOT_WORK_CITY	-0.221387
REGION_POPULATION_RELATIVE	-0.007198
DAYS_BIRTH	0.626114
DAYS_EMPLOYED	1.000000

## Defaulters

	DAYS_EMPLOYED
CNT_CHILDREN	-0.192864
CNT_FAM_MEMBERS	-0.186515
AMT_GOODS_PRICE	0.006648
AMT_ANNUITY	-0.081207
AMT_INCOME_TOTAL	-0.014977
AMT_CREDIT	0.001930
LIVE_CITY_NOT_WORK_CITY	-0.198484
REGION_POPULATION_RELATIVE	0.015532
DAYS_BIRTH	0.582185
DAYS_EMPLOYED	1.000000

# Multivariate Analysis

38



# Previous Data: Cleaning

1. 37 columns were reduced to 26 as 11 columns had null values which were more than 40%.
2. 7 columns had less than 23% null values .
  1. Used median statistics for the following columns : AMT\_GOODS\_PRICE, AMT\_ANNUITY, AMT\_CREDIT, CNT\_PAYMENT, PRODUCT\_COMBINATION
  2. Did not impute with anything or capped the null rows for DAYS\_FIRST\_DUE and DAYS\_FIRST\_DRAWING.

# Previous Data: Anomalies

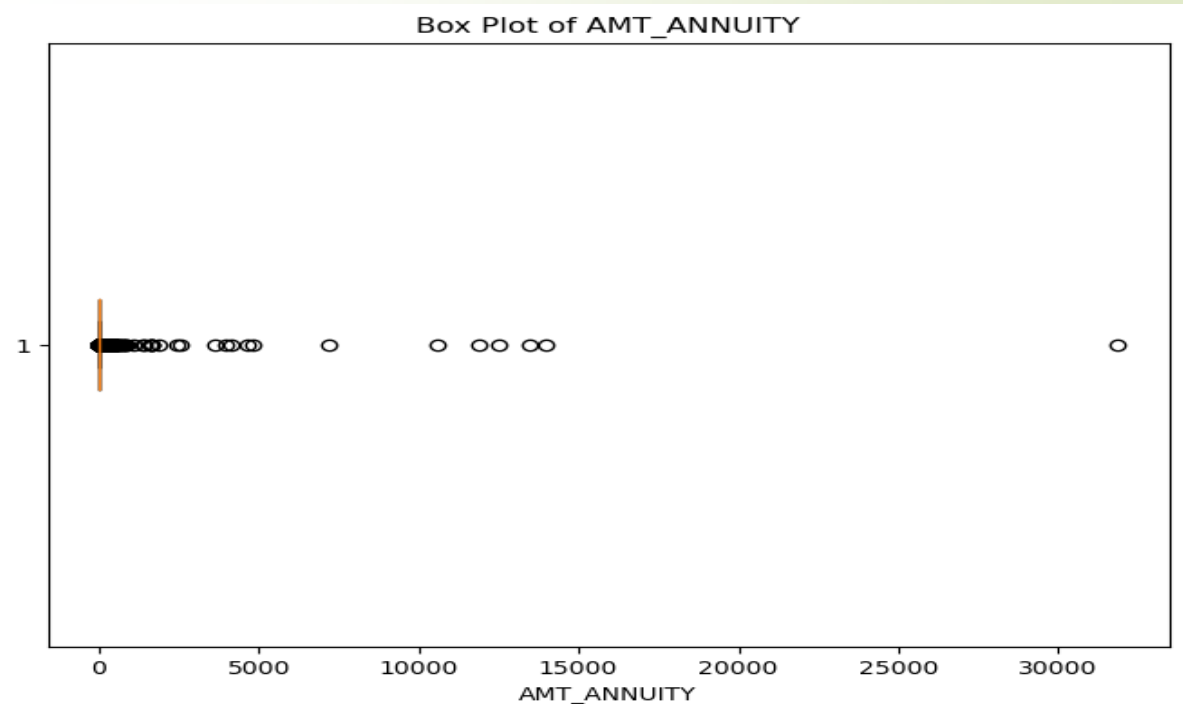
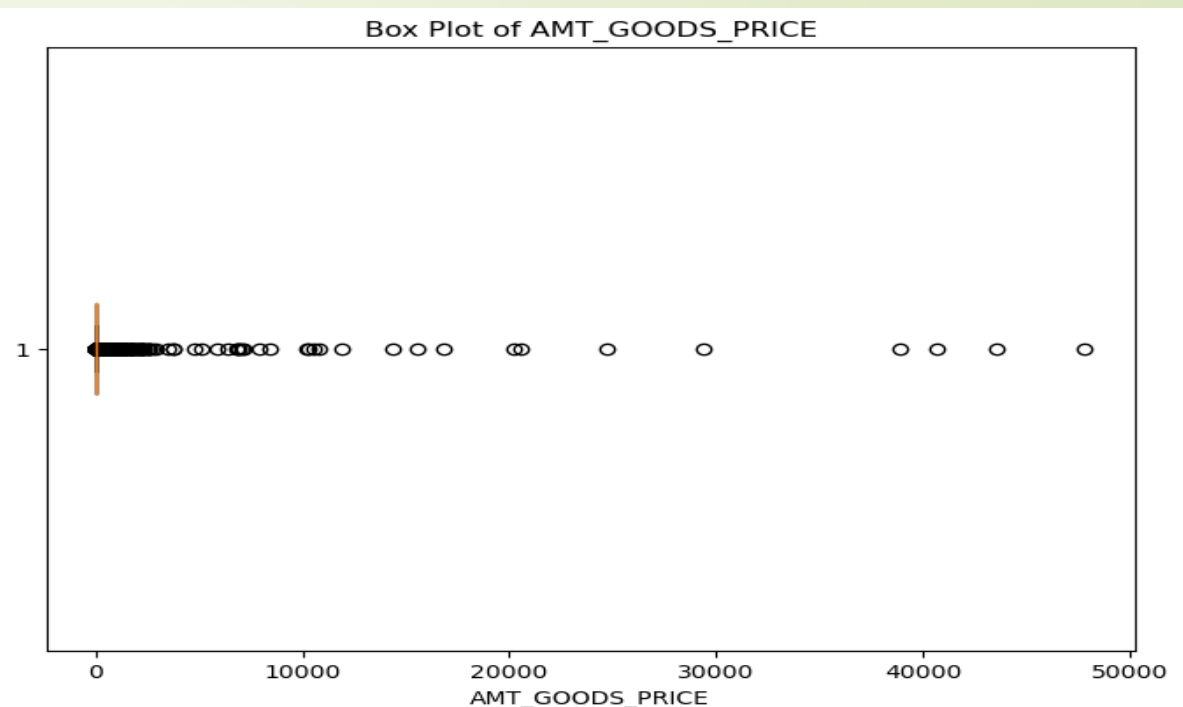
1. NAME\_CONTRACT\_TYPE and NAME\_CLIENT\_TYPE columns had “XNA” rows replaced it with “Unknown”
2. NAME\_PAYMENT\_TYPE column had “XNA” rows with significantly large numbers , hence did not make any changes.

# Univariate Analysis

41

## AMT\_GOODS\_PRICE AMT\_ANNUITY:

Both the columns have outliers and data points are beyond 75<sup>th</sup> percentile.

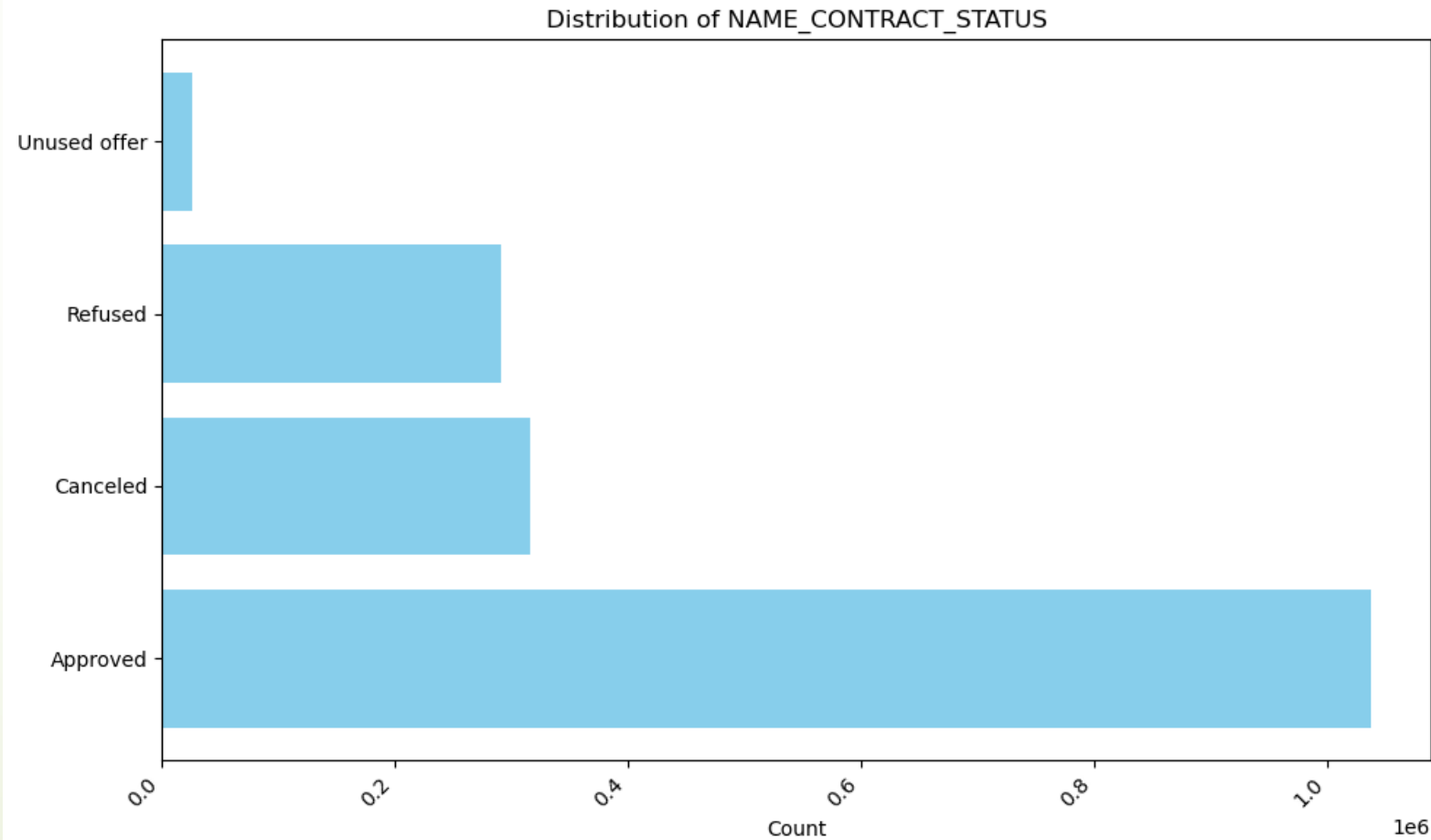


# Univariate Analysis

42

## NAME\_CONTRACT\_STATUS

For most of the clients loan was **Approved** in the previous application data



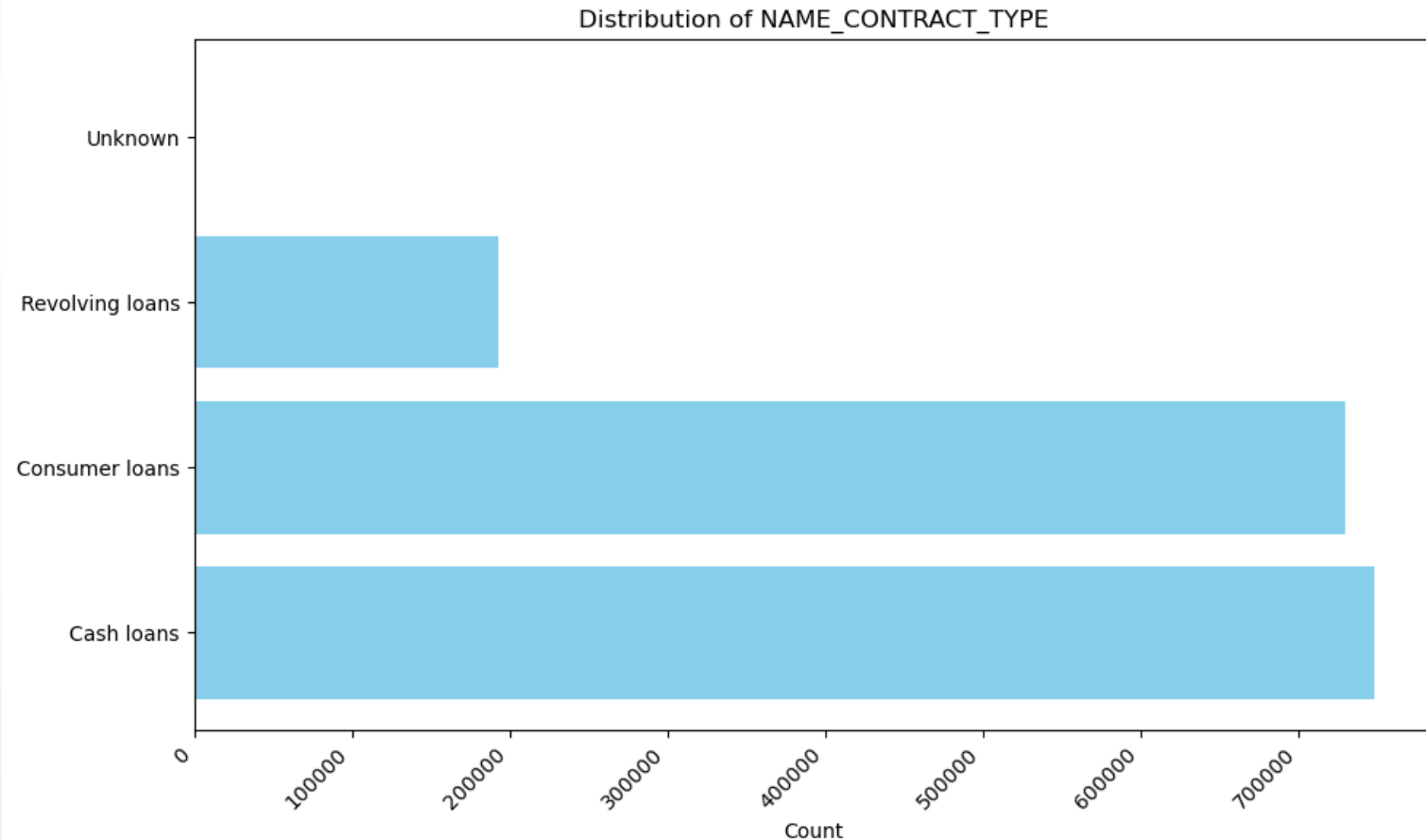


# Univariate Analysis

43

## NAME\_CONTRACT\_TYPE

**Cash and Consumer loans** are the most common contract product of the previous clients.

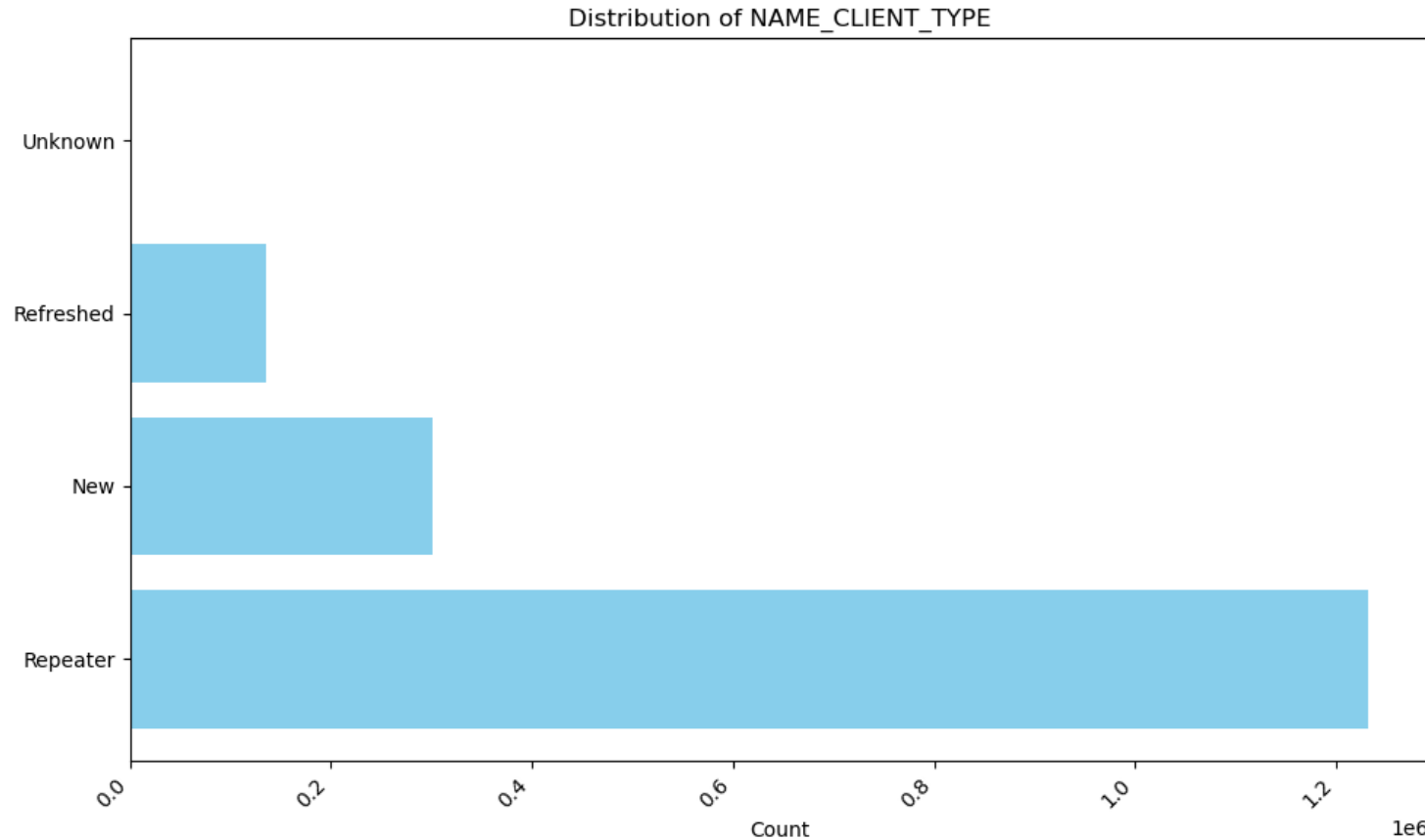


# Univariate Analysis

44

## NAME\_CLIENT\_TYPE

**Repeaters** are the most common clients seen . They have come back for more loans.

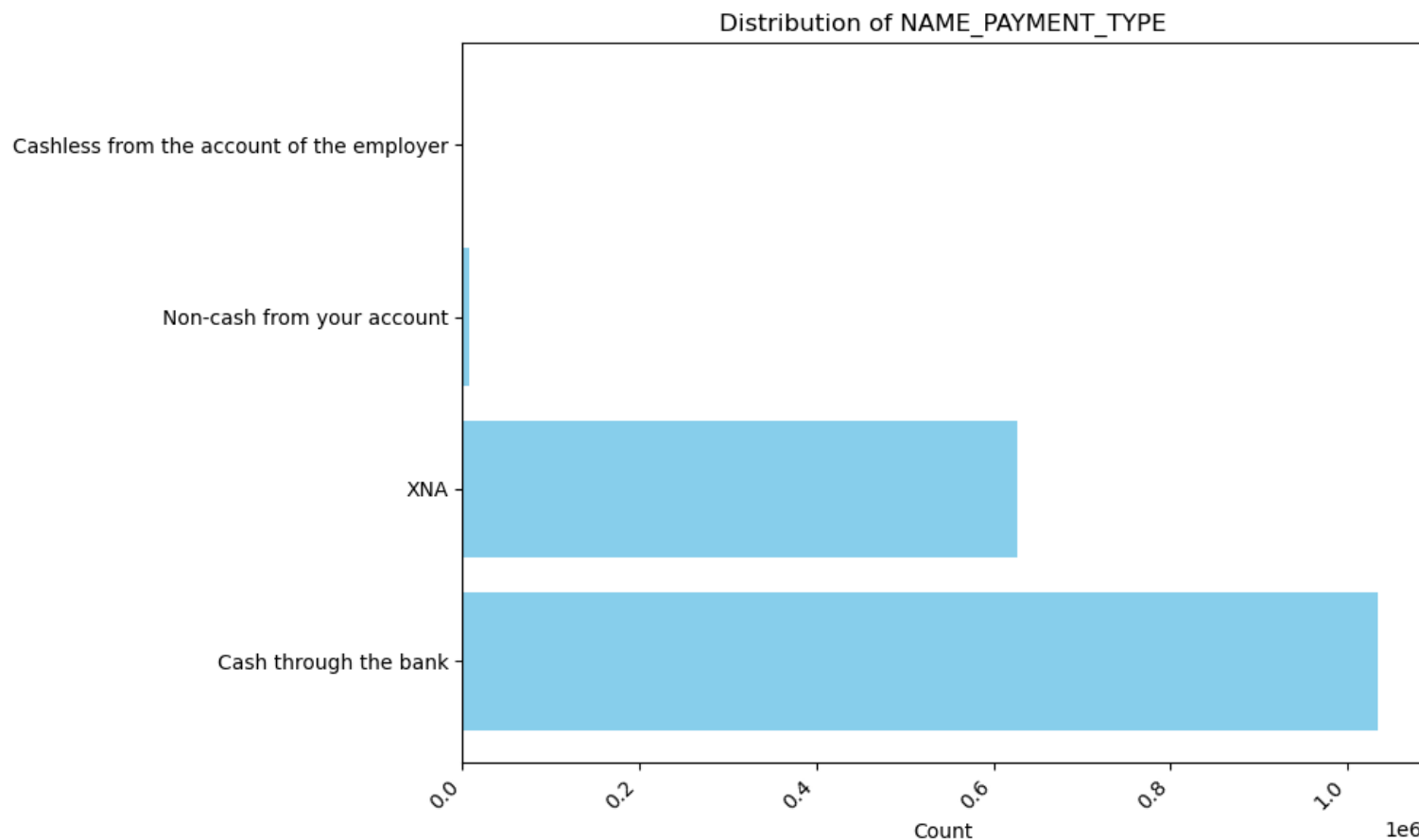


# Univariate Analysis

45

## NAME\_PAYMENT\_TYPE

**Cash through the bank** is the most preferred method used by the previous clients



# MERGED DATA

46

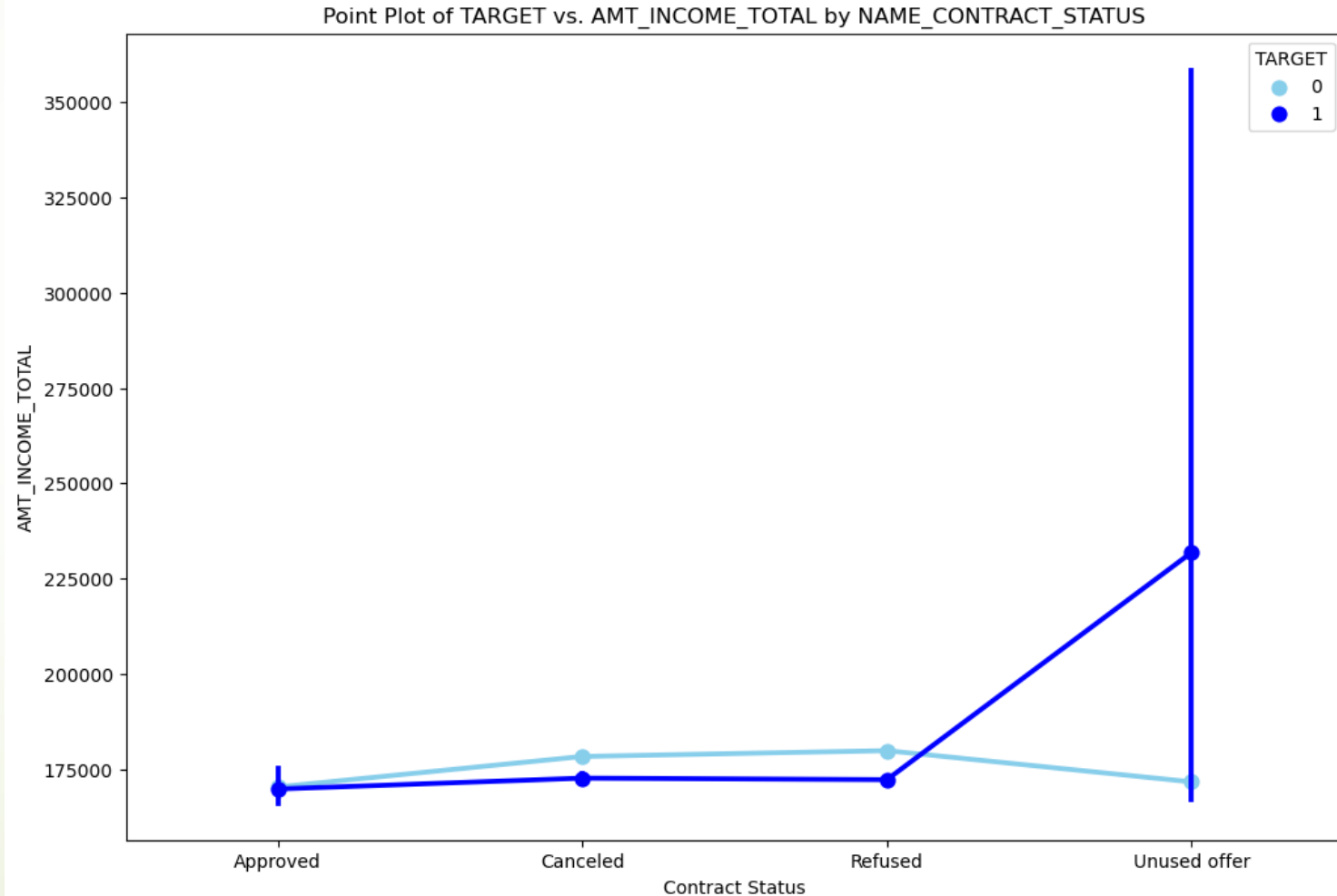
Both the data frames were merged using SK\_ID\_CURR.  
With this we have 78 rows and 1,413,701 columns

# Multivariate Analysis

47

## TARGET column comparison with NAME\_CONTRACT\_STATUS & AMT\_INCOME\_TOTAL

There are significant number of “**Unused offers**” by the defaulted clients who have higher income.

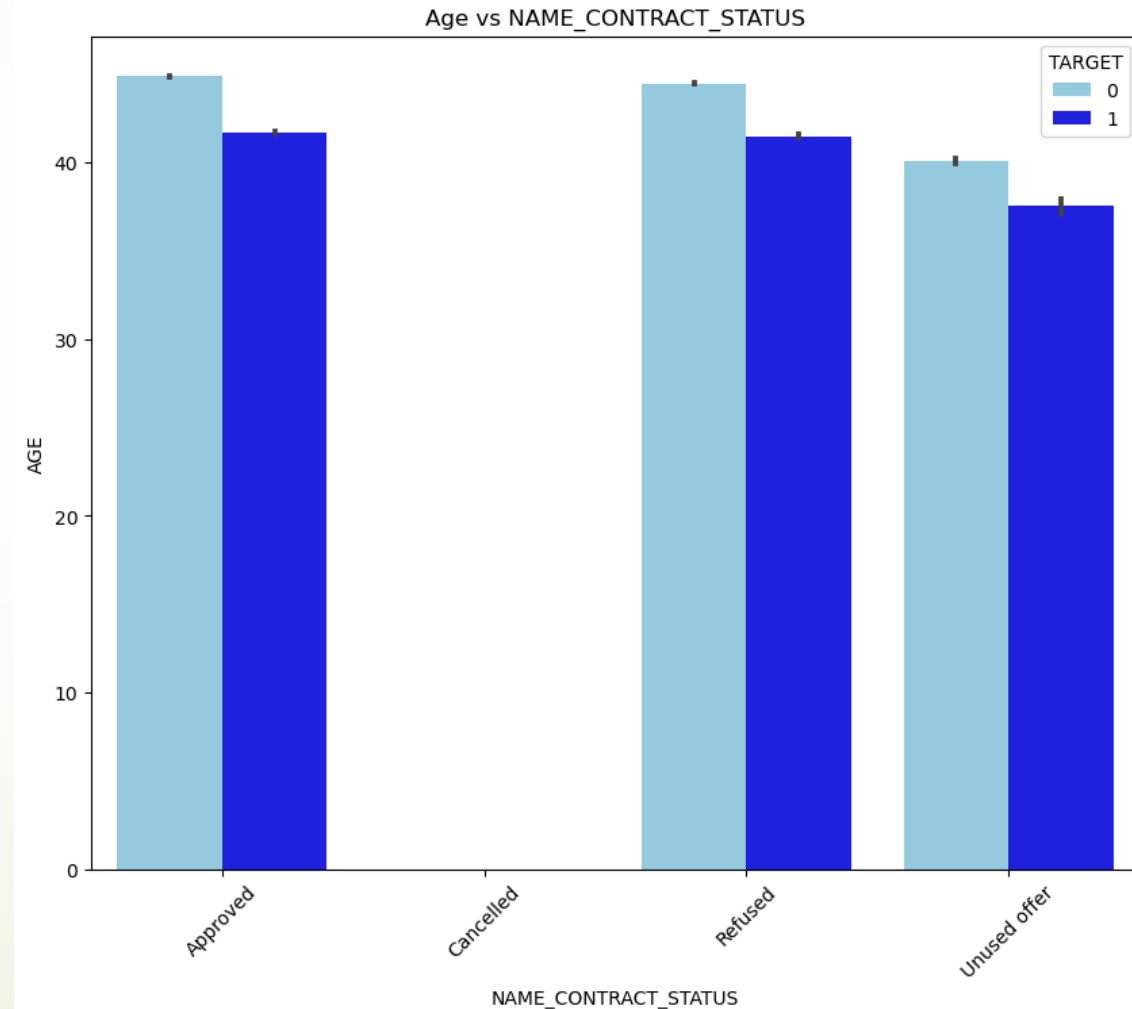


# Multivariate Analysis

48

## TARGET column comparison with NAME\_CONTRACT\_STATUS & AGE

Loan Approval rate is more than Cancelled.  
Clients who have Refused or Unused the loan offer is more among both defaulters and non-defaulters.





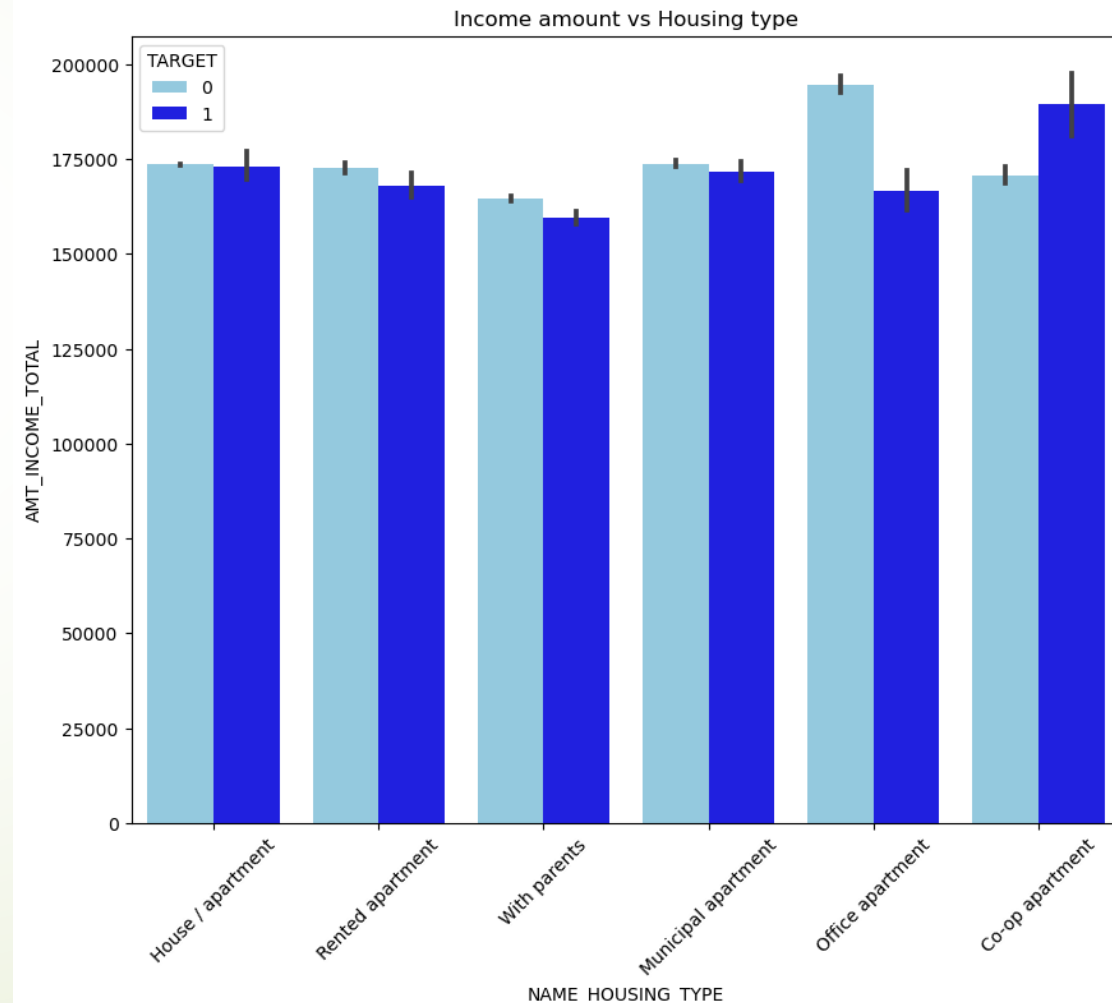
# Multivariate Analysis

49

## TARGET column comparison with NAME\_HOUSING\_TYPE & AMT\_INCOME\_TOTAL

There are significant number of “**Office apartment**” housing type by **the non-defaulted** clients who have **higher income**.

Client who lived in “**Co-op apartment**” housing type by the **defaulters** also had **higher income**.

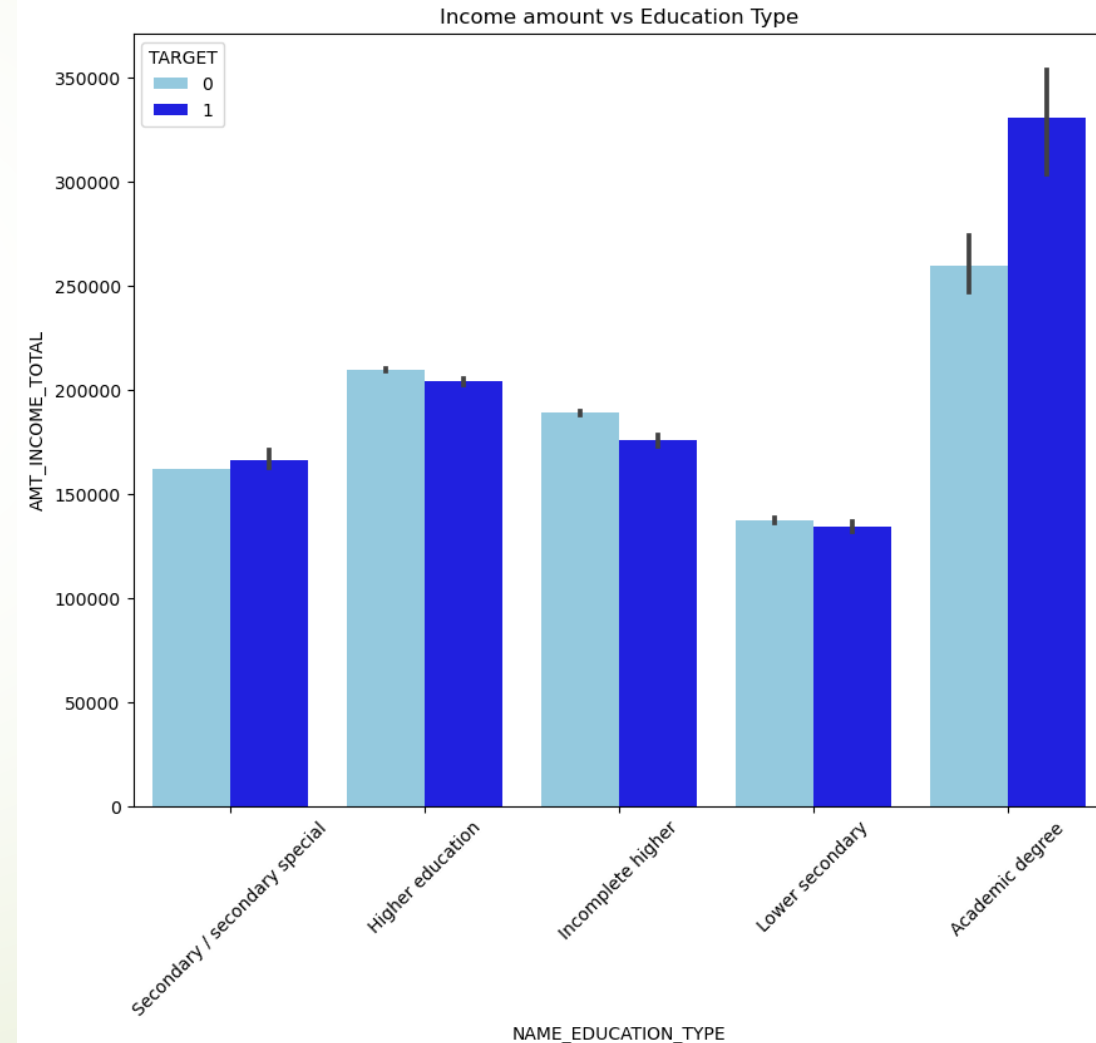


# Multivariate Analysis

50

## TARGET column comparison with EDUCATION TYPE & AMT\_INCOME\_TOTAL

There are significantly higher number of “**Academic degree**” holders by **both defaulter and non-defaulter** clients who have **higher income**.

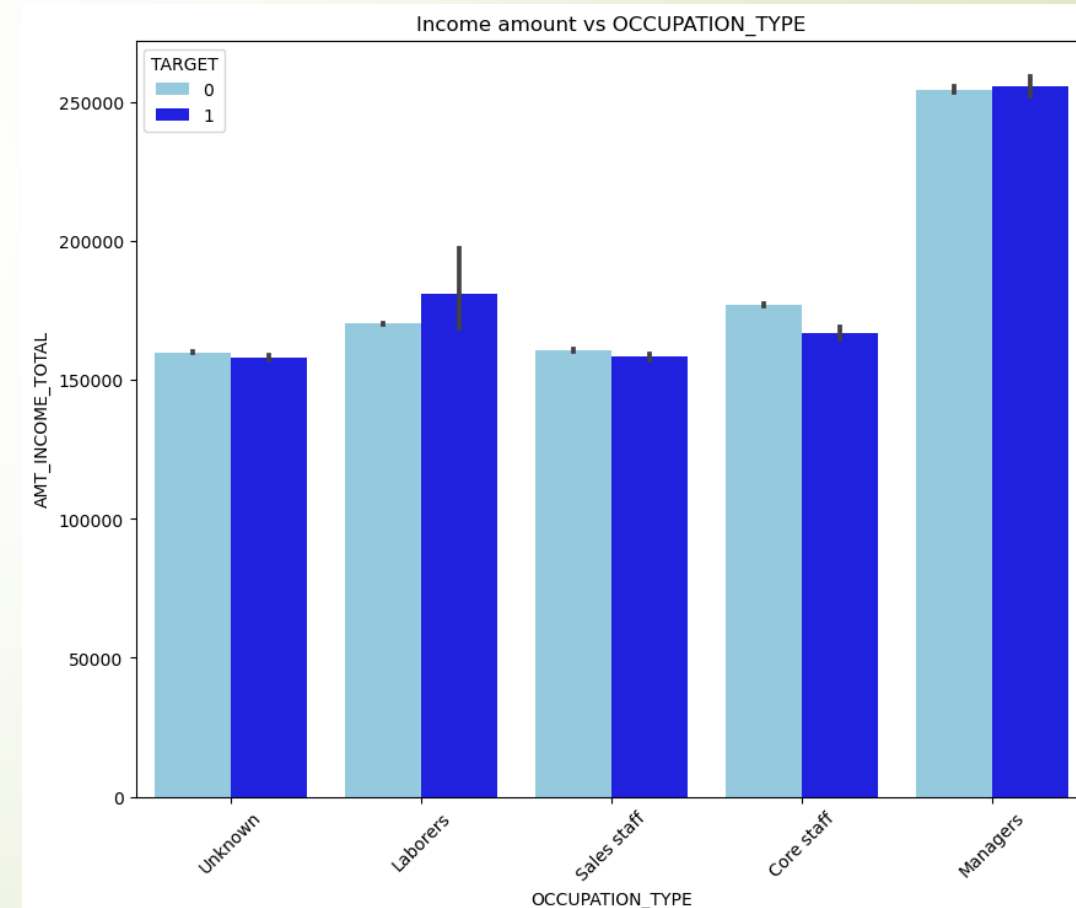


# Multivariate Analysis

51

## TARGET column comparison with OCCUPATION\_TYPE & AMT\_INCOME\_TOTAL

There are significantly higher number of “**Managers**” in **both defaulter and non-defaulter** clients who have **higher income**.

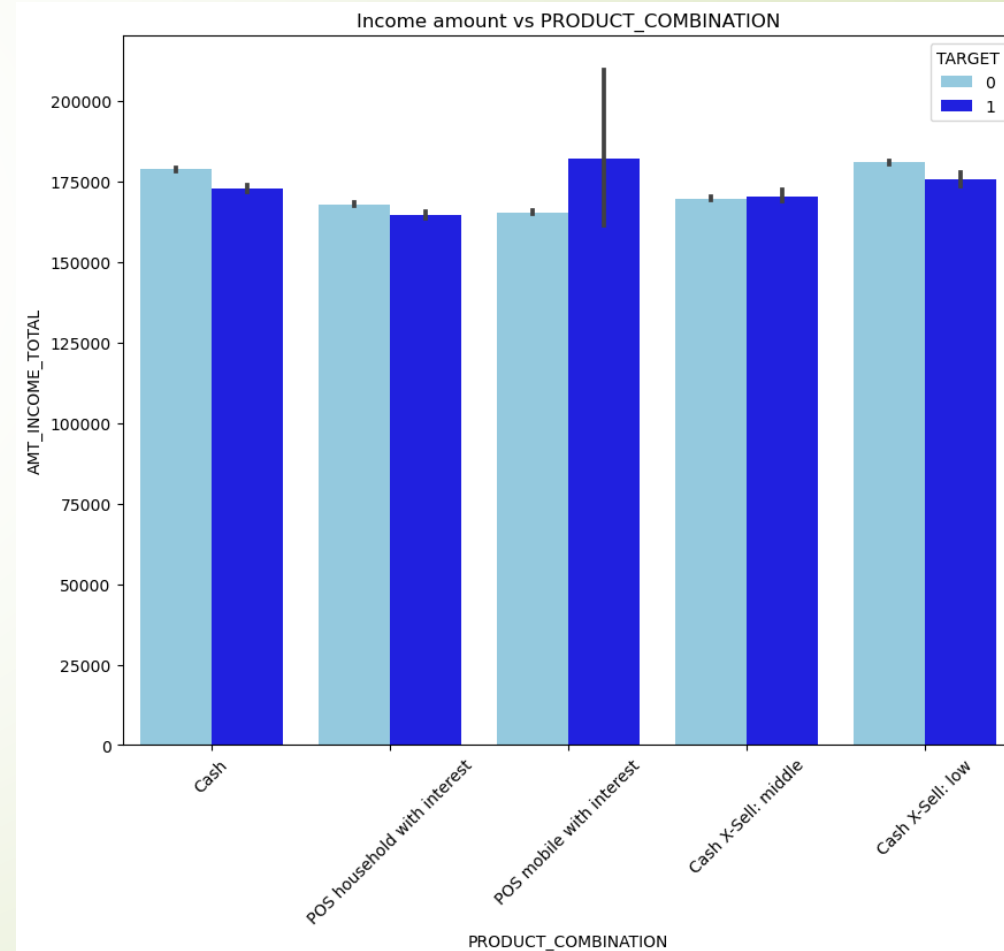


# Multivariate Analysis

52

## TARGET column comparison with PRODUCT\_COMBINATION & AMT\_INCOME\_TOTAL

Population who purchased “**POS Mobile with interest**” in the **Defaulter** clients have higher income and population who bought “**Cash x-sell low**” among the **Non-defaulter**s have higher income.

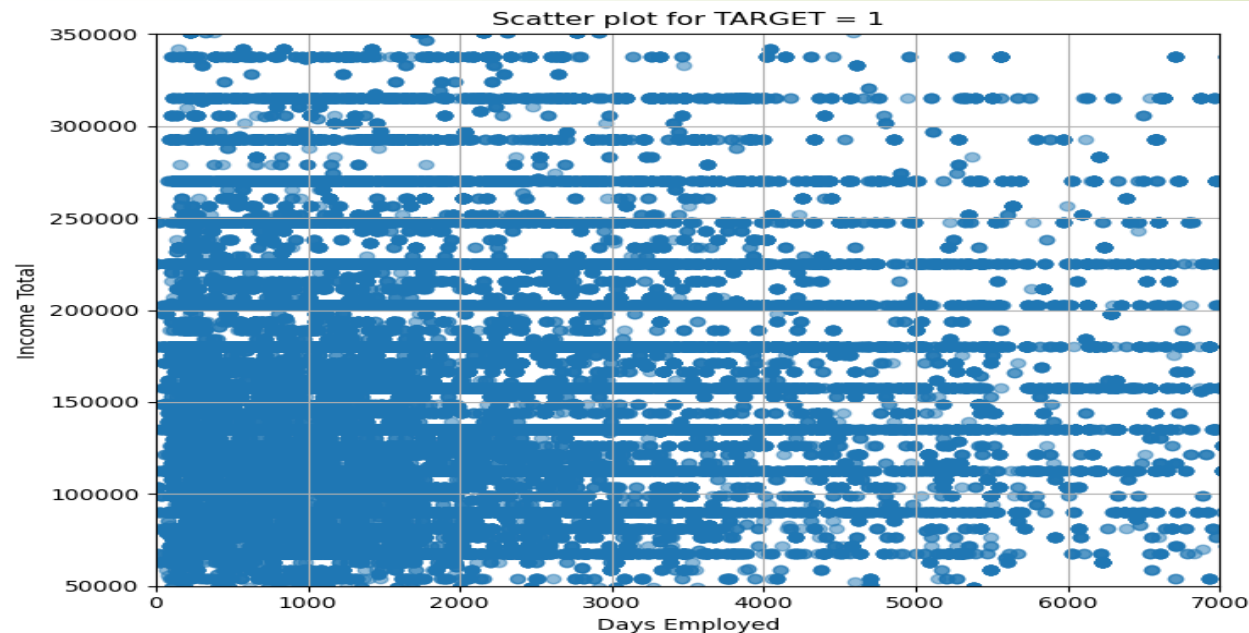
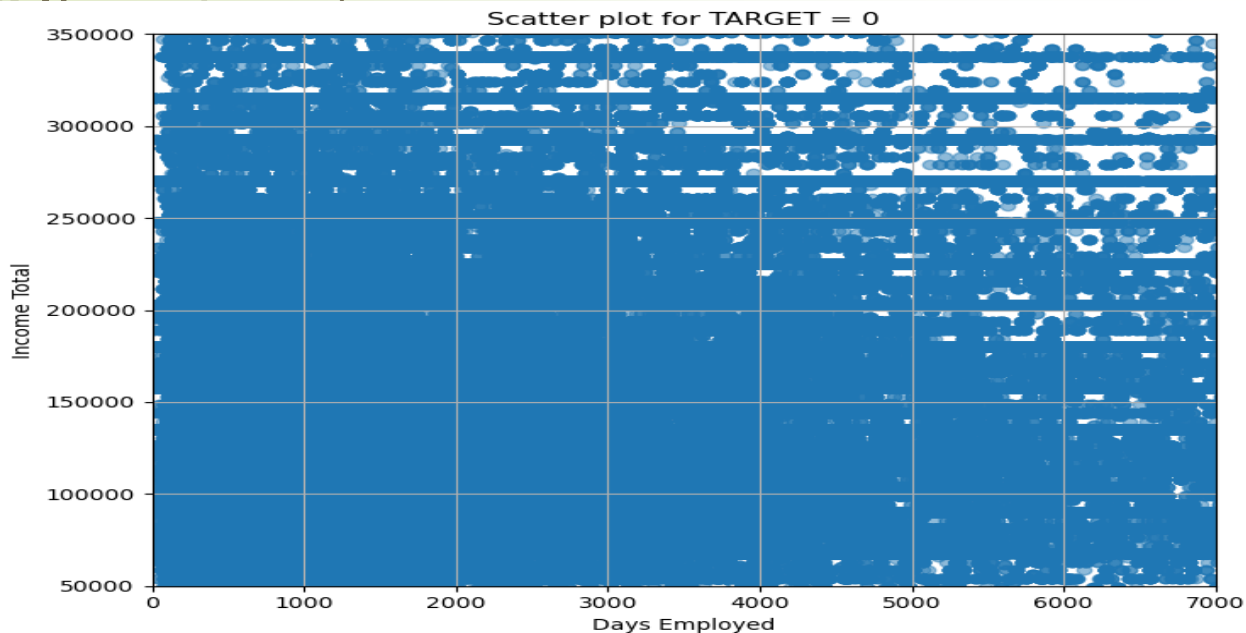


# Multivariate Analysis

53

## TARGET column comparison with DAYS\_EMPLOYED & AMT\_INCOME\_TOTAL

Clients with higher DAYS\_EMPLOYED had no payment difficulties, they also had higher income. Whereas clients' who had payment difficulties we see less population who were more than 3000 DAYS\_EMPLOYED.

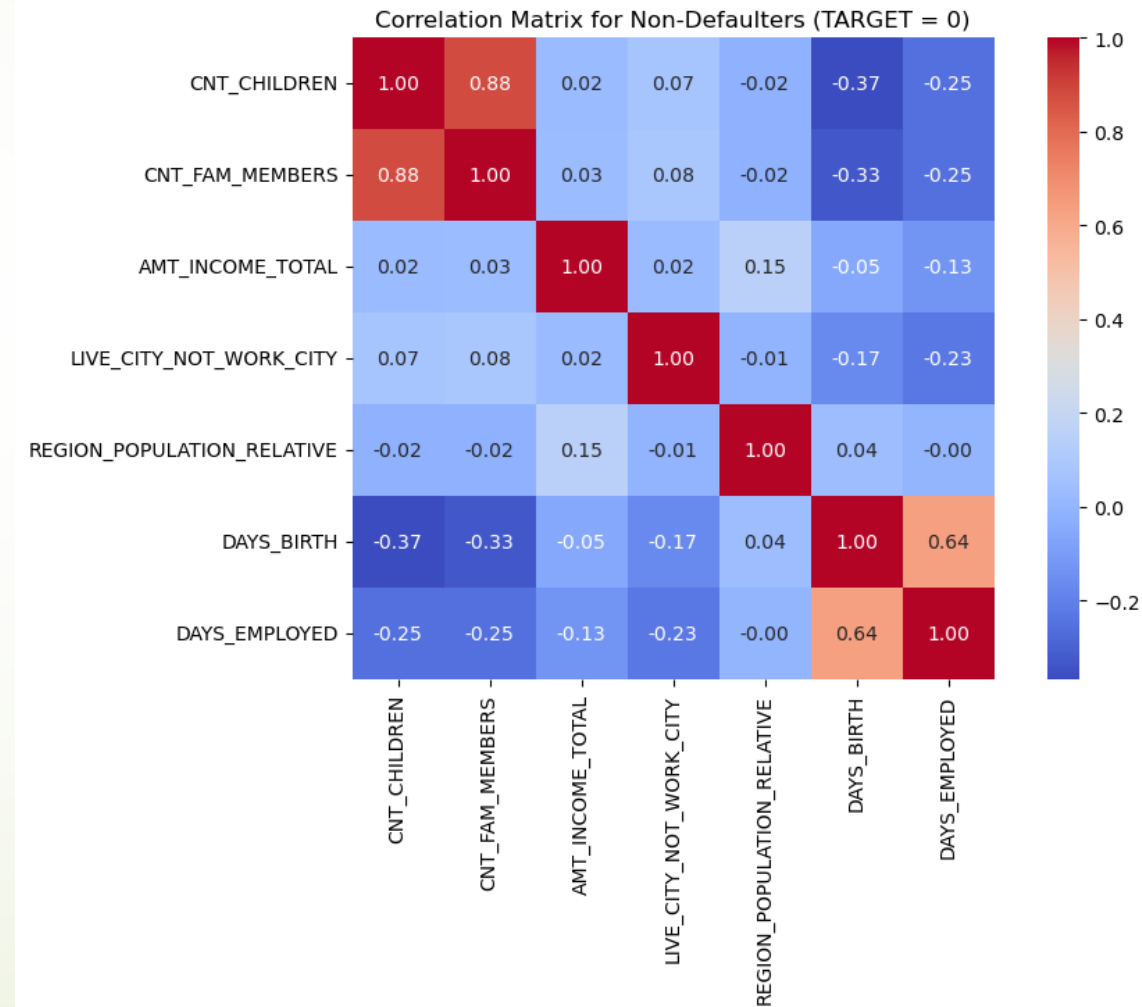


# Multivariate Analysis

54

## CORRELATION MATRIX

From the merged data we see strong correlation between **DAYS\_BIRTH** and **DAYS\_EMPLOYED** which is around **0.64** in **Non-Defaulters**



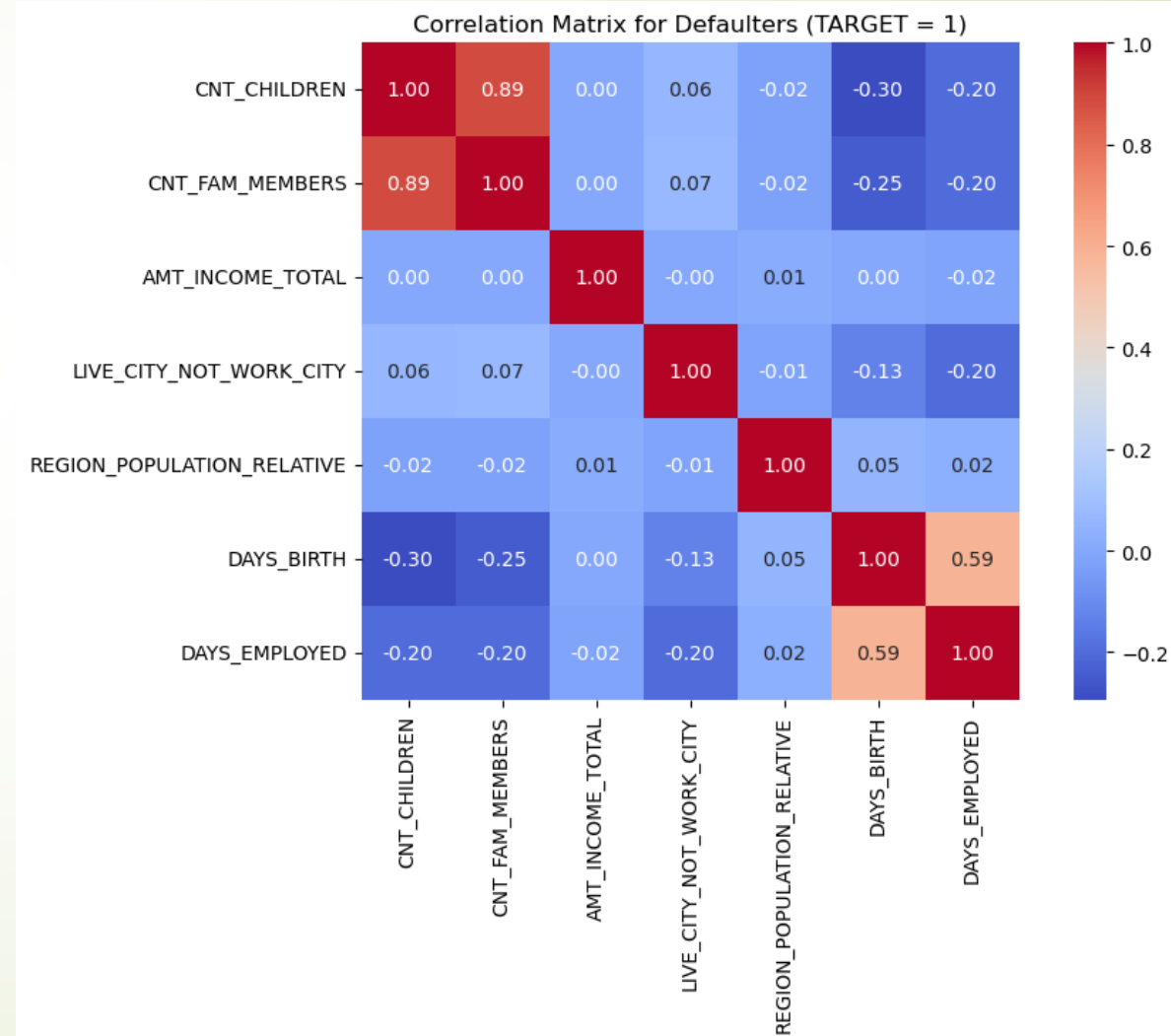


# Multivariate Analysis

55

## CORRELATION MATRIX

From the merged data we see strong correlation between **DAYS\_BIRTH** and **DAYS\_EMPLOYED** which is around **0.59** in **Defaulters**





**NAME\_CONTRACT\_TYPE:** Most of the clients choose Cash loans or Consumer loans .

**CODE\_GENDER :** Female populations took more loan then male, though they earned lesser than males. Hence, they are likely to face more payment difficulties.

**FLAG\_OWN\_CAR :** most of the population did not own a car.

**FLAG\_OWN\_REALTY:** three fourth of the population did own a flat or house.

**CNT\_CHILDREN :** population with no children were mostly seen taking loans.

**AMT\_INCOME\_TOTAL :** Majority of the population had income range in 76 to 113K .

**AMT\_CREDIT :** Most of them took credit in the range of 222 to 259K.

**AMT\_ANNUITY:** 80 to 90% of the population had annuity in the range of 40 to 76K.

**AMT\_GOODS\_PRICE:** Several of the population bought goods priced with in the range of 222 to 259K.

**NAME\_TYPE\_SUITE:** Many of them were “Unaccompanied” while taking loans.

**NAME\_INCOME\_TYPE:** More number of “working” group have taken the loan and are likely to default too.

**OCCUPATION\_TYPE :** Laborers are found to default more .

**CNT\_FAM\_MEMBERS:** population who have around 2 members are more likely to default.

# INSIGHTS

57

**NAME\_EDUCATION\_TYPE:** Secondary and secondary special have taken more loans and are more likely to default.

**NAME\_FAMILY\_STATUS:** More number of married population have taken loan and are likely to default too.

**NAME\_HOUSING\_TYPE:** population who have house/apartment have taken loans and more likely to default.

**DAYS\_EMPLOYED:** More number of population have worked less than 3000 days or 8 years.

**ORGANIZATION\_TYPE:** people who work in “Business Entity 3” have taken more loan and are more likely to default.

**NAME\_PAYMENT\_TYPE:** lot of them pay the bank by “Cash through the bank”

**NAME\_CLIENT\_TYPE:** Majority of Repeaters have taken loan.

**AGE or DAYS\_BIRTH :** Majority of the population within the age group of 31-40 have taken loans and they are found to default more when compared to 41-50 and 51 – 60.

Clients who can pay back loan on time OR Bank to target these population for **approving** loans.

1. **Clients who have income range 76 to 186K**
2. **Clients who are in the age group of 21-30 and 41-50 .**
3. **Clients who have higher education and incomplete higher.**
4. **Clients who are employed in an organization for more than 4000days or 10 years .**
5. **Clients who come back after paying the initial loans , Repeaters.**
6. **Clients who have occupation such as Core staff and Sales staff .**
7. **Clients who work in organization such as Medicine, Others and are Self Employed.**
8. **Clients who are single or had a civil marriage .**
9. **Clients who do not own a house or flat.**
10. **Male population with more income range.**

Clients who can not pay back loan on time OR Bank to target these population for **refusing loans**.

1. **Clients who have more than 2 number of children or family members.**
2. **Academic degree holders are more likely to default.**
3. **Clients who are managers and laborers.**
4. **Clients who work in Business Entity 3 organization type .**
5. **Clients who are married or separated.**
6. **Female population with less income.**

**Thank you**