

EDA CREDIT CASE STUDY

BY –

SMARANIKA SAHOO & NITISH RATHORE

LOAD REQUIRED LIBRARY & FILTER WARNINGS

- THIS STEP INVOLVE INCLUDES IMPORTING OF SOME LIBRARIES OF PYTHON ALONG WITH LIBRARY TO IGNORE WARNINGS.

LOAD THE DATA

- THIS STEP INCLUDES READING THE .CSV FILE THROUGH PANDA LIBRARY .

CHECK THE DATA

GLIMPSE THE DATA

- USING SHAPE, DESCRIBE, INFO FUNCTIONS TO CHECK NO. OF COLUMN, NO. OF ROWS, DATA TYPES AND FEW MATHEMATICAL STATS LIKE STANDARD DEVIATION, MEAN, MAX, VARIOUS QUANTILES.

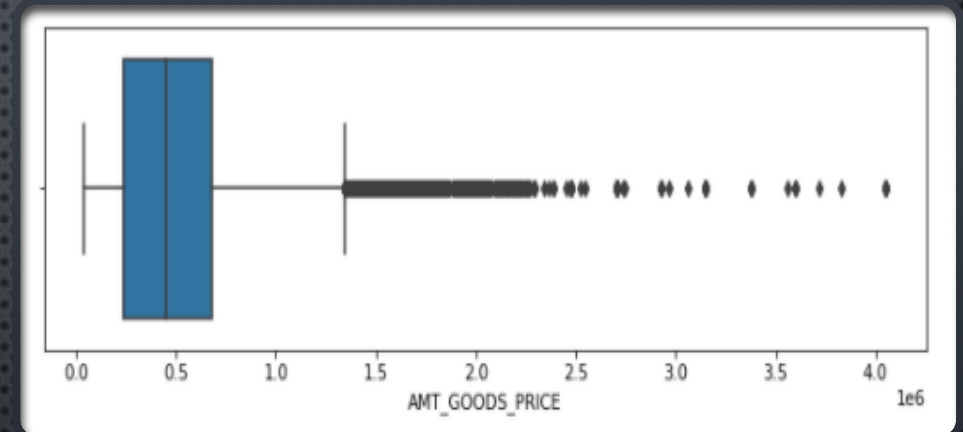
CHECK MISSING VALUES

- THIS STEPS INVOLVE CHECKING THE NULL VALUES COUNT & PERCENTAGE OF NULL VALUES IN EVERY COLUMNS.
- MOVING AHEAD, WE'LL REMOVE THE COLUMNS WHERE PERCENTAGE OF NULL VALUES IS GREATER THAN 50 %.
- CHECK THE DATA AND REMOVE COLUMNS WHICH ARE NON OF USE FOR TO MAKE INFERENCE FURTHER.

IMPUTATION OF NULL VALUES

Basically, there are two ways to Impute Null Values from Respective Columns.

- For Continuous variables :- To impute Null Values from Column with Continuous variables, Generally Median is used if There is Outlier in Data otherwise Mean can be used.
- For Categorical variables :- To impute Null Values from Column with Categorical variables, Generally Mode is used.



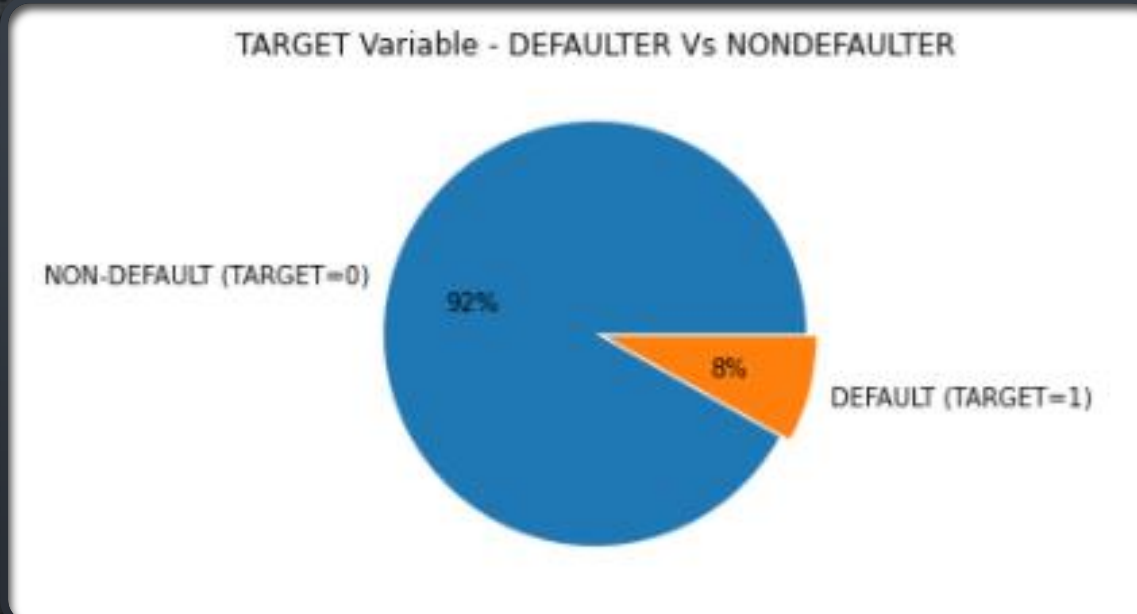
We can see in the Above Diagram, there are Outliers, so Here Median will be used for Imputation.

Note :- For the Given Set of Data, we Have not done Null Value Imputation, As Mentioned in Problem Statement.

CORRECTION IN COLUMNS

- BINNING OF AMT_INCOME_TOTAL AND CREATED NEW COLUMN AS INCOME_GROUP WHICH INCLUDE SUB-POINT TO DENOTE INCOME RANGE. SIMILARLY BINNED AGE COLUMN AND CREATED AGE_GROUP.
- CREATED NEW COLUMN AS AGE FROM DAYS_BIRTH THEN DROPPED DAYS_BIRTH COLUMN.
- THERE ARE FEW COLUMN WITH DAYS AS UNIT, THAT HAVING SOME NEGATIVE VALUES, SO CONVERTED THOSE COLUMNS TO POSITIVE.

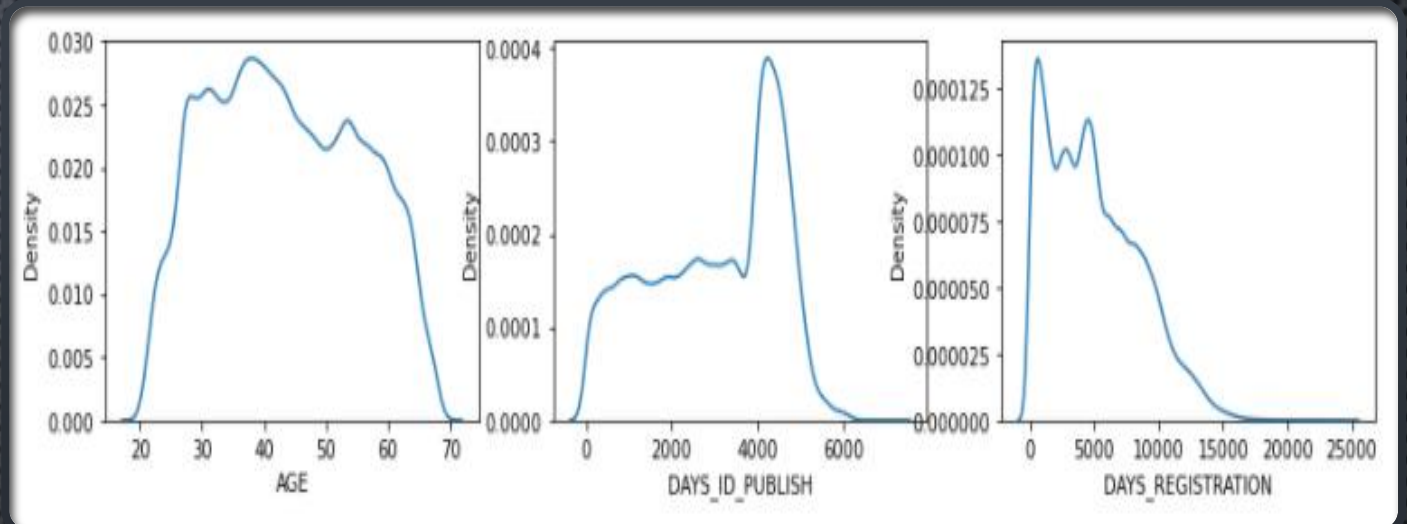
CHECK DATA UNBALANCE



- WE ARE CHECKING THE DATA UNBALANCE FOR THE **TARGET** VARIABLE.
- THIS IS TO FIND THE RATIO BETWEEN THE NON-DEFAULT CLIENTS & DEFAULT CLIENTS .
- **TARGET 1** MEANS DEFAULT CLIENTS WHICH IS 8%
- **TARGET 0** MEANS NON- DEFAULT CLIENTS WHICH IS 92%
- RATIO BETWEEN BOTH THE CASES IS 11.39

UNIVARIATE ANALYSIS

(DISTPLOT FOR **AGE**, **DAYS_ID_PUBLISH** & **DAYS_REGISTRATION**)



INFERENCES

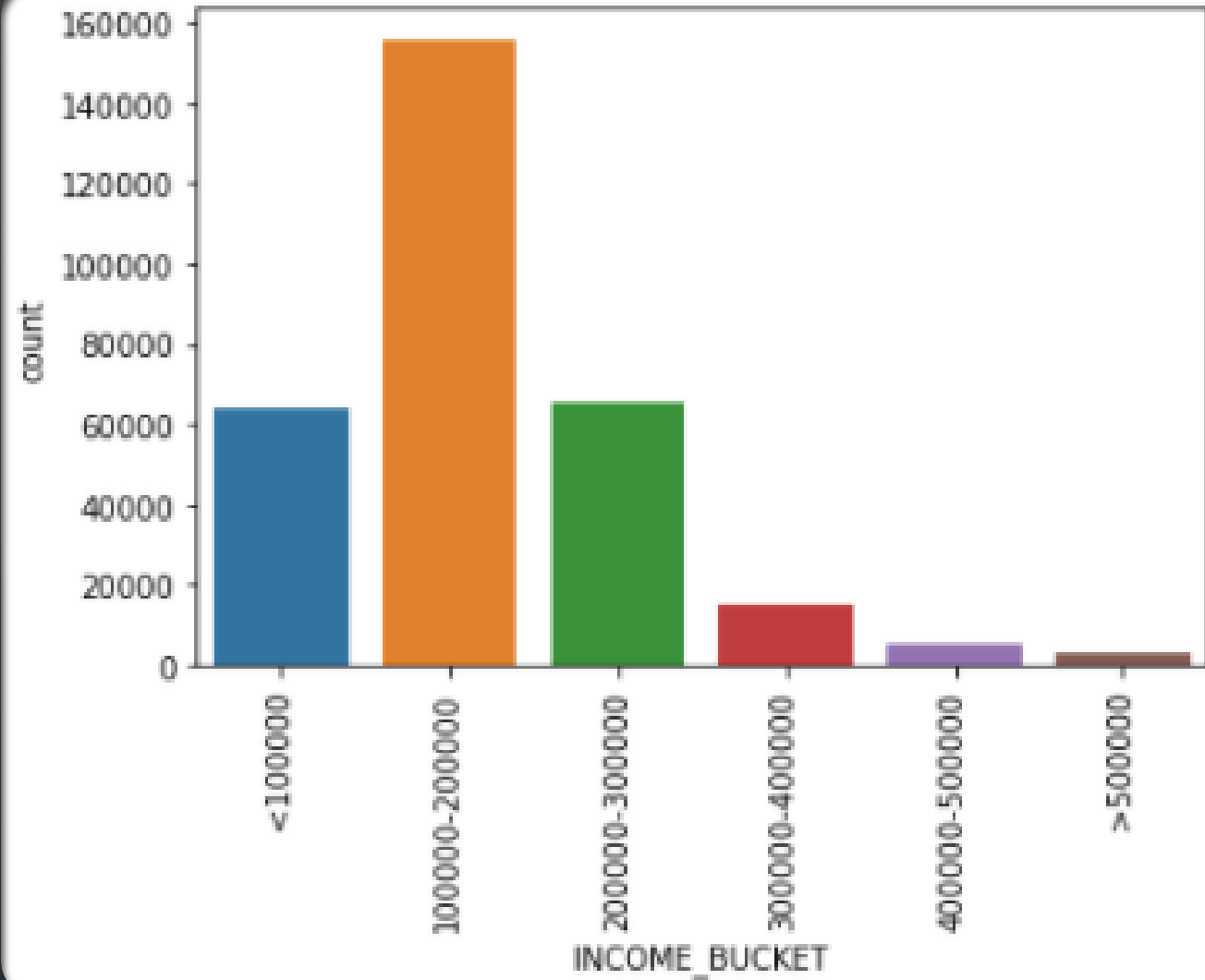
1. MOST PEOPLE WHO HAVE TAKEN LOANS ARE IN B/W 30-40 YEARS OF AGE.
2. CLIENTS WHO HAVE CHANGED THEIR ID DOCUMENTS DID IT MOSTLY 4 000-5000 DAYS BEFORE APPLYING THE LOAN.
3. CLIENTS CHANGING THEIR REGISTRATION WITHIN 5000 DAYS ARE TAKING MOST LOANS.

UNIVARIATE ANALYSIS

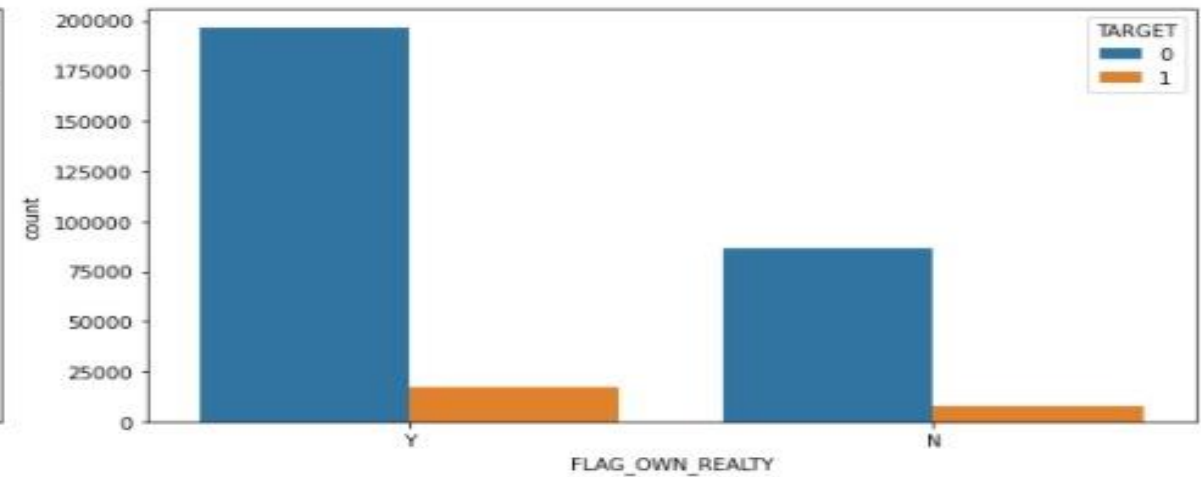
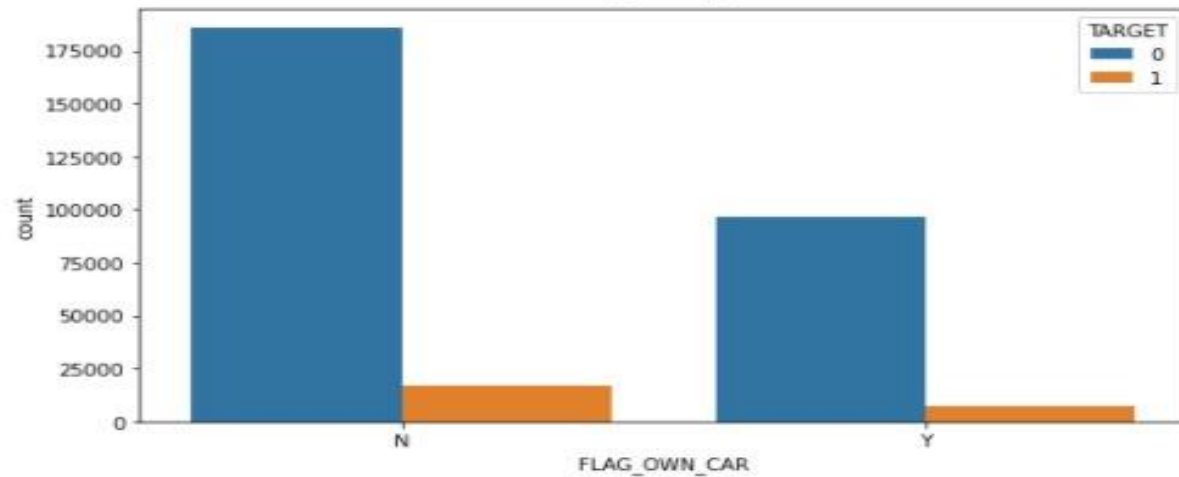
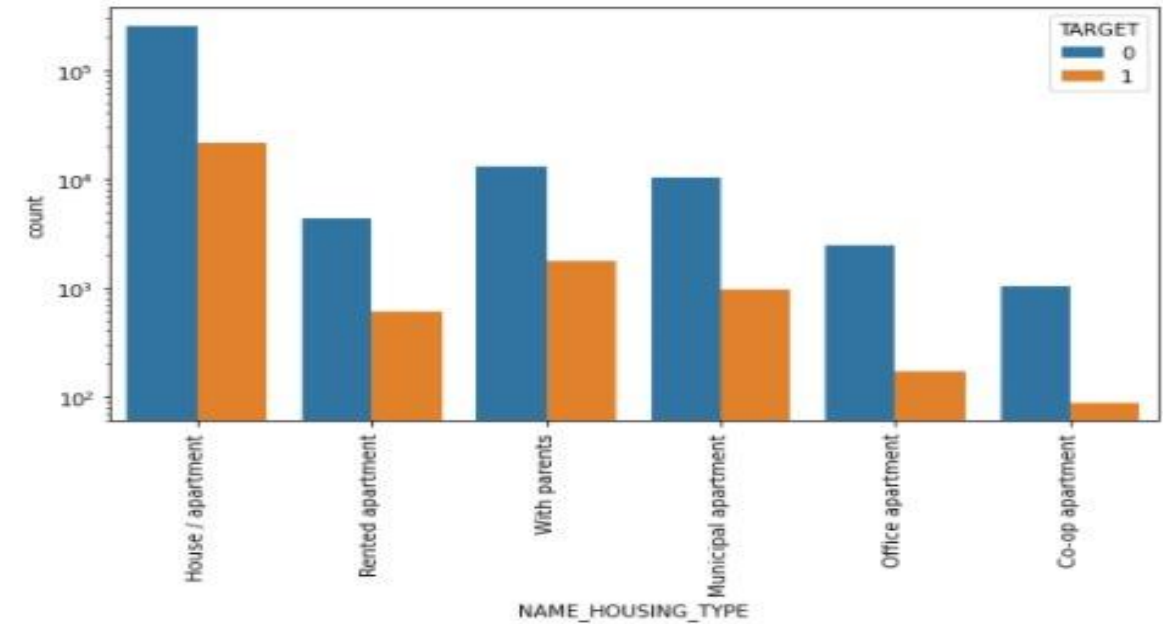
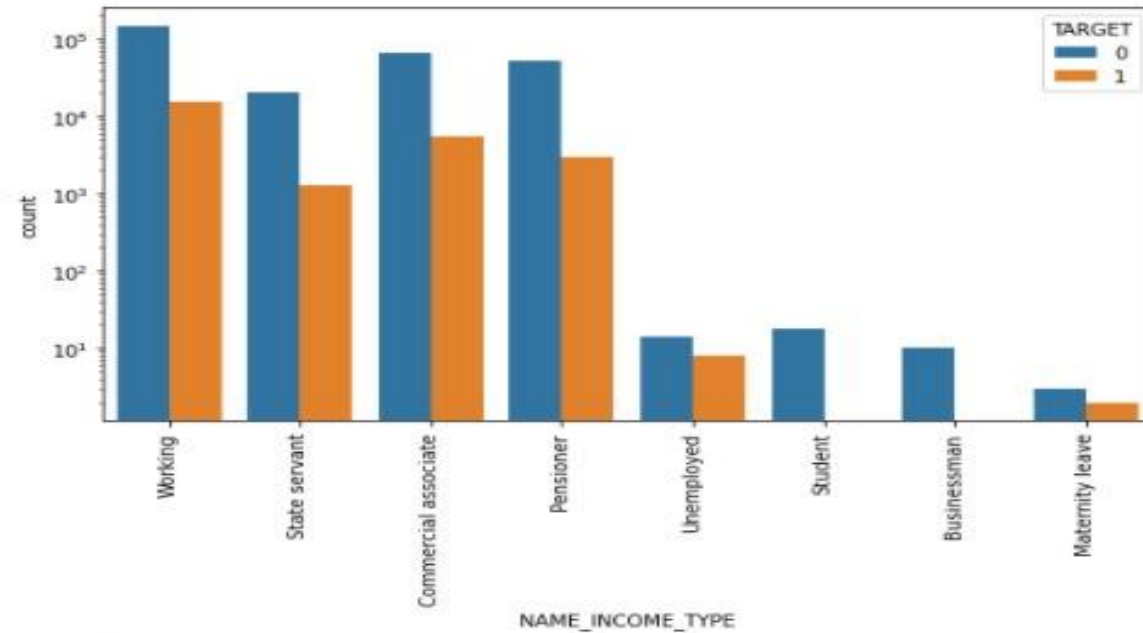
(PLOTTING SUBPLOT FOR
INCOME_BUCKET)

INFERENCE

- PEOPLE IN 100000-200000
HAVE TAKEN MOST LOANS
- PEOPLE WITH SALARY 200000 AND
MORE ARE LESS LIKELY TO TAKE LO
ANS COMPARISON TO PEOPLE WITH
SALARY LESS THAN 200000



COMPARISON WITH DEFAULT AND NON-DEFAULT



INFERENCES

1. THE INCOME BUCKET 100000-200000 HAVE HIGHEST DEFAULTERS ALSO HIGHEST NON DEFAULTERS.
2. THE INCOME BUCKET >500000 HAVE COMPARATIVELY LESS DEFAULTERS WHICH MEANS BANK WILL BE SAFER WITH THESE PEOPLE.
3. WORKING PEOPLE ARE MOST INCLINED TO TAKING LOANS ALSO, THEY ARE HIGHEST IN DEFAULTER GROUP AS WELL
4. STUDENT AND BUSINESSMAN ARE THE SAFER GROUP TO GIVE LOANS.

HOUSING TYPE:

1. PEOPLE WHO LIVE IN HOUSE/APARTMENT HAVE TAKEN MAXIMUM LOANS AND THEIR DEFAULTER PERCENTAGE IS HIGH.
2. PEOPLE WHO LIVE IN CO-OP APARTMENT HAVE SIGNIFICANTLY LESS DEFAULTERS

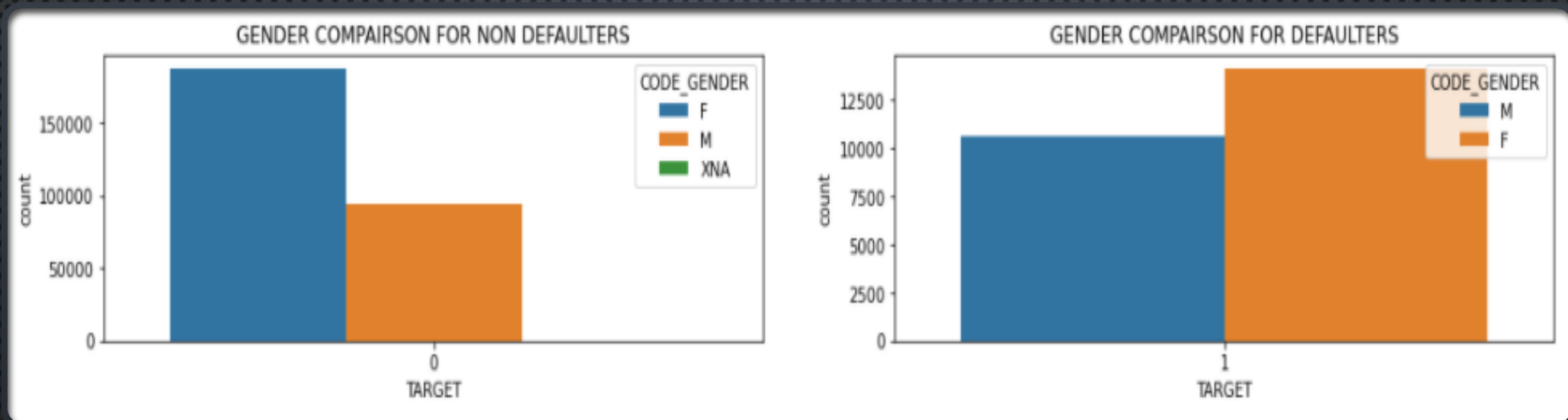
OWNING CAR

1. PEOPLE WHO DOESN'T OWNS A CAR HAVE LESS POSSIBILITY TO BECOME DEFAULTER.

OWNING A FLAT

1. PEOPLE WHO OWNS A FLAT HAVE LESS POSSIBILITIES TO BECOME DEFAULTERS.

BIVARIATE ANALYSIS (BETWEEN TARGET AND CODE_GENDER)

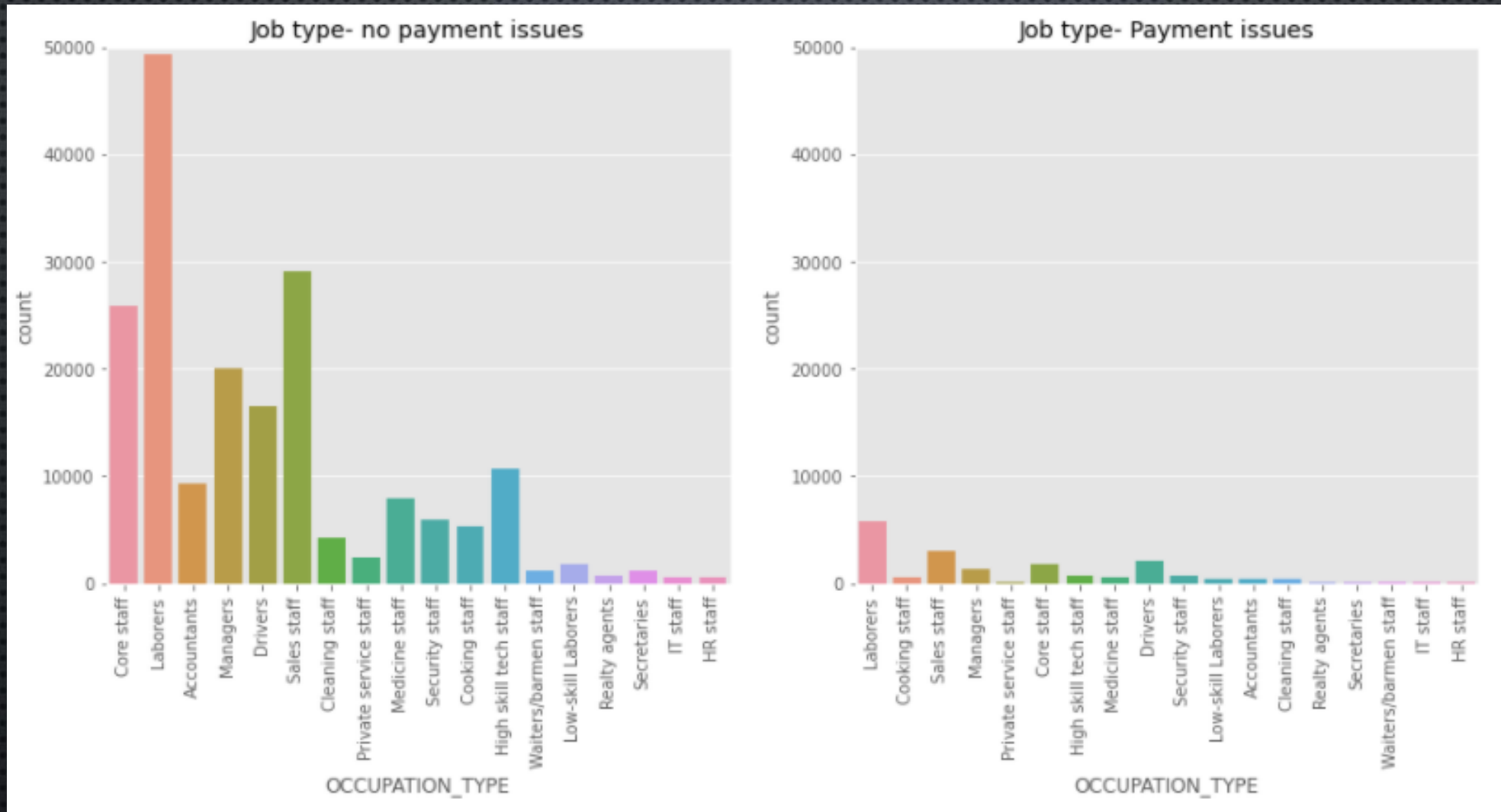


INFERENCES-

FEMALES ARE THE MAJORITY IN BOTH DEFAULTERS AND NON-DEFAULTERS. UNDOUBTEDLY, MOST FEMALES HAVE APPLIED FOR LOAN THAN MALES.

MALES ARE MORE IN DEFAULT THAN IN RATIO WITH MALES IN NON-DEFAULT.

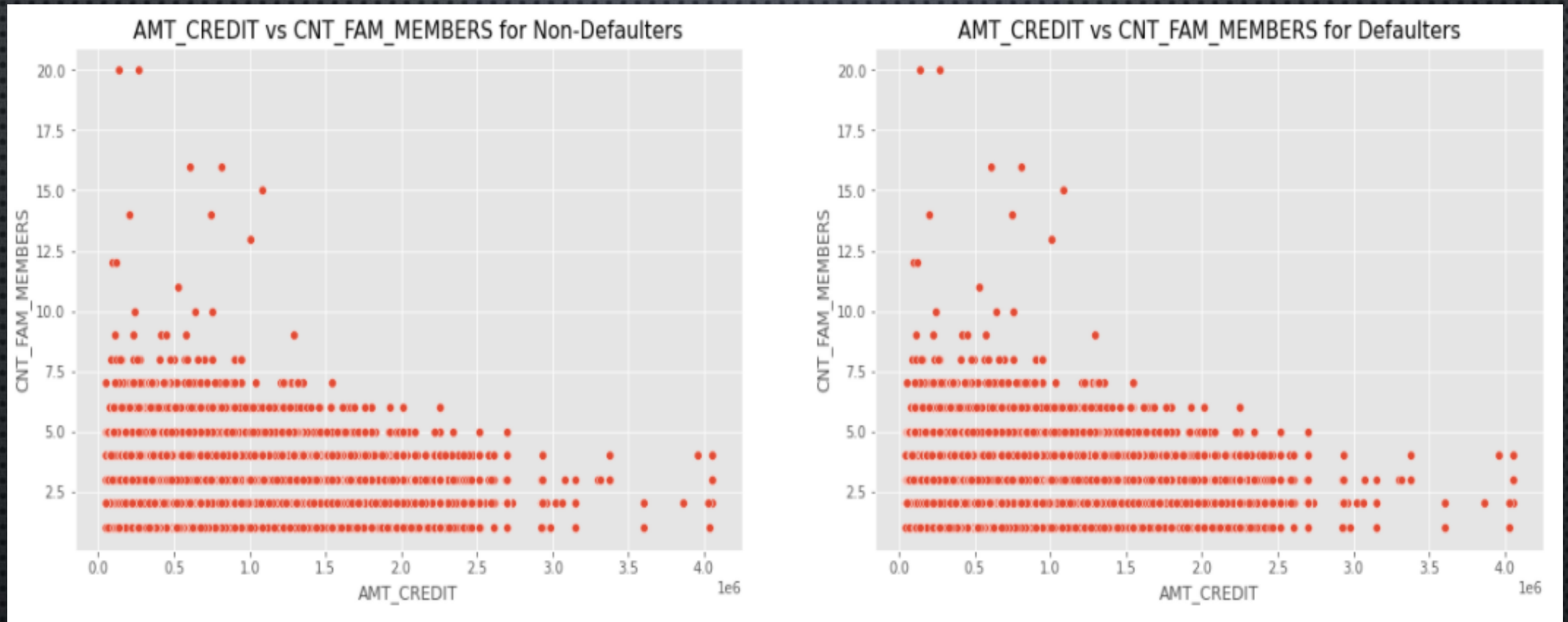
BIVARIATE ANALYSIS (BETWEEN TARGET AND OCCUPATION_TYPE)



INFERENCE-

FROM THE ABOVE GRAPH, WE COULD SEE THAT LABORERS ARE MOST LIKELY TO MAKE PAYMENT ON TIME WHEREAS HR STAFF ARE LESS LIKELY TO MAKE PAYMENT ON TIME.

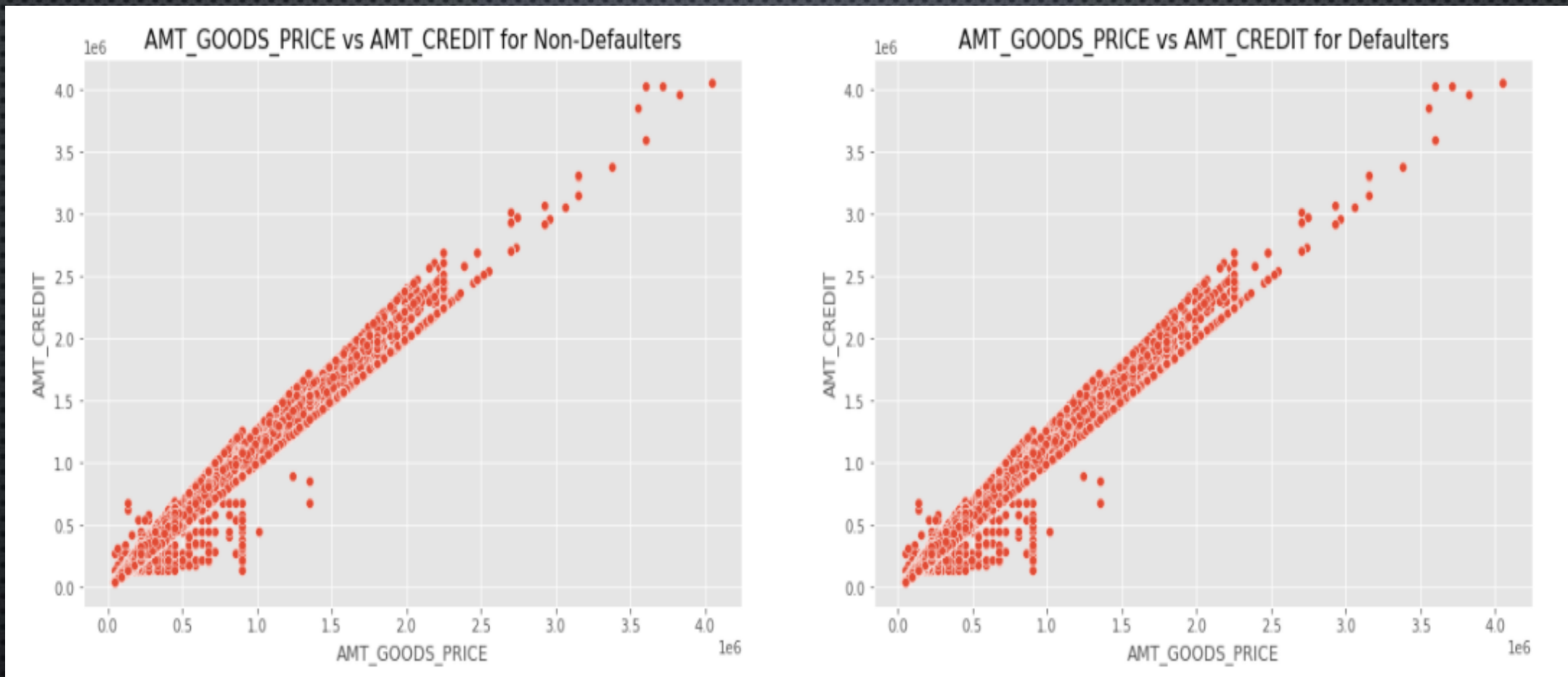
BIVARIATE ANALYSIS (WITH AMT_CREDIT & CNT_FAM_MEMBERS)



INFERENCES

- We can see that the density is in the lower left corner, similar in both the case. So, the people are equally likely to default if the family is small and the AMT_CREDIT is low. We can observe that larger families and people with larger AMT_CREDIT default less often.

BIVARIATE ANALYSIS (WITH AMT_GOODS_PRICE & AMT_CREDIT)



INFERENCE

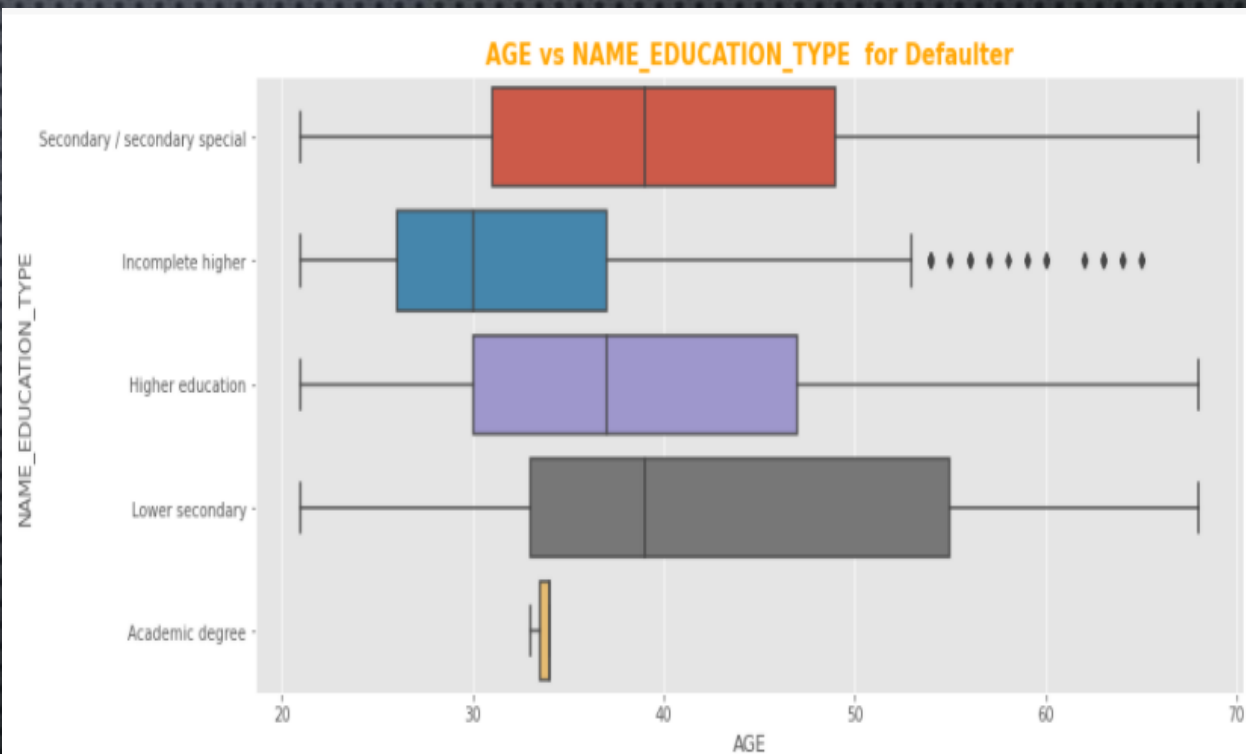
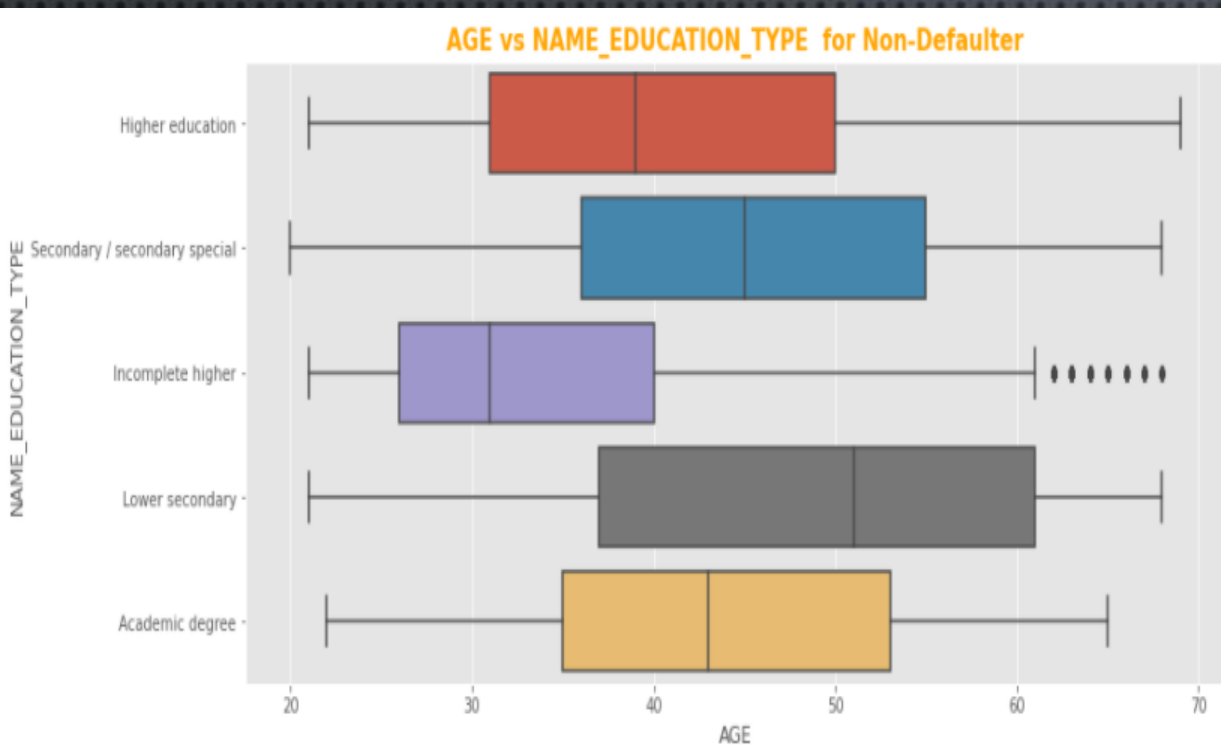
THE DEFAULTER SEEMS TO BE IN LESS NUMBER WHENEVER THE PRICE OF GOODS IS UPTO 50K AND CREDIT AMOUNT IS EVEN LESSER THAN 50K.

MULTIVARIATE ANALYSIS (FOR VARIOUS COLUMNS IN NON-DEFAULTER DATASET)



TOP3 CORRELATION IN GIVEN DATASET
('AMT_GOODS_PRICE', 'AMT_CREDIT')
('CNT_FAM_MEMBERS', 'CNT_CHILDREN')
('AMT_ANNUITY', 'AMT_CREDIT')

BOXPLOT BETWEEN AGE AND NAME_EDUCATION_TYPE (FOR DEFAULTER AND NON-DEFAULTER)



INFERENCES

- FOR NON-DEFAULTERS:

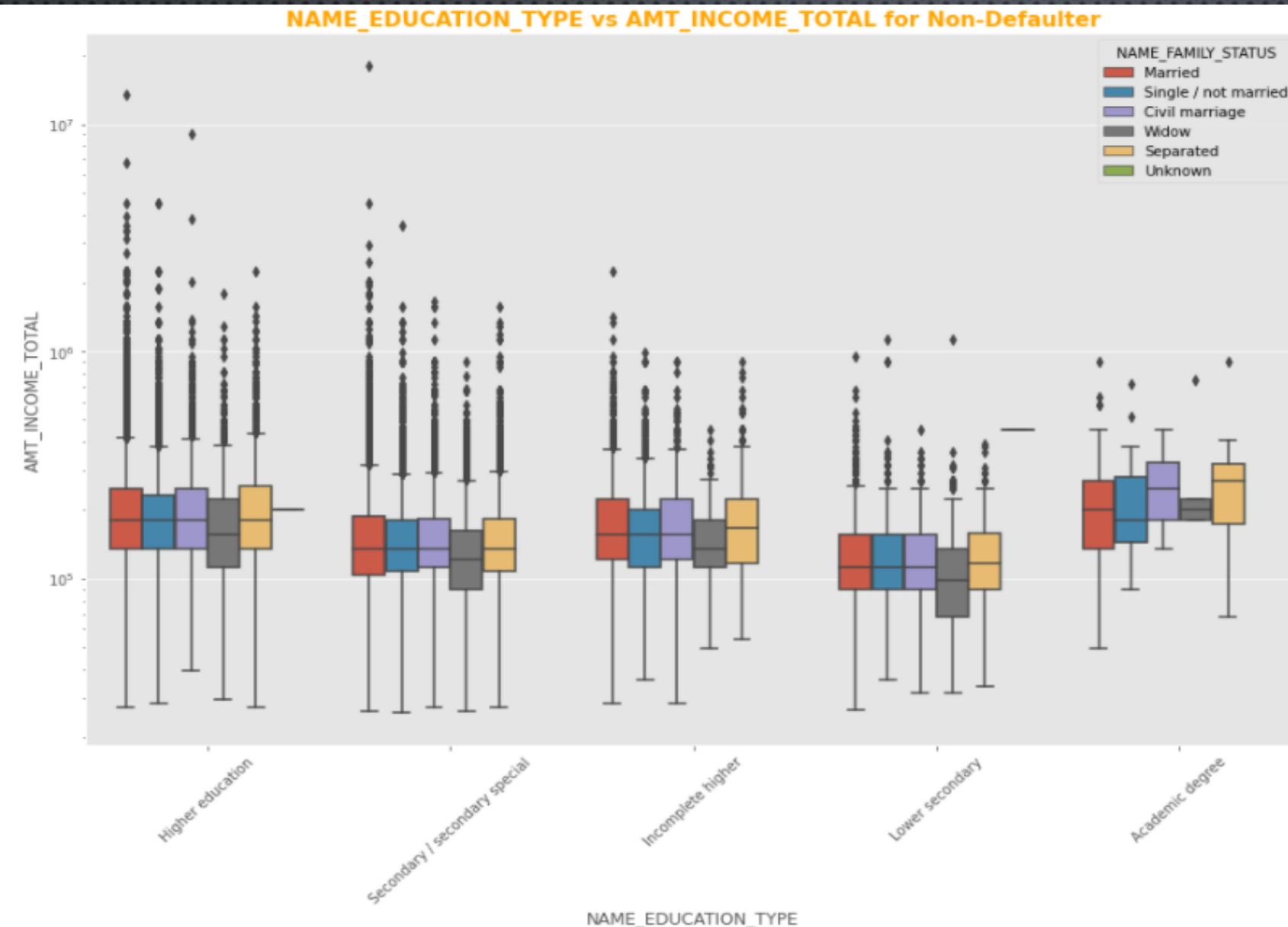
THERE IS AN OUTLIERS IN INCOMPLETE HIGHER IN BETWEEN AGE 60 TO 70. APART FROM INCOMPLETE HIGHER WHERE FIRST QUARTILE LIES ON AGE 40, OTHER EDUCATION TYPES AREN'T FACING MUCH DIFFICULTIES IN LOAN REPAYMENT.

- FOR DEFAULTERS:

THERE IS A VISIBILITY OF AN OUTLIERS IN INCOMPLETE HIGHER FROM THE AGE GROUP BETWEEN 50-70. PEOPLE WITH AN AGE GROUP 30-40 AND EDUCATION TYPE AS ACADEMIC DEGREE AND INCOMPLETE HIGHER SEEMS TO BE FACING DIFFICULTIES IN LOAN REPAYMENT.

MULTIVARIATE ANALYSIS

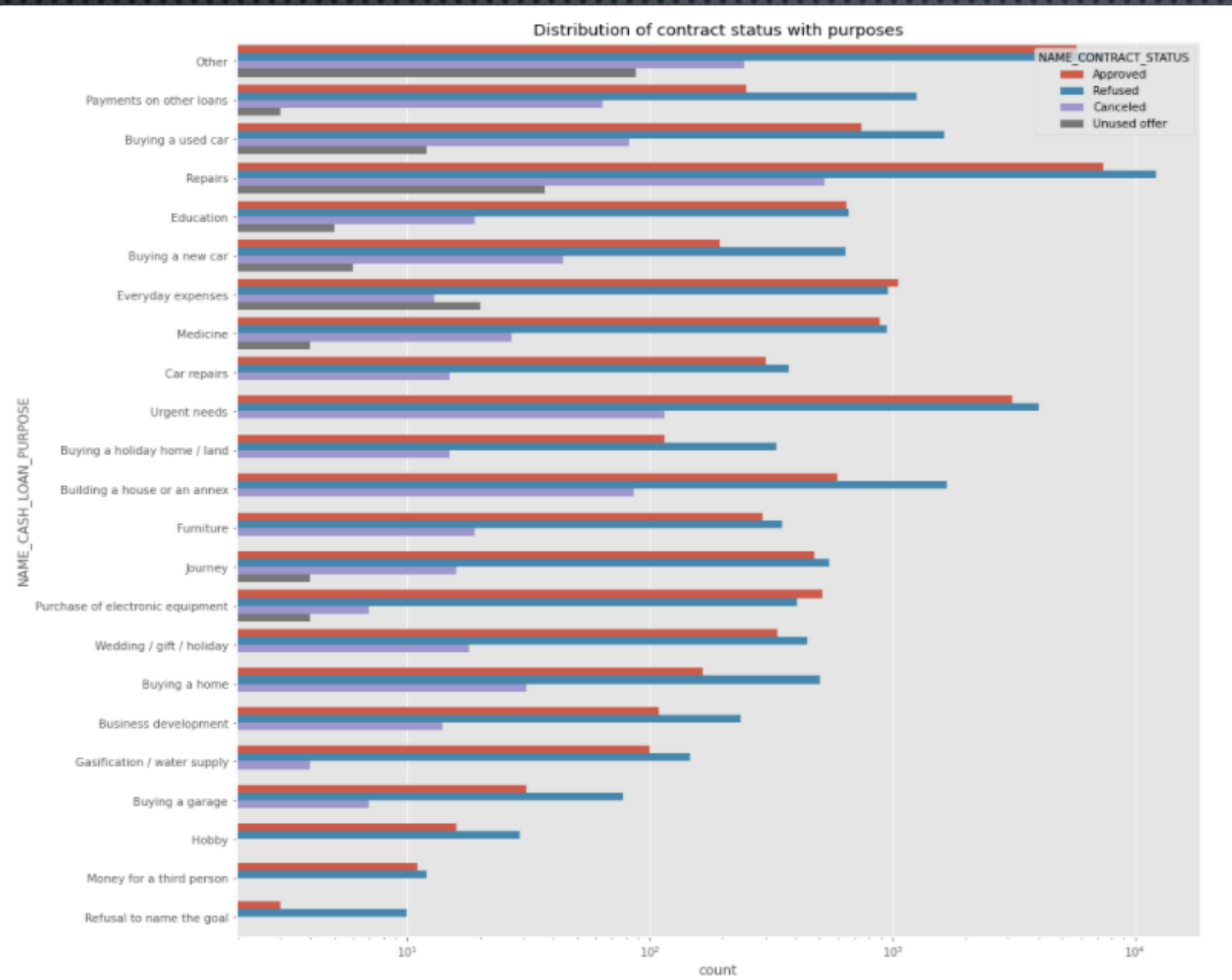
(B/W NAME_EDUCATION_TYPE & AMT_CREDI & NAME_FAMILY_STATUS FOR NON-DEFAULTER)



INFERENCES

- FOR EDUCATION 'HIGHER EDUCATION' TYPE THE INCOME AMOUNT IS MOSTLY EQUAL WITH FAMILY STATUS.
- HIGHER EDUCATION HAS MORE NUMBER OF OUTLIERS.
- ACADEMIC DEGREE TYPE HAVING LESS NUMBER OF OUTLIERS BUT THEIR INCOME AMOUNT IS LITTLE HIGHER THAN HIGHER EDUCATION.
- LOWER SECONDARY OF WIDOW FAMILY STATUS ARE HAVE LESS INCOME AMOUNT THAN OTHERS.

DISTRIBUTION OF CONTRACT STATUS WITH PURPOSE



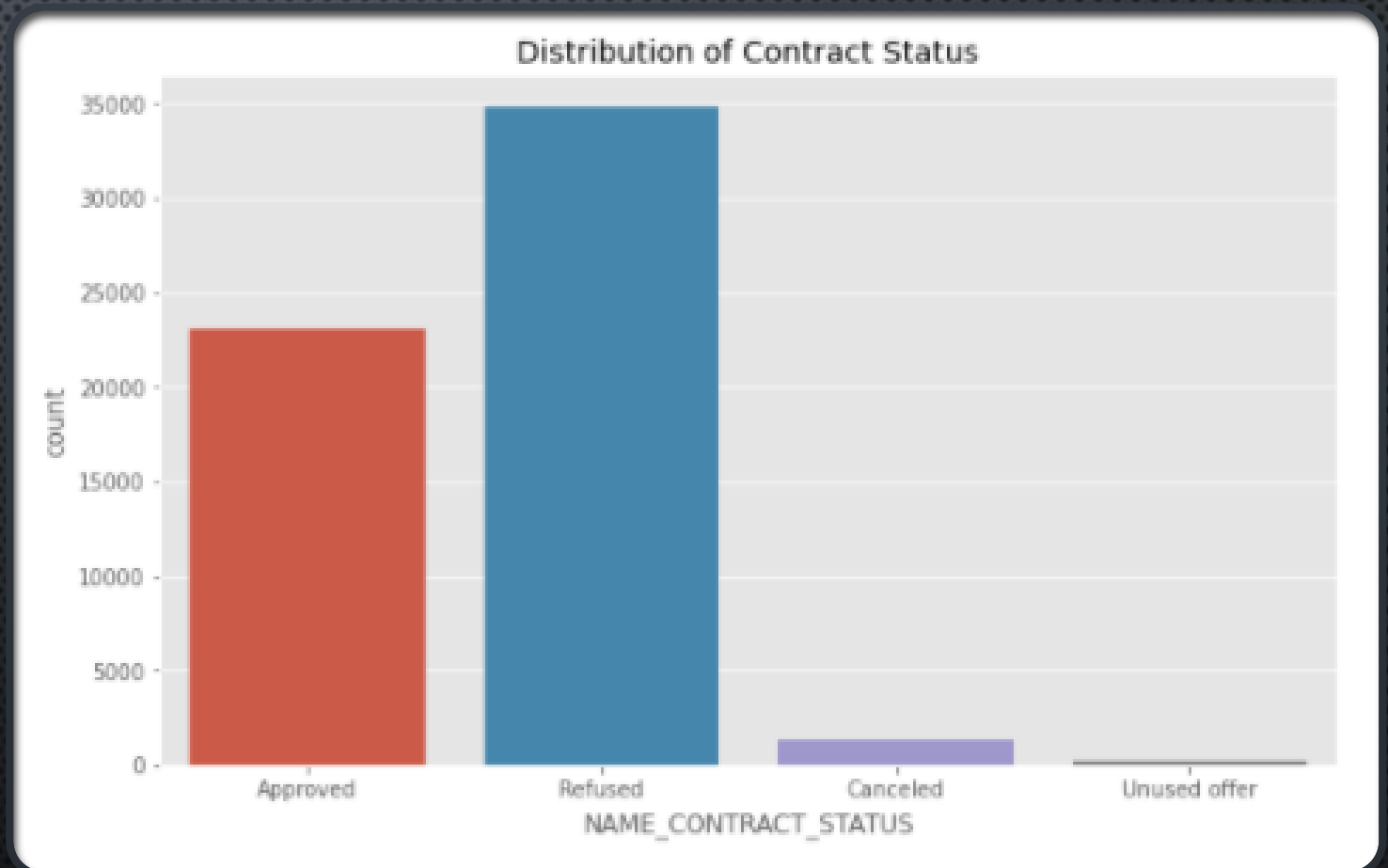
INFERENCES

- MOST OF THE REFUSED APPLICATIONS IS FROM REPAIRS.
- IN PURCHASE OF ELECTRONIC EQUIPMENT AND EVERYDAY EXPENSES, APPROVED APPLICATION STATUS IS MORE THAN REFUSED ONE.
- APPROVED AND REFUSED STATUS SEEMS TO BE ALMOST EQUAL IN EDUCATION.

DISTRIBUTION OF CONTRACT STATUS

INFERENCE

- FROM THE PLOT, REFUSED APPLICATION SEEMS TO BE MORE THAN OTHER STATUS OF APPLICATIONS.



AMT_APPLICATION VS AMT_INCOME_ORIGINAL



INFERENCES

- LOAN REQUEST HIGHER THAN 200K HAD A HIGHER REJECTION RATE.
- LOAN REJECTION RATE WAS MUCH LOWER IF THE INCOME WAS HIGHER.

CONCLUSION

1. BANK SHOULD FOCUS MORE ON STUDENT AND BUSINESSMAN FOR SUCCESSFUL PAYMENT.
2. BANK SHOULD FOCUS LESS ON WORKING TYPE CANDIDATES AS THEY ARE FAILING TO PAY DUES.
3. BANK SHOULD LOOK TO GIVE MORE REVOLVING LOANS AS THEY ARE SAFER.
4. BANK SHOULD AVOID GIVING LOANS TO THE HOUSING TYPE OF CO-OP APARTMENT AS THEY ARE HAVING DIFFICULTIES IN PAYMENT. BANK CAN FOCUS MOSTLY ON HOUSING TYPE HOUSE/APARTMENT, RENTED APARTMENT AND MUNICIPAL APARTMENT FOR SUCCESSFUL PAYMENTS.
5. REPAIRS SECTION IN LOAN PURPOSE SHOULD BE AVOIDED BY THE COMPANY. LOANS FOR BUYING LAND, NEW CAR, GARAGE AND NEW LAND SHOULD BE APPROVED MORE.
6. THE COMPANY SHOULD HAVE A CLOSER LOOK AT THE PEOPLE WHO ARE CHANGING THEIR ID DOCUMENTS RECENTLY AS THEY ARE MORE PRONE TO DEFAULT.

THANK YOU