# ☐ **Problem Statement**

- It is important to note that this data carries two risks

1. The Bank can suffer if the customer is unable to handle the load and default it
2. BANK will lose money if it does not give the loan to the customer who is able to repay it.

# ❑ Assumptions

- I am getting only one approximation in this given set, as I checked

- According to my observation, XNA is present in each of these columns (ORGANIZATION_TYPE, NAME_GOODS_CATEGORY, NAME_SELLER_INDUSTRY, NAME_CASH_LOAN_PURPOSE, NAME_PORTFOLIO, NAME_PAYMENT_TYPE, CODE_GENDER)

# ❑ **Data Cleaning & Fill Missing Value**

- First of all in EDA I cleaned the null values and get its percentage, and removed from that which were more than 35% NULL values, because if we fill that value, we may have problem in visualization.

- After that, I removed the low percentage NULL value (less than 1%) because we have a lot of data,

- I filled the mean value in place of NA, which were column integers, floating, and I filled the mode value in the categorical (object) column.

- I removed outliers with the help of boxplot

- In standardizing numerical values, I checked the datatype of the column, all the columns were correct, but many columns had negative values but they should have been positive, I converted it into positive.
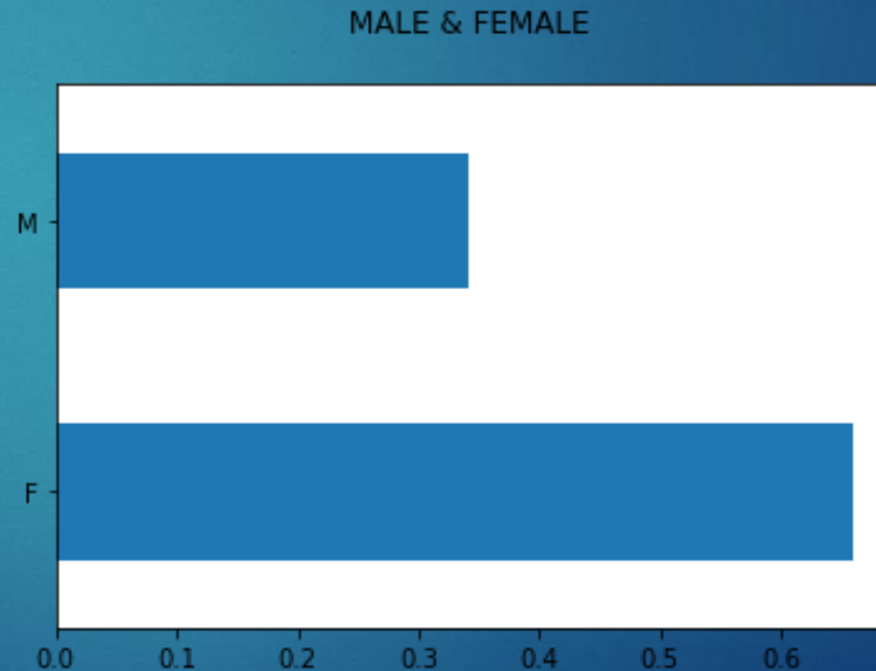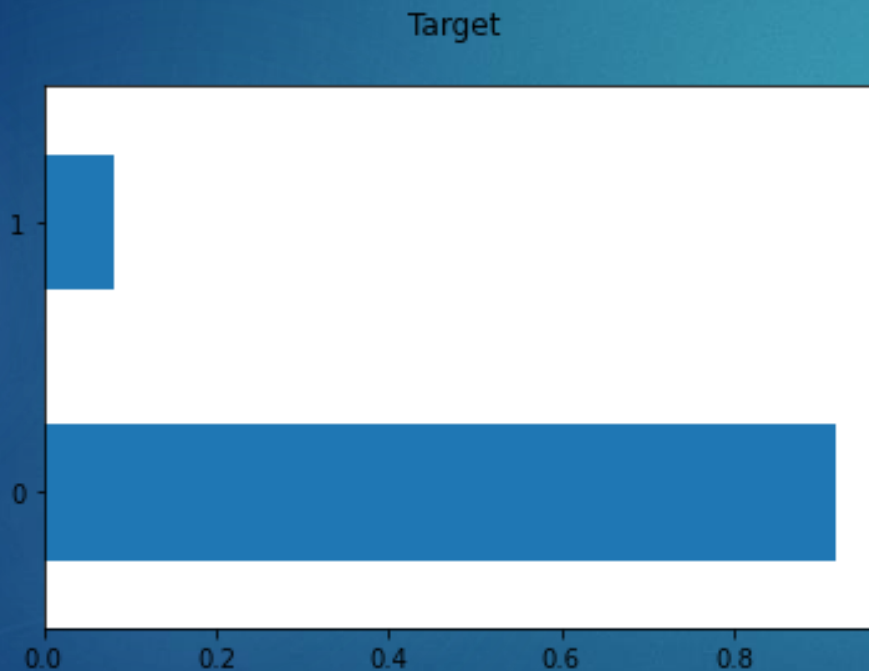
# ❑ Data Balancing

- I checked the data balance, In the data I check data describe 25%, 50%, and 75%.

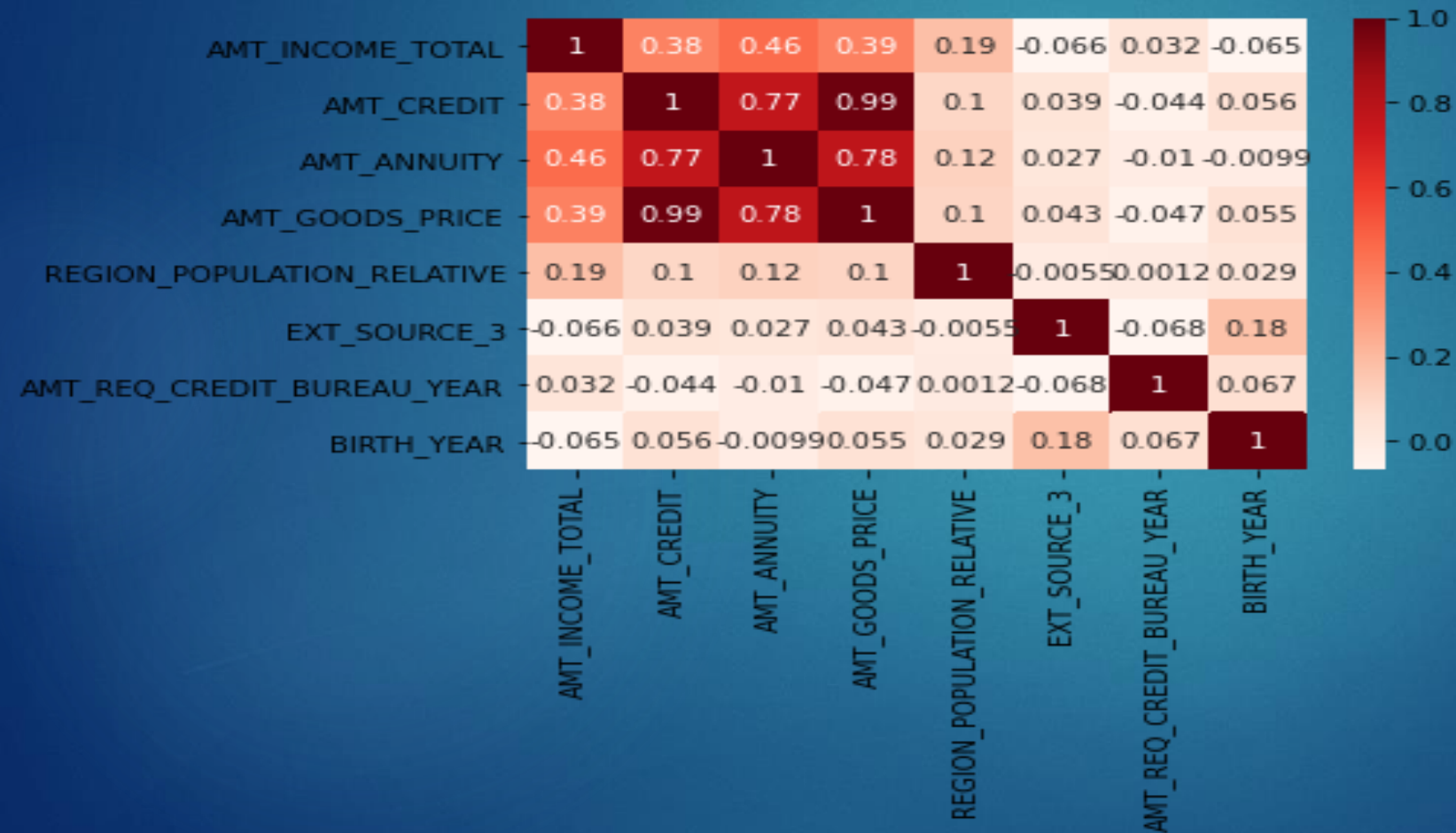- In order to make sure the columns were balanced, I removed them

# ❑ Analysis of Variables

- A Barplot analysis has been performed on all columns

# ❑ inside

- I checked with correlation plot how many columns are more related to each other
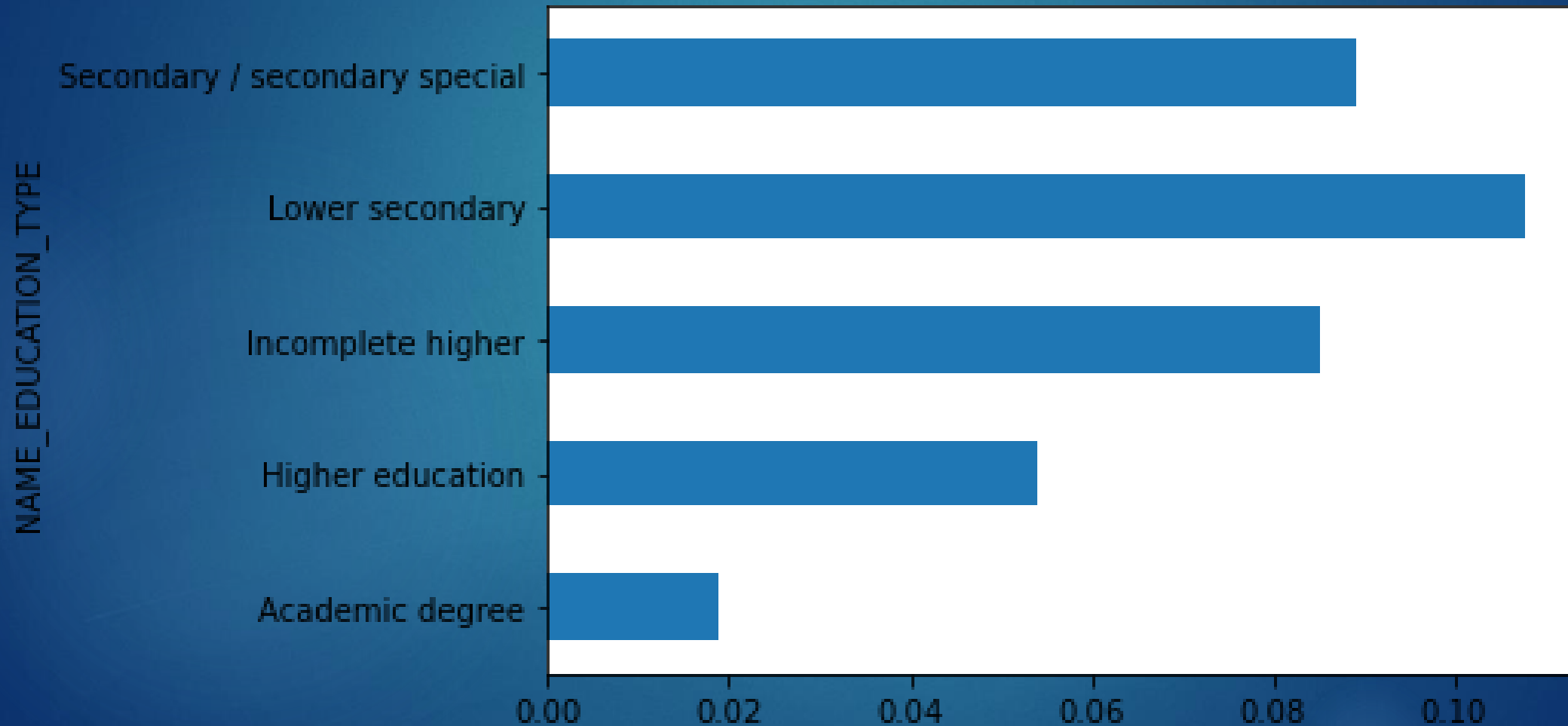
# ❑ inside

- As a group, I formed a column for education, and another for those who were unable to pay their loans, and plotted mean values on a bar graph

- And in that plot I have come to know that the customers with lower secondary education are not able to pay the loan much.
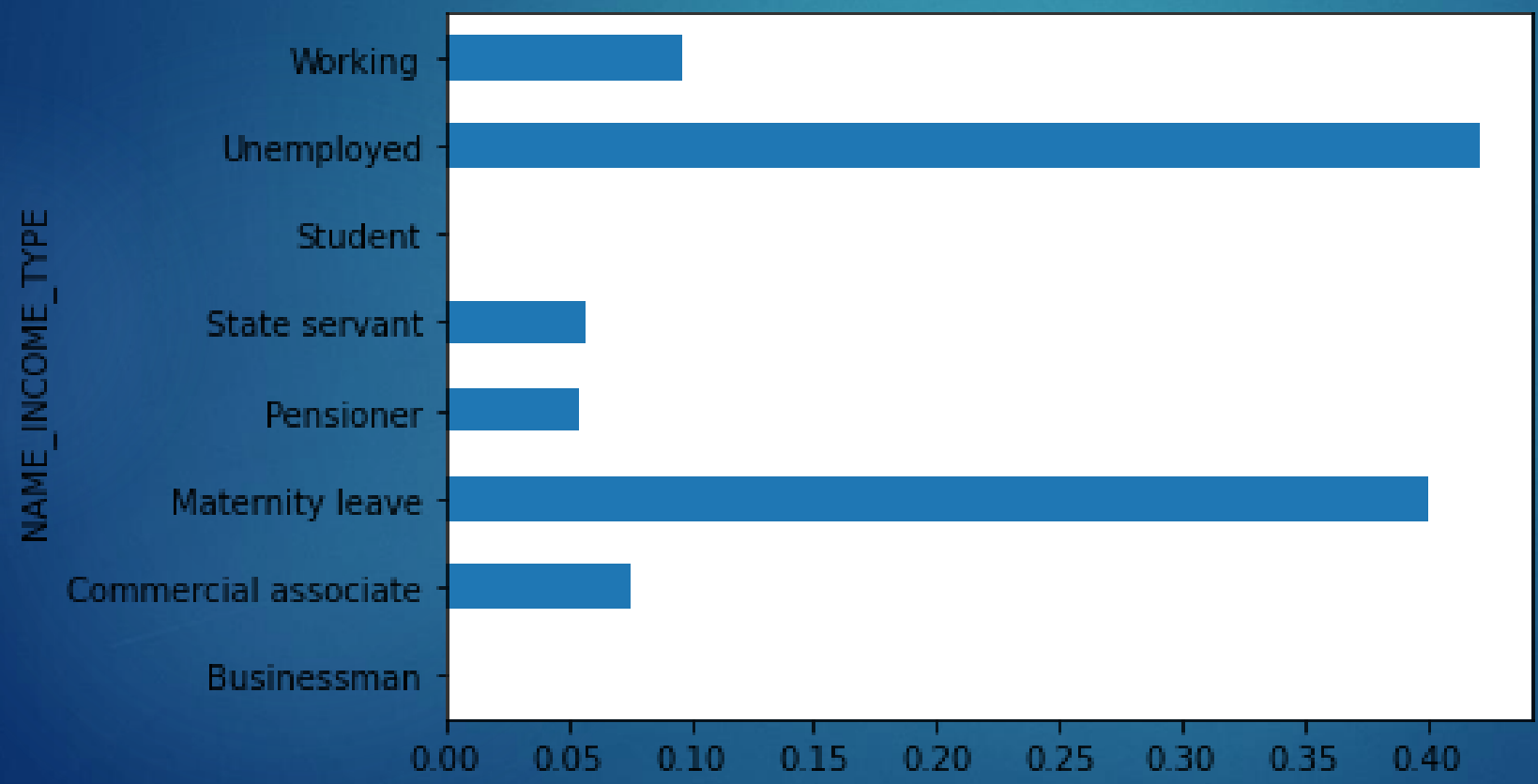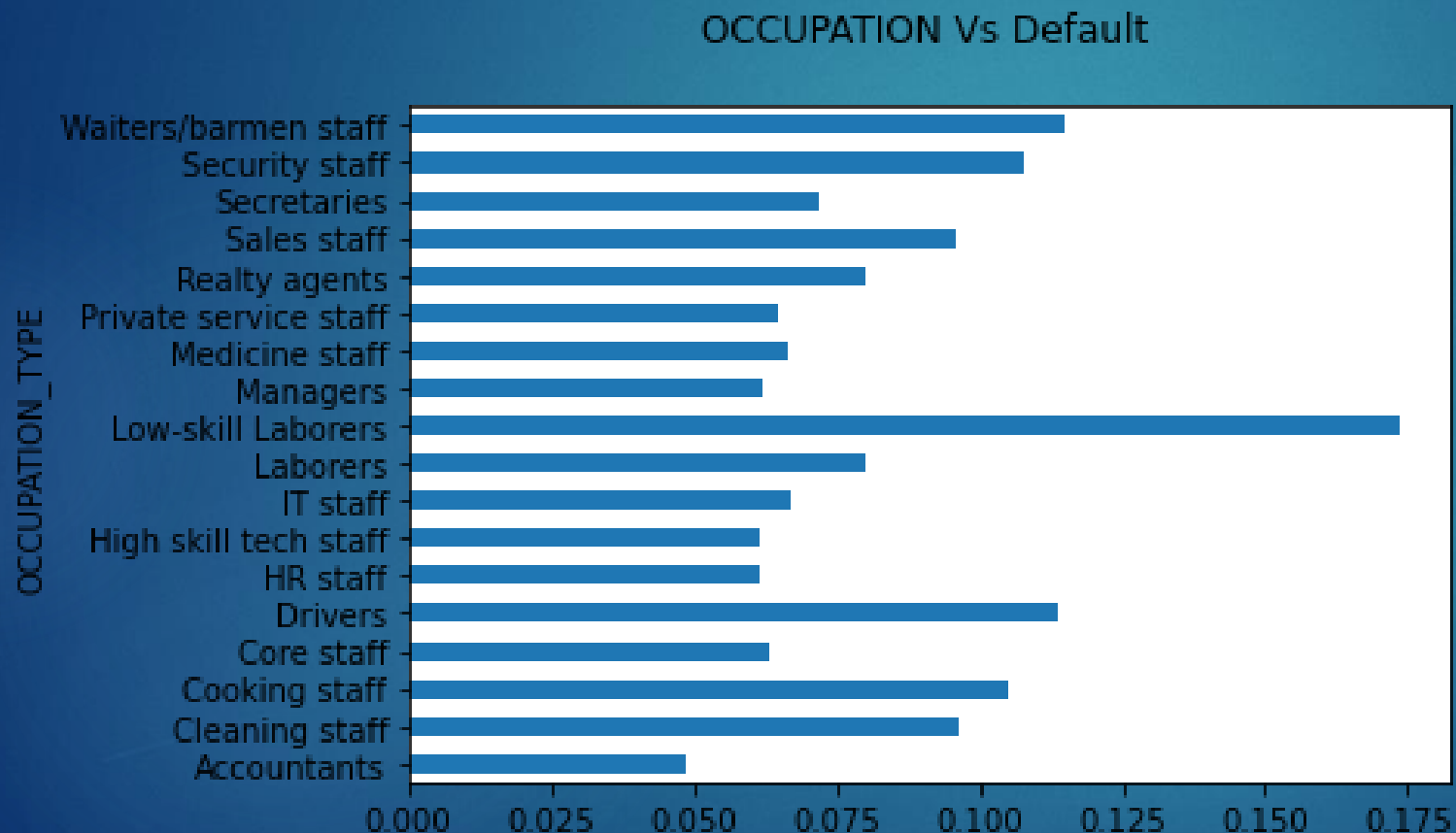
Similarly, I plotted the income type and default customer, and I got the higher rate Maternity leave and Unemployed, unable to pay.
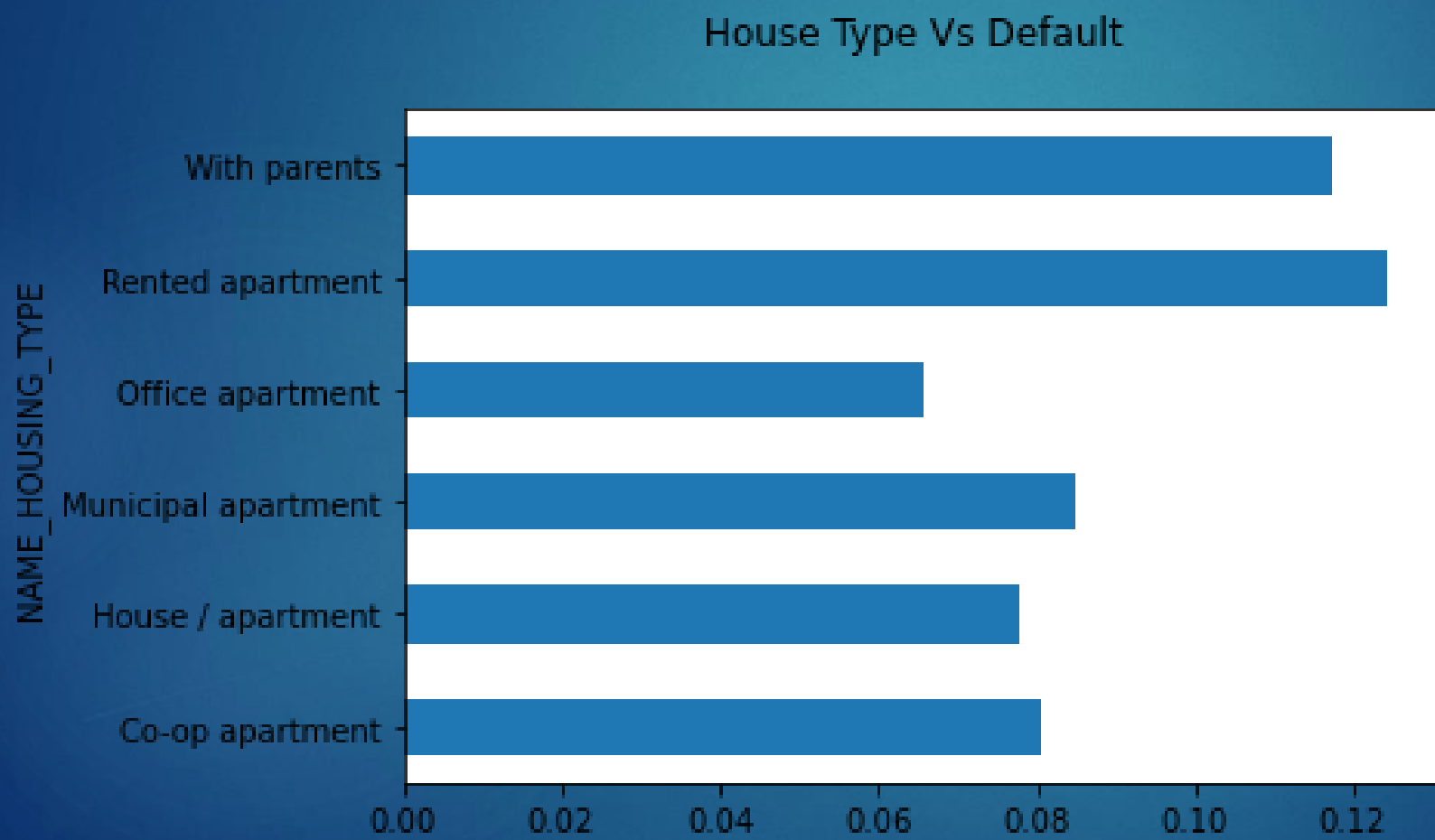


INCOME_TYPE Vs Default

- Similarly, I plotted the OCCUPATION TYPE and default customer, and I got the higher rate Low-skill Laborers, unable to pay.



OCCUPATION Vs Default

- Similarly, I plotted the HOUSING TYPE and default customer, and I got the higher rate in Rented apartment and With parents, unable to pay.


House Type Vs Default

First I split the age column in the difference bucket, and Similarly, I plotted the Age group and default customer, and I got the higher rate in between 20 to 30, unable to pay.



Age Groups Vs Default