# CREDIT EDA CASE STUDY

SUJITH M TOM
SHARIE R NATH

# AGENDA

- Give an idea of applying EDA in a real business scenario.

- Develop a basic understanding of risk analytics in banking and financial services

- Understand how data is used to minimise the risk of losing money while lending to customers

# Business Understanding

When the company receives a loan application, the company has to decide for loan approval based on the applicant's profile. Two types of risks are associated with the bank's decision:

- If the applicant is likely to repay the loan, then not approving the loan results in a loss of business to the company

- If the applicant is not likely to repay the loan, i.e. he/she is likely to default, then approving the loan may lead to a financial loss for the company.

# TARGET VARIABLE

The dataset contains the information about the loan application at the time of applying for the loan. It contains two types of scenarios:

- **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,

- **All other cases:** All other cases when the payment is paid on time.

# TYPES OF DECISION

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

- **Approved:** The Company has approved loan Application

- **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.

- **Refused:** The company had rejected the loan (because the client does not meet their requirements etc.).

- **Unused offer:** Loan has been cancelled by the client but on different stages of the process.

# PURPOSE OF CASE STUDY

- ■ Using EDA techniques analyse the dataset to understand how consumer attributes and loan attributes influence the tendency of default

# DATA UNDERSTANDING

- *'application_data.csv'* contains all the information of the client at the time of application.
  The data is about whether a **client has payment difficulties.**


- *'previous_application.csv'* contains information about the client's previous loan data. It contains the data whether the previous application had been **Approved, Cancelled, Refused or Unused offer.**

# PROBLEM STATEMENT

- Identify patterns which indicate if a client has difficulty paying their installments

- Ensuring consumers capable of repaying the loan are not rejected

- Find out client variables and loan variables that are high indicators of defaulting

# Exploratory Data Analysis Approach

The dataset is being analysed using following steps:

- Data Sourcing
- Data Cleansing
- Derived Metrics
- Univariate Analysis
- Segmented Univariate Analysis
- Bivariate Analysis
- Correlation Analysis
- Arriving  at Insights

# Data Sourcing

- *'application_data.csv' and 'previous_application.csv' datasets are meged on current application ID*

- *The combined dataset is used for analysis*

- *The columns are divided into numerical and catrgorical for ease of analysis*

# Data Cleansing

- Fixing Rows and Columns
- Dealing with missing values
- Detecting Outliers
- Finding Data Imbalance

# Fixing Rows and Columns

- Duplicate Rows and Columns are removed

- Certain columns are renamed

- Datatype of certain columns changed to category

- Columns segregated based on datatype as numerical,object and categorical

# Dealing with missing values

■ Null values in categorical columns replaced with appropriate values

■ Numerical Columns with more than 50% missing values identified

■ Insignificant columns are removed

■ RATE_INTEREST_PRIMARY and RATE_INTEREST_PRIVILEGED removed

List of Columns & NA counts where NA values are more than 50%

# Detecting Outliers

- Numerical columns are analysed for identifying outliers using **Box plots**

- Outlier Datapoints are detected using **IQR method**

- Datapoints that are beyond 1.5 times Inter Quartile Regions are considered as outliers

- Rows with Significant Outliers for particular variables are dropped

# Detected and Removed Outliers

SELLERPLACE_AREA :[4000000]

AMT_INCOME_TOTAL: [117000000.0, 18000090.0, 13500000.0]

# Box Plots for outliers

# Pseudo Outliers

Highly Significant outliers were found as follows:

DAYS_EMPLOYED:[ 365243]

DAYS_FIRST_DUE :[365243.0]

DAYS_LAST_DUE_1ST_VERSION:[365243.0]

DAYS_LAST_DUE :[365243.0]

DAYS_TERMINATION :[365243.0]

Since the value is spread over many rows, it is assumed as indication of unemployment or other category

# Data Imbalance

- Dataset is segmented into Repayers with **Target=0** and Defaulters with **Target=1**

- Countplot of Repayers and Defaulters indicates DATA IMBALANCE

- Count of Repayers = 1291326

- Count of Defaulters =122357

- Ratio of Data Imbalance is  1291326 :122357 ie. **10.554:1**

# Univariate Analysis

# TOTAL INCOME

# CURRENT CREDIT AMOUNT

# PREVIOUS CREDIT AMOUNT

# APPLICATION AMOUNT

# SELLER PLACE AREA

# CLIENT AGE
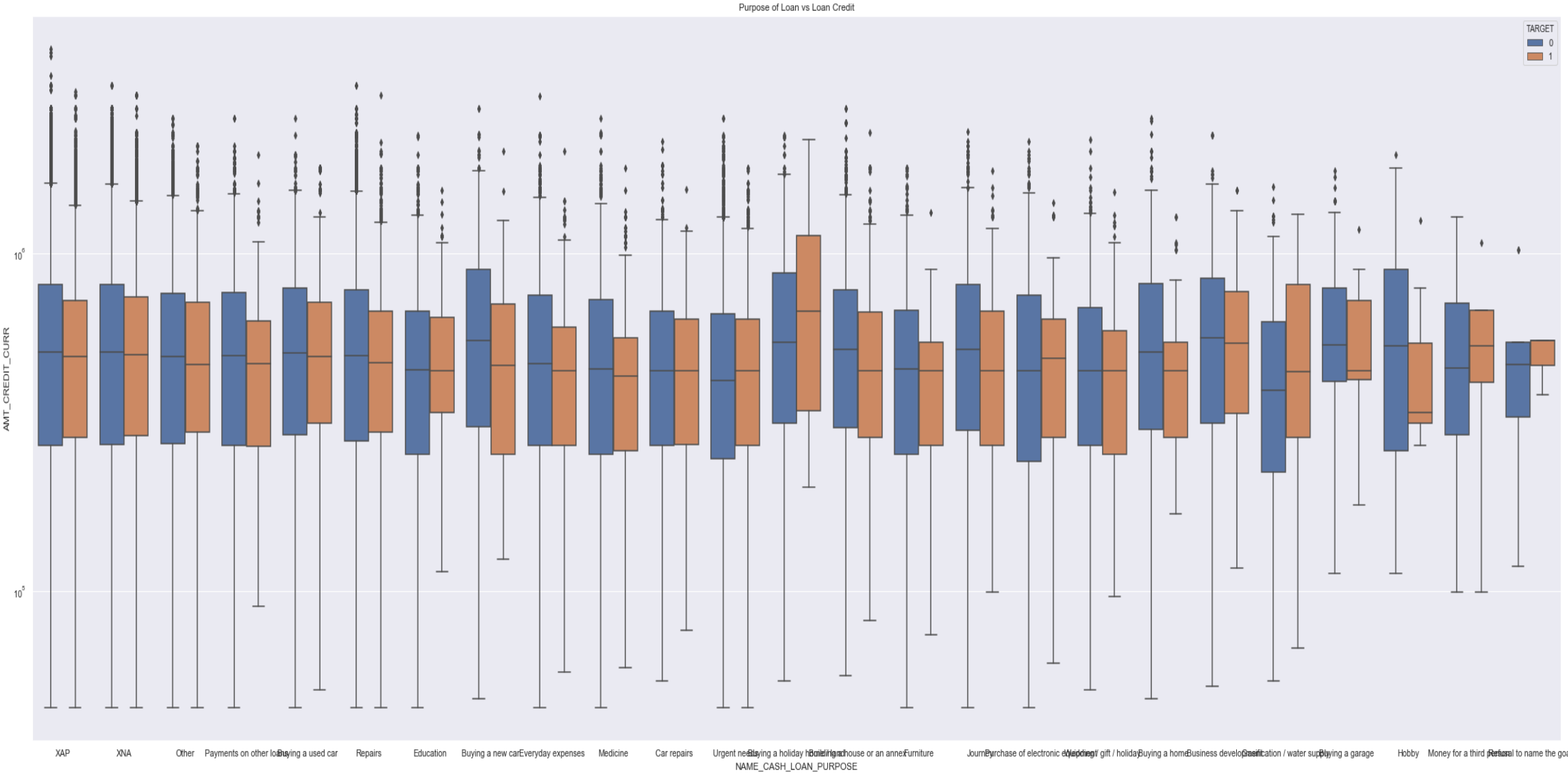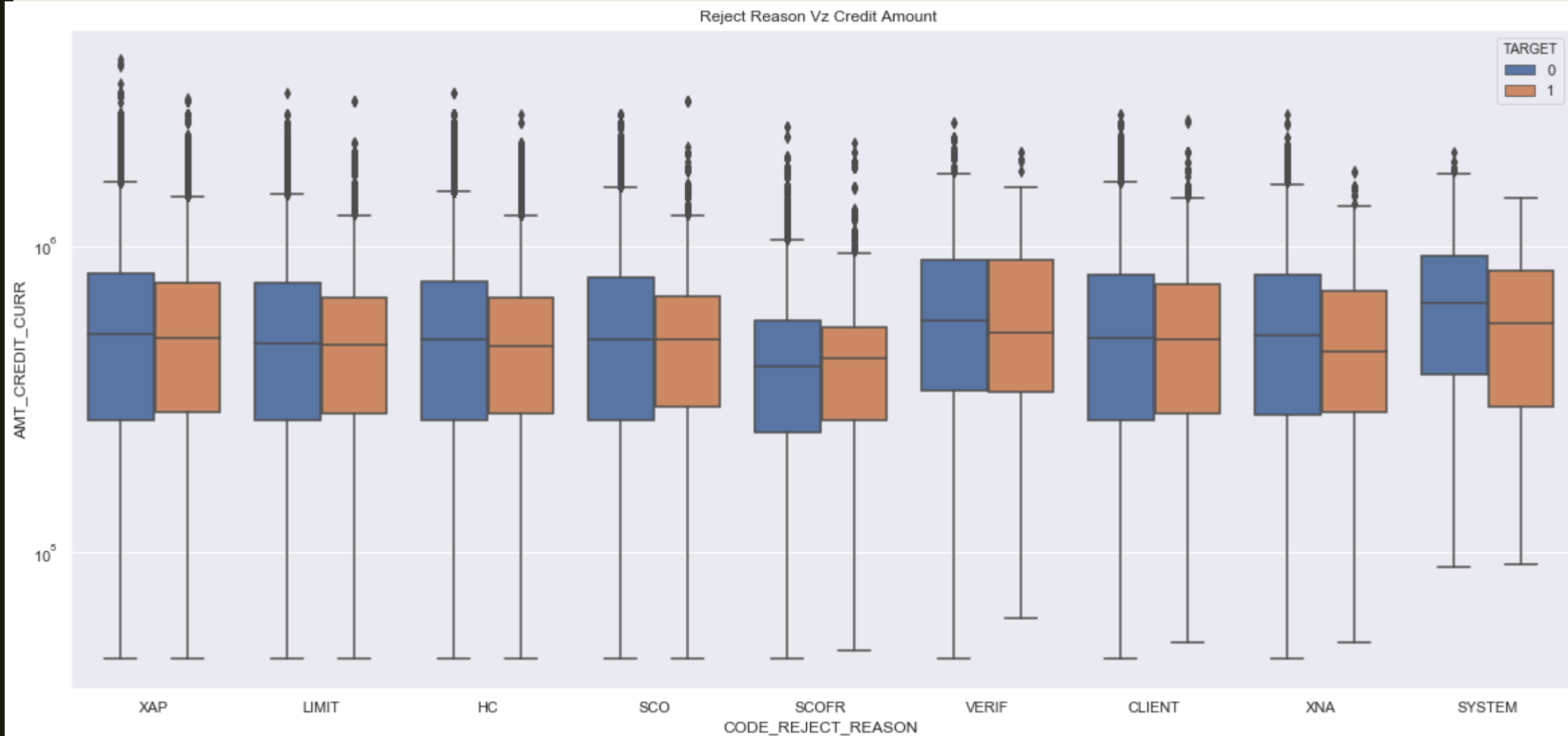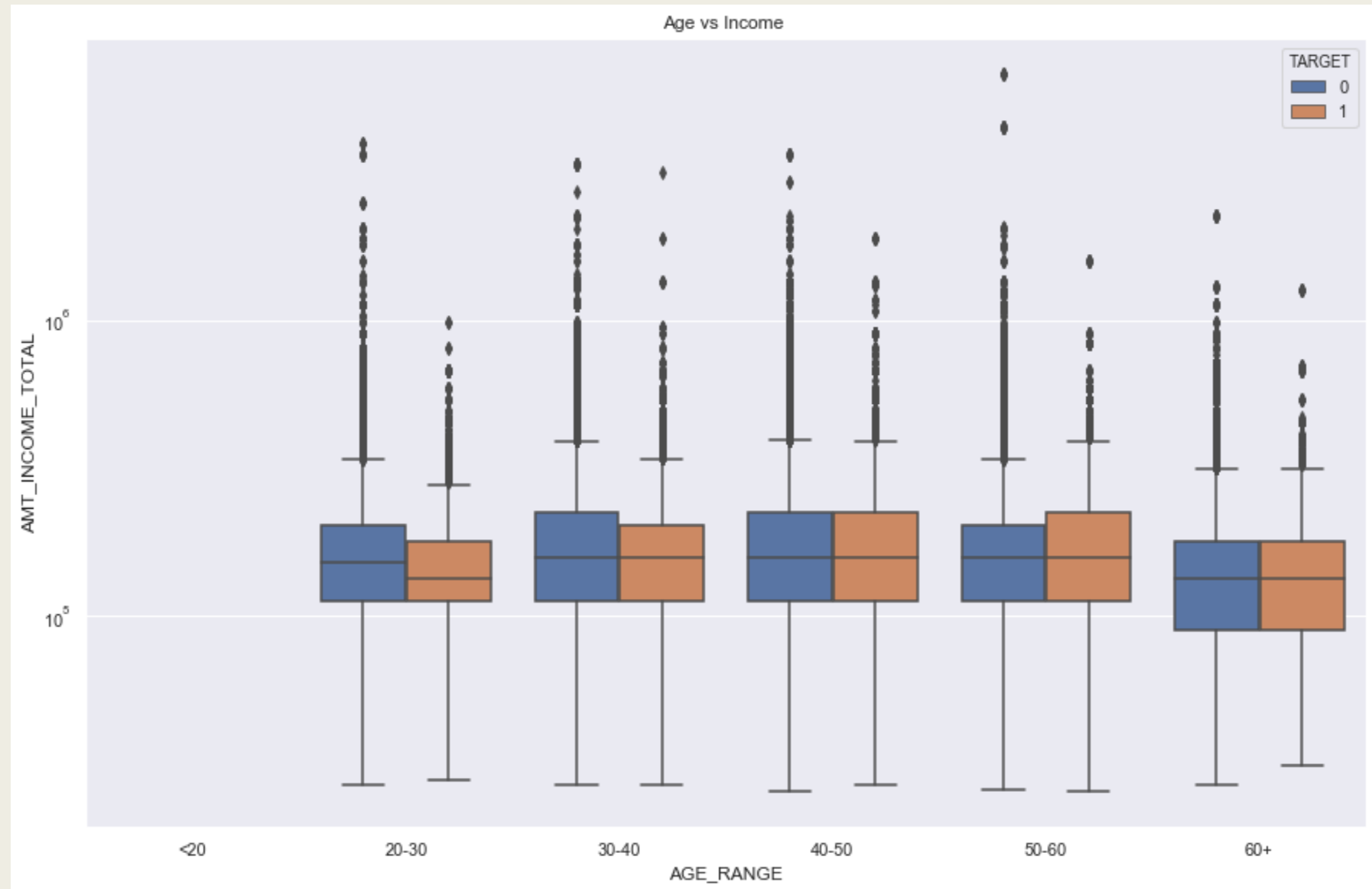
# Segmented Univariate Analysis

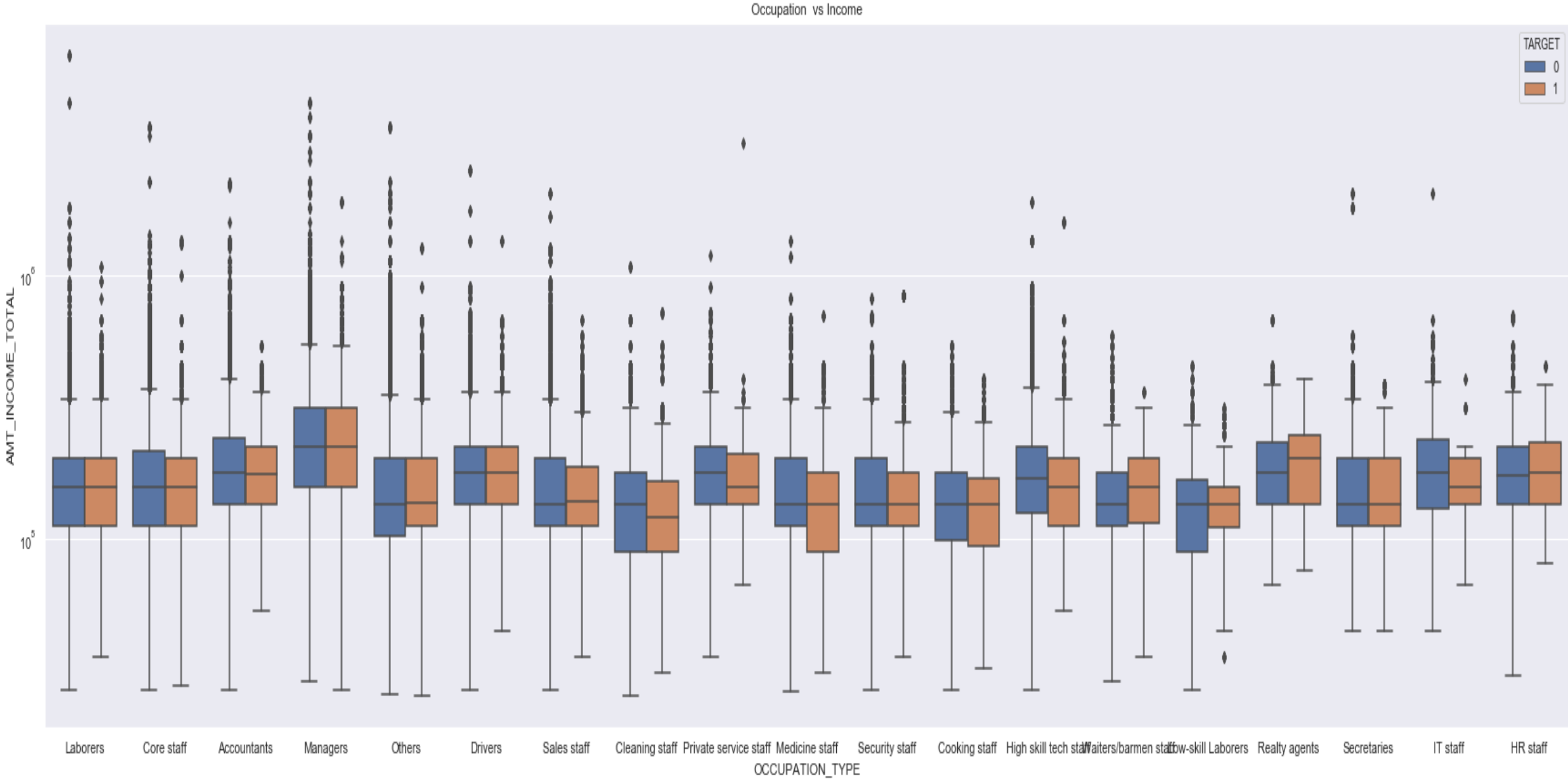# Bivariate Analysis

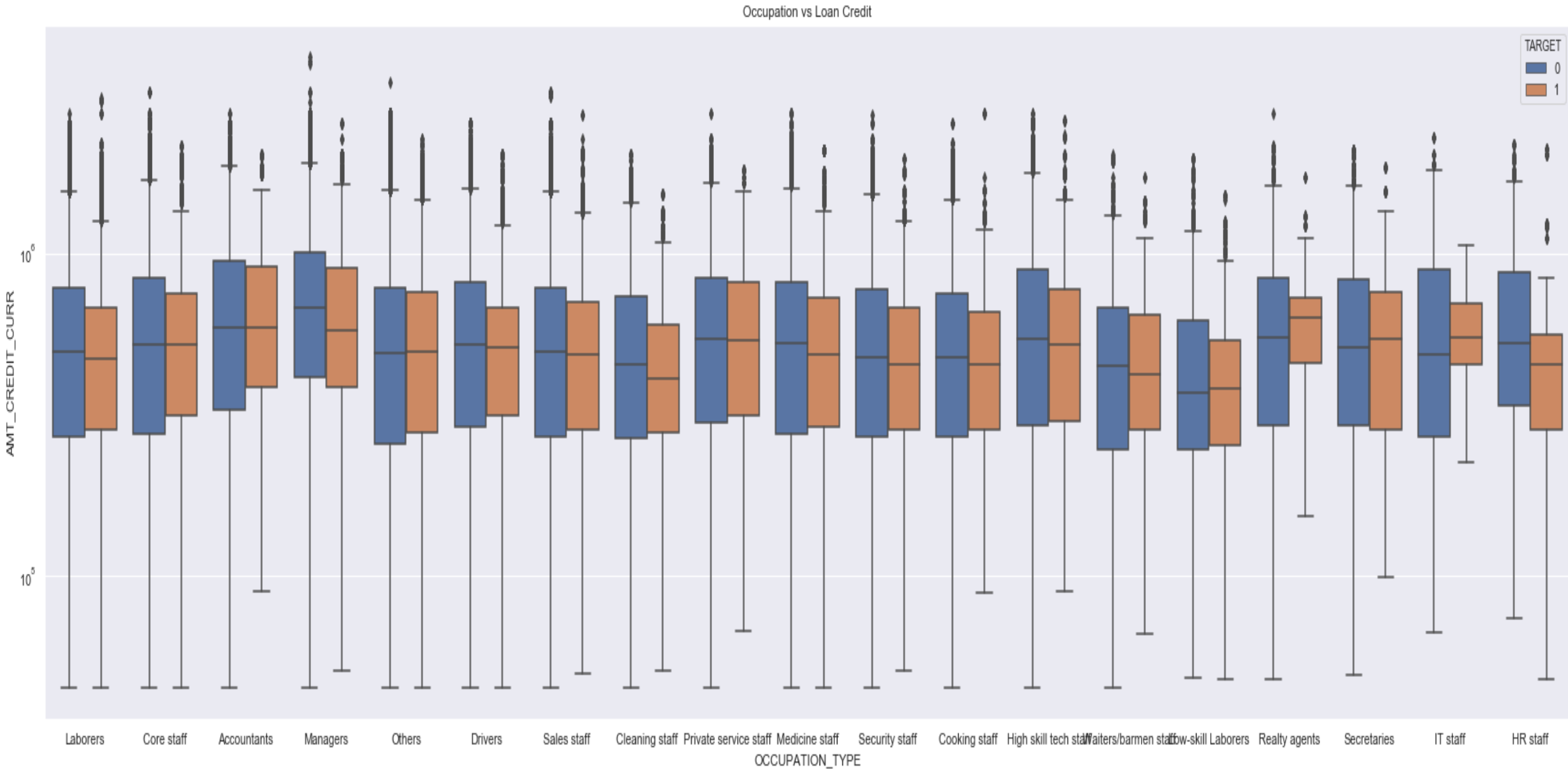# Purpose of Loan vs Loan Credit

# REJECT REASON VZ CREDIT AMOUNT



Reject Reason Vz Credit Amount

# AGE VZ INCOME AMOUNT



Age vs Income

# OCCUPATION VZ INCOME AMOUNT



Occupation vs Income

# OCCUPATION VZ LOAN AMOUNT



Occupation vs Loan Credit

# WORKING CITY VZ LOAN AMOUNT

# REPAYERS CORRELATION MATRIX-HEATMAP

# TOP 10 correlation for REPAYERS

| | | | |
|---|---|---|---|
| 1 | AMT_CREDIT_CURR | AMT_GOODS_PRICE_CURR | 0.9998 |
| 2 | AMT_APPLICATION | AMT_GOODS_PRICE_PREV | 0.9926 |
| 3 | YEARS_TERMINATION | YEARS_LAST_DUE | 0.9902 |
| 4 | AMT_CREDIT_PREV | AMT_GOODS_PRICE_PREV | 0.9902 |
| 5 | AMT_APPLICATION | AMT_CREDIT_PREV | 0.989 |
| 6 | AMT_ANNUITY_PREV | AMT_GOODS_PRICE_PREV | 0.964 |
| 7 | AMT_ANNUITY_PREV | AMT_CREDIT_PREV | 0.9584 |
| 8 | AMT_ANNUITY_PREV | AMT_APPLICATION | 0.9583 |
| 9 | CREDIT_TO_INCOME | ANNUITY_T0_INCOME | 0.9449 |
| 10 | AMT_ANNUITY_CURR | AMT_GOODS_PRICE_CURR | 0.944 |

# DEFAULTERS CORRELATION MATRIX-HEATMAP

# TOP 10 correlation for DEFAULTERS

| 1 | AMT_CREDIT_CURR | AMT_GOODS_PRICE_CURR | 0.9998 |
|---|---|---|---|
| 2 | YEARS_TERMINATION | YEARS_LAST_DUE | 0.9942 |
| 3 | AMT_APPLICATION | AMT_GOODS_PRICE_PREV | 0.993 |
| 4 | AMT_APPLICATION | AMT_CREDIT_PREV | 0.989 |
| 5 | AMT_CREDIT_PREV | AMT_GOODS_PRICE_PREV | 0.989 |
| 6 | AMT_ANNUITY_PREV | AMT_GOODS_PRICE_PREV | 0.969 |
| 7 | YEARS_FIRST_DRAWING | YEARS_LAST_DUE_1ST_VERSION | 0.9688 |
| 8 | AMT_ANNUITY_PREV | AMT_CREDIT_PREV | 0.9665 |
| 9 | AMT_ANNUITY_PREV | AMT_APPLICATION | 0.9646 |
| 10 | CREDIT_TO_INCOME | ANNUITY_T0_INCOME | 0.9508 |

# Business Insights

- Majority Loanees are more in Region_Rating_Client=2
- Majority Loanees are more in Region_Rating_Client_W_City=2
- No one loanee without mobile number
- Majority Loanees both repayers and defaulters  are married and working with secondary education
- Majority of the loanee are employed within 2-5 years
- Majority of loanees are in the age range 30-40
- Defaulters are very less in the age range of 60+
- Majority of loanees in the range of income 1LA-2LA
- Majority of loanees are of organization type -Busniess Entity
- Majority of loanees are getting credit range 4LA-10LA and Annuity range 20k-30K
- Majority of loanees have current good price in the range 2la-6la
- Majority of loanees  application amount is between 50k-4LA
- Majority of loanees Down payment in the range 2,5k-10K

- Majority of the repayers last phone change is 5+ years before application

- Majority of the defaulters last phone change is 1-2 years before application

- Defaulters are more in Region_Rating_Client_W_city=2

- Defaulters occupation is majorly Others and Laborers

- Defaulters are very less in the age range of 60+

- Deafulting high in Cash Loans and Approved (Contract Status),and cash through bank payment gtype

- Defaulting tendency more in Client type : Repeaters

- Defaulers are very less in IT SECTOR

- Defaulters are more on high yield group than middle and lower

- Defaulters are more in Loan Credited for the purpose: Buying a Holiday Home

- Loan with SCOFR reject reason defaulted less

- Main code reject reason is XAP then, HC

# THANK YOU