

Module 8: Credit EDA Assignment

# Credit EDA Case Study

An Entry into Data Analysis

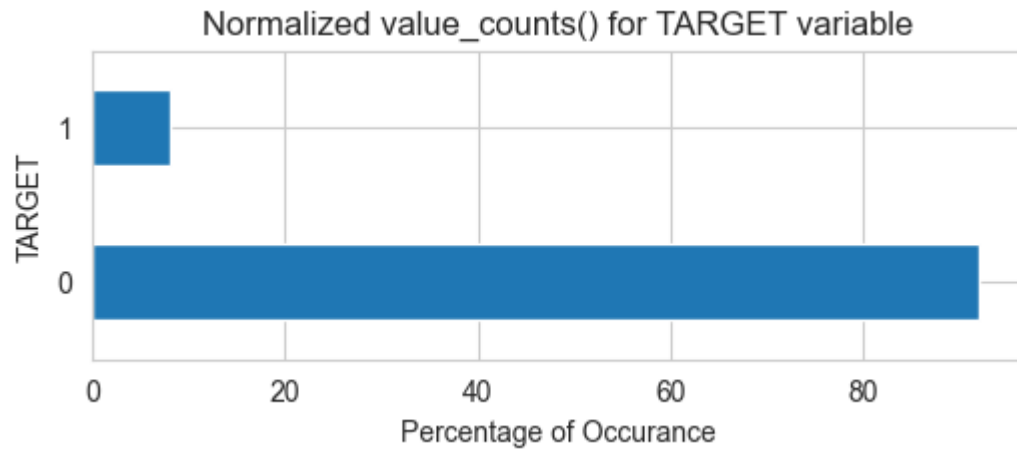
Submitted by: **SHUBHAM B. KANHEKAR**

Batch : DSC73 (Oct'24)

Date: 31/12/2024

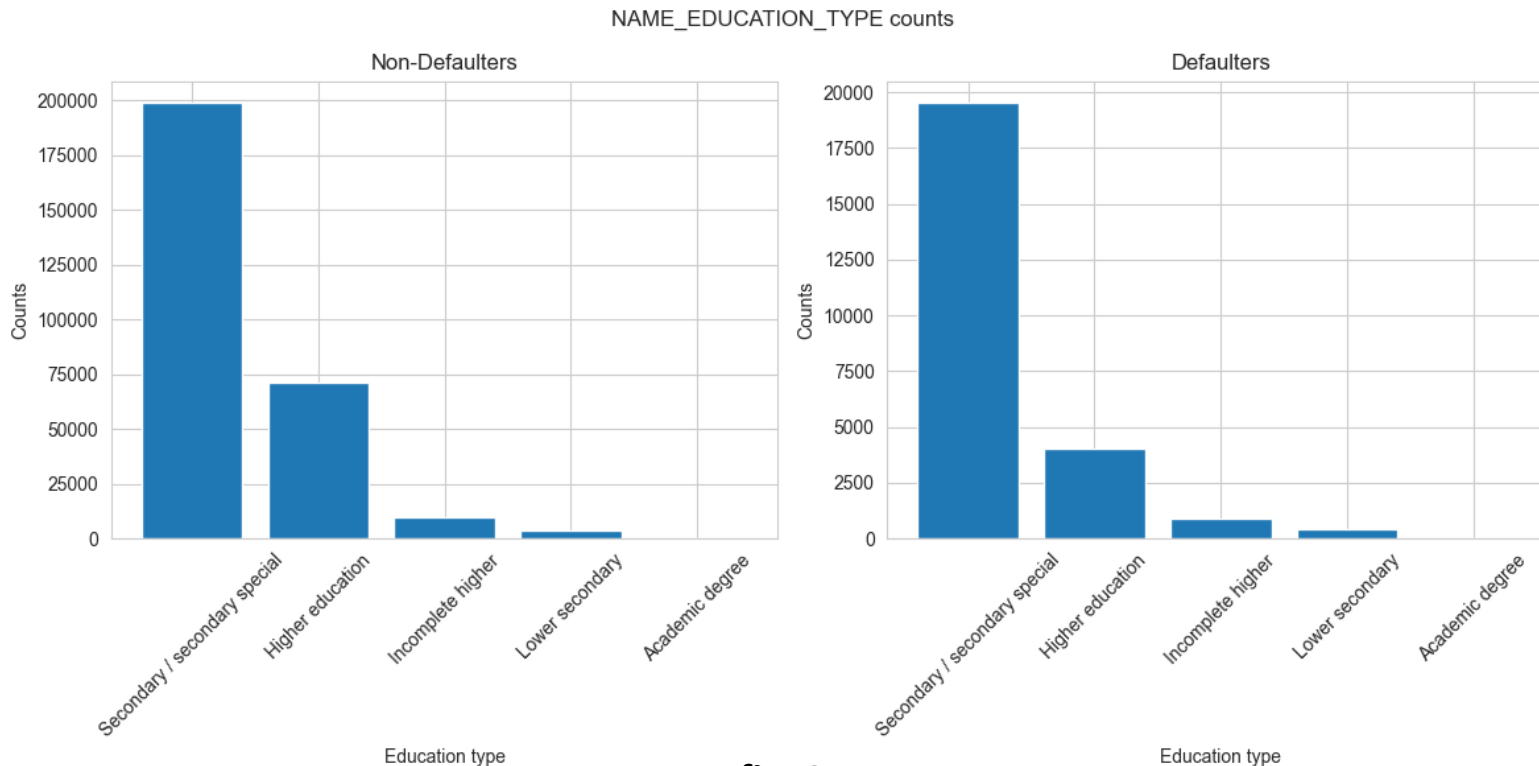
# Goal of the Analysis

The goal of the Credit EDA Assignment is to get us acquainted with various cleaning and analysis techniques while also gathering Insights from the dataset provided to us. In this presentation, we will be looking at some of the many plots created during the analysis and a brief note about each plot has be added with the figures. Also, the recommended variables which should be observed by the financial institute/bank are listed in the second to last slide.



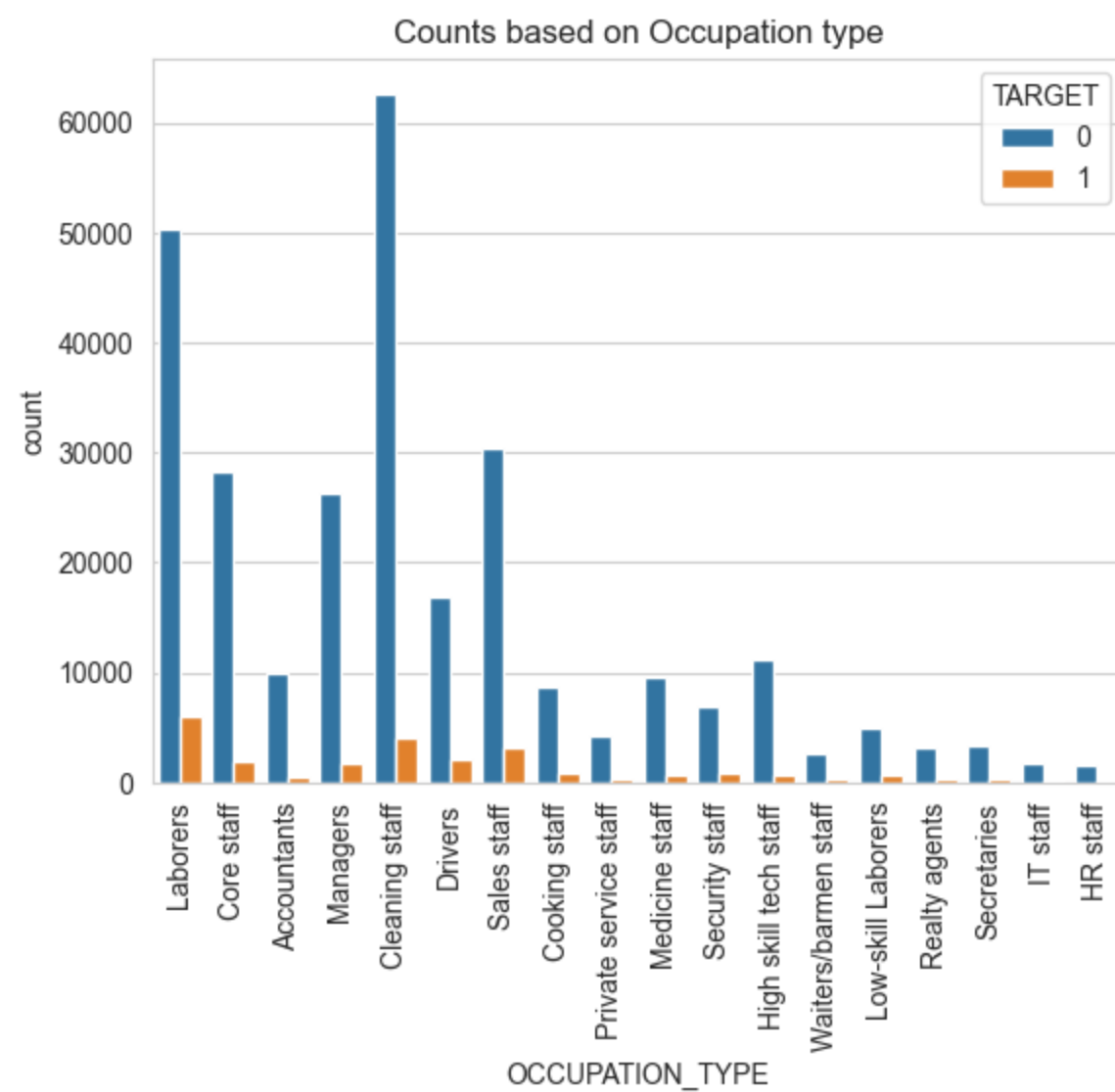
**fig\_1**

fig\_1: 'TARGET' variable value counts show that 8.07% of applicants have trouble with payments

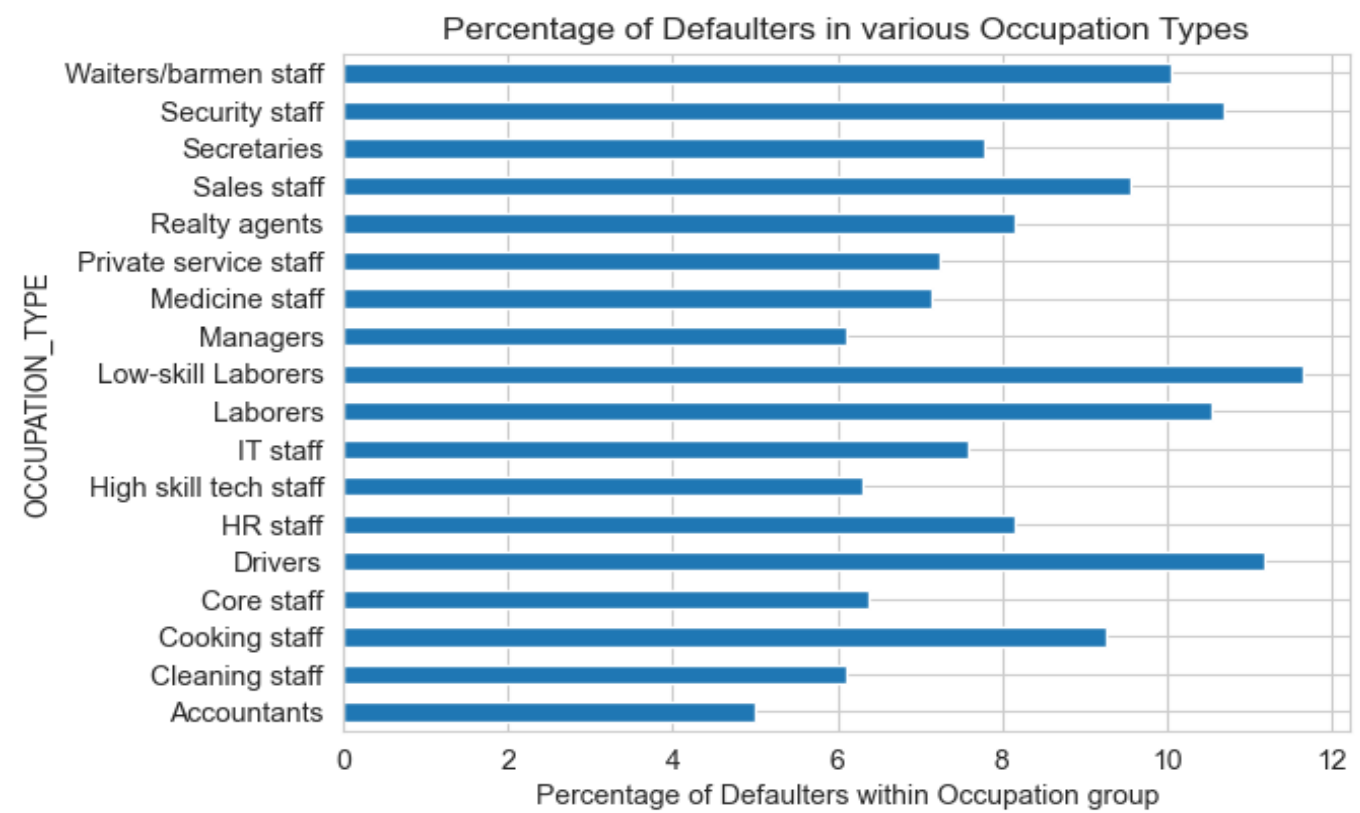


**fig\_2**

fig\_2: 'Education type: Higher education' shows the most promising clients who are least probable to default on payments as we can see in fig\_2; All bars are same height (proportion) except 'Higher Education' has lower proportion in defaulters plot

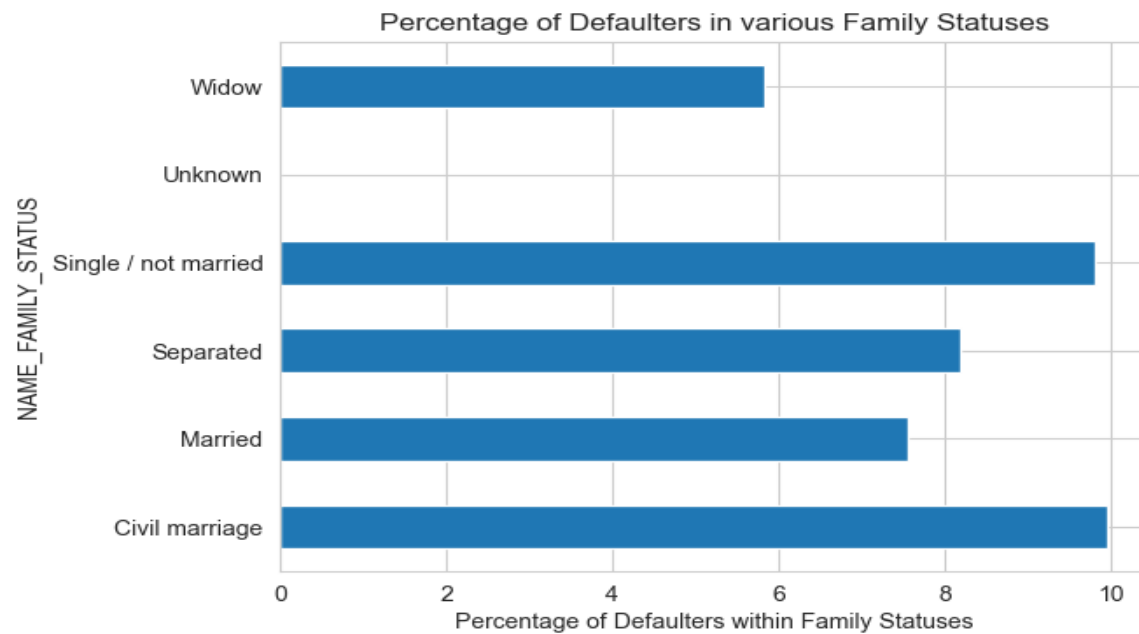


**Fig\_3**



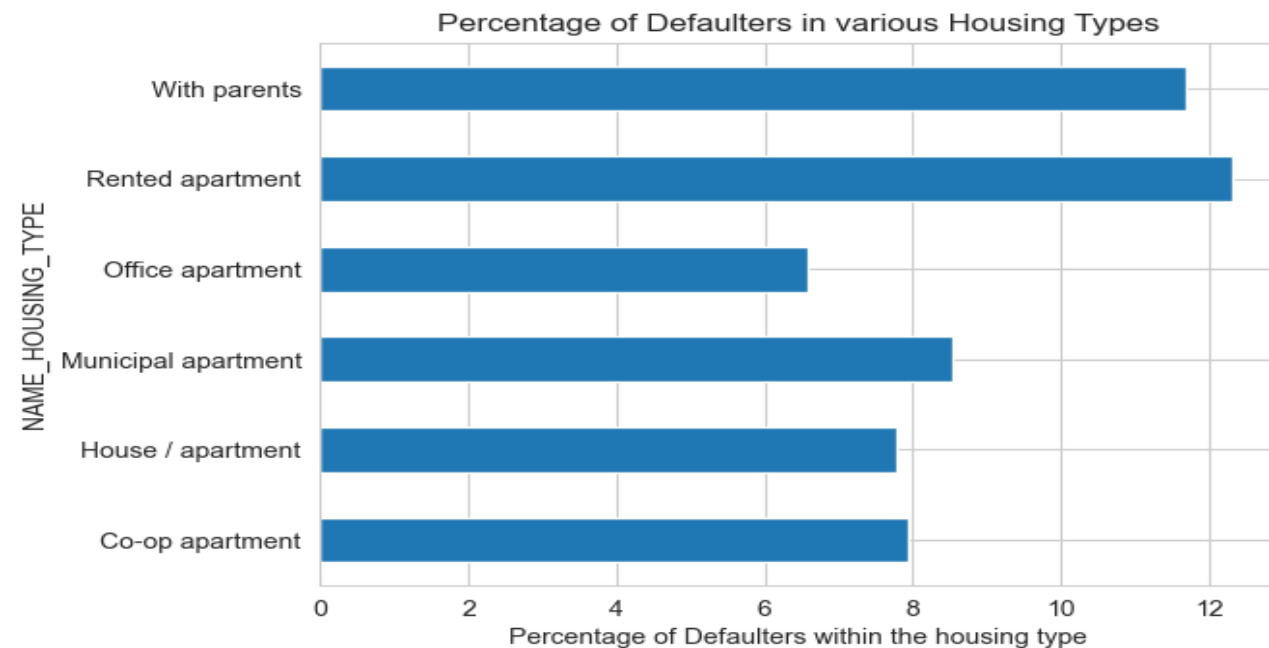
**Fig\_4**

Fig\_3 and fig\_4 are essentially showing that cleaning staff have most number of occurrences but Laborers, Low-skilled Laborers, Drivers and Security staff are categories which have most defaulters percentage within their Occupation category.



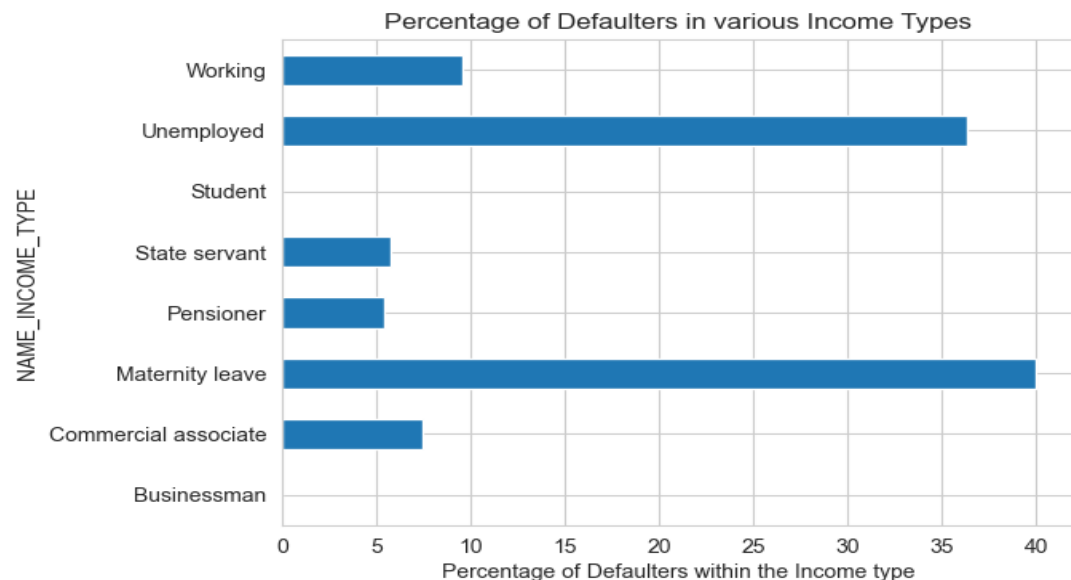
**Fig\_5**

Fig\_5 shows that single and civil marriage categories of Family statuses exhibit most probability to default on payments



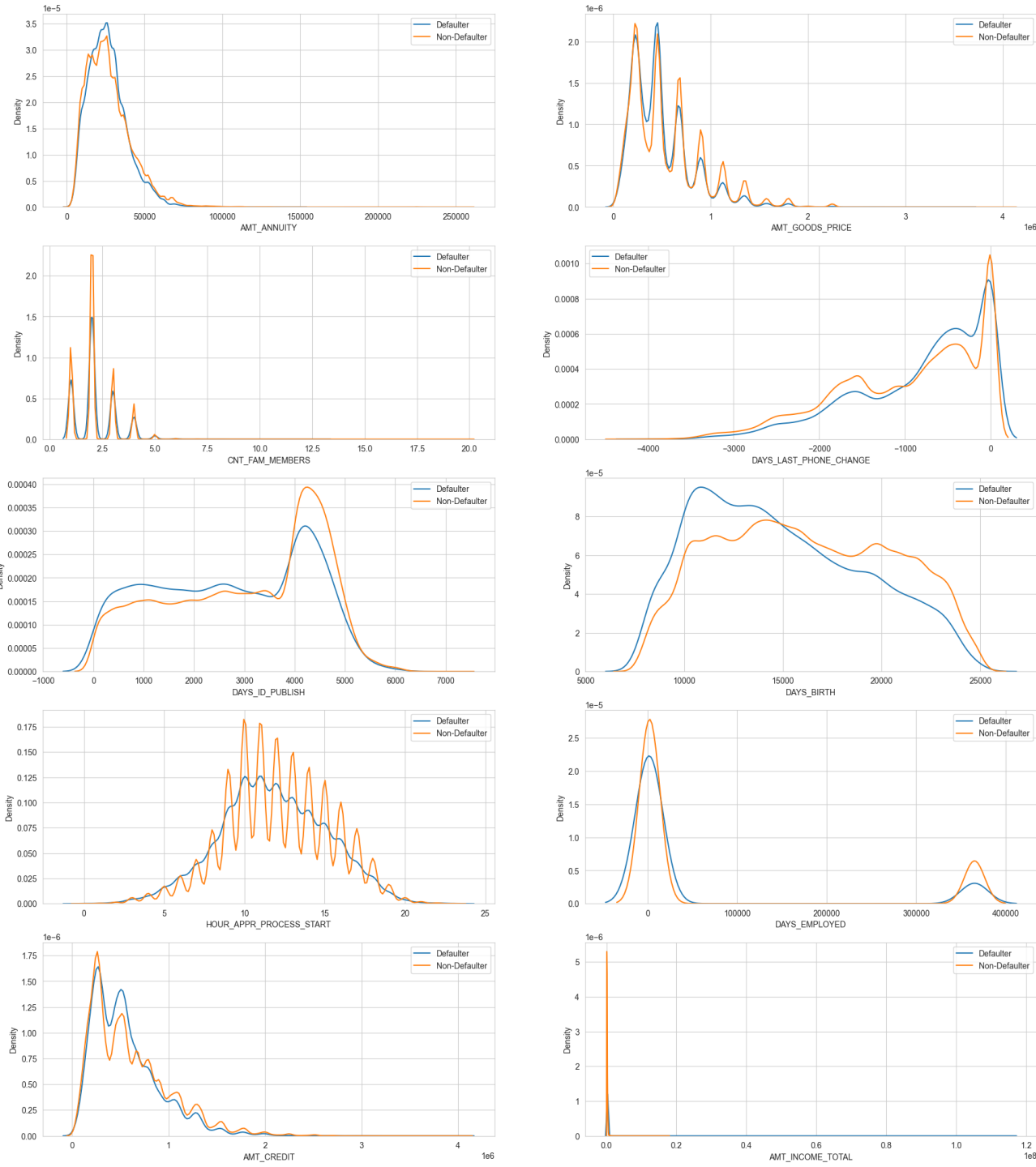
**Fig\_6**

Fig\_6 shows that People living with parents and rented apartments are most likely to default, while those living in Office apartment are least probable to default



**Fig\_7**

fig\_7 : Unemployed and maternity leave categories are most probable to have payment difficulties while Pensioners, State servant, Commercial associate and working Income Types have least probability to default

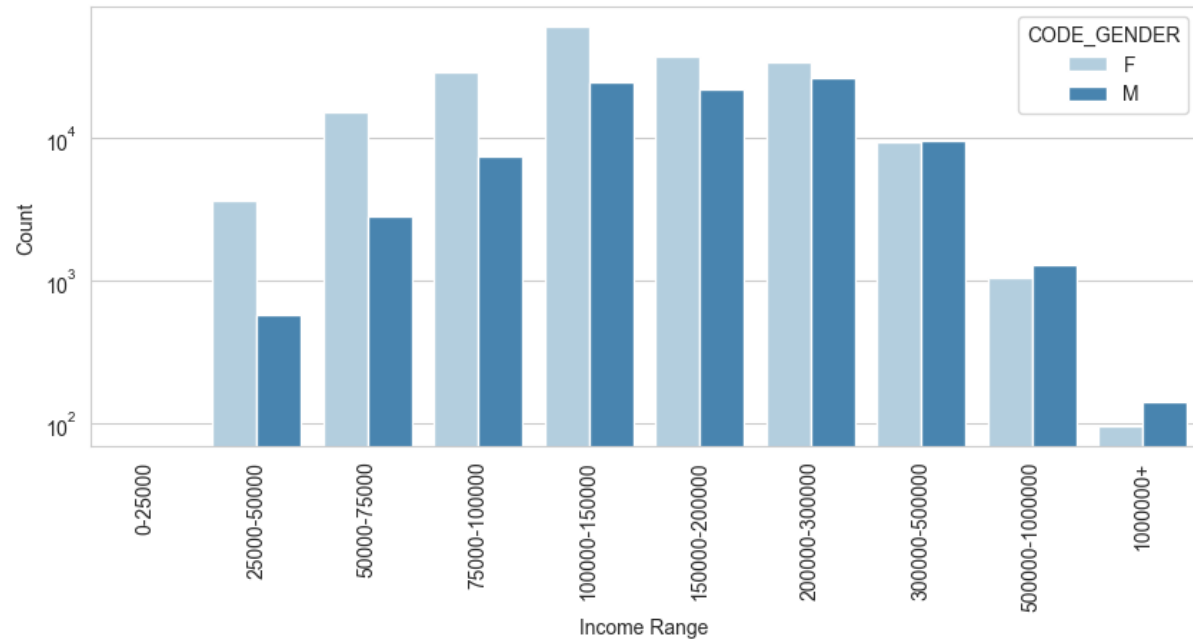


**Fig\_8:** 10 continuous numerical columns:  
'AMT\_ANNUITY', 'AMT\_GOODS\_PRICE', 'CNT\_FAM\_MEMBERS',  
'DAYS\_LAST\_PHONE\_CHANGE', 'DAYS\_ID\_PUBLISH', 'DAYS\_BIRTH',  
'HOUR\_APPR\_PROCESS\_START', 'DAYS\_EMPLOYED', 'AMT\_CREDIT',  
'AMT\_INCOME\_TOTAL'

### Insights:

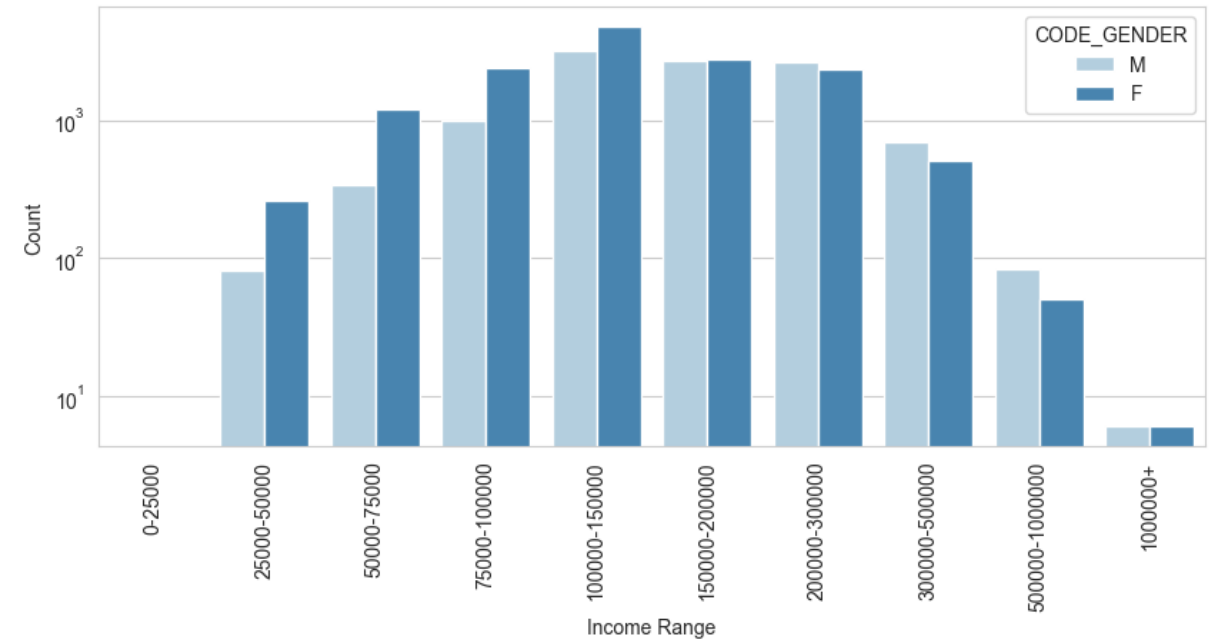
- A significant number of applications are submitted between 9 AM and 2 PM in both the Current and Previous datasets, indicating that these are the bank's busiest hours.
- Families of size 1 to 4 (i.e. nuclear families) tend to apply for more loans compared to other family types.
- Most of the applicants received lower credits based on distribution of 'AMT\_CREDIT' variable
- Looking at 'DAYS\_EMPLOYED's kdeplot, applicants are either at the starting of their careers or are very senior professionals.
- from 'DAYS\_BIRTH' we can say that younger people are more likely to default while older applicants defaulted lesser.
- 'AMT\_ANNUITY' kdeplot shows that lower annuity shows more possibility of defaulting.

Distribution of Income Range- Non-Defaulters



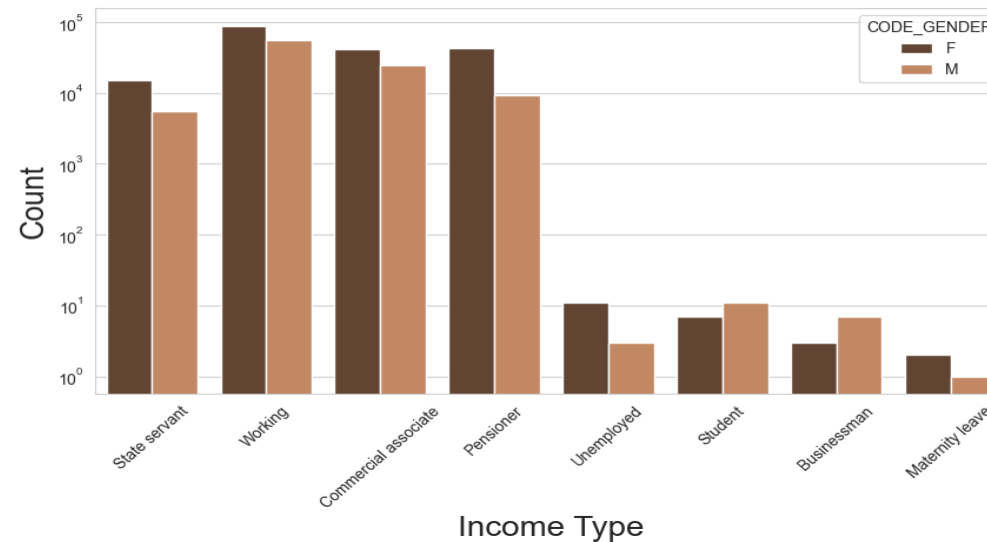
Fig\_9

Distribution of Income Range- Defaulters

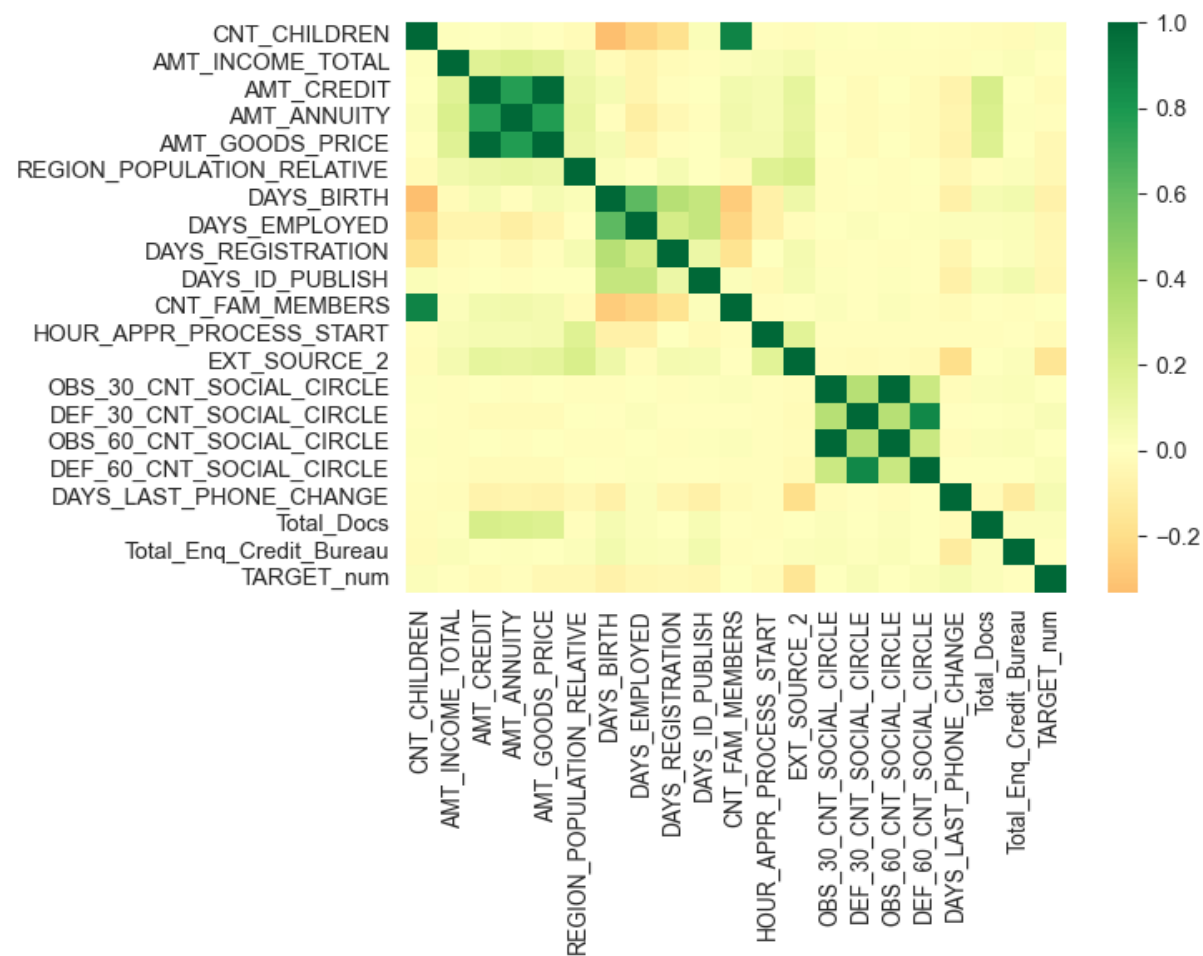


Fig\_10

Fig\_11 Distribution of Income Type



In Fig\_11, State servant, working, Commercial associate and Pensioner are the most common Income types among applicants. And Females dominate the number of applications in all these four Income types



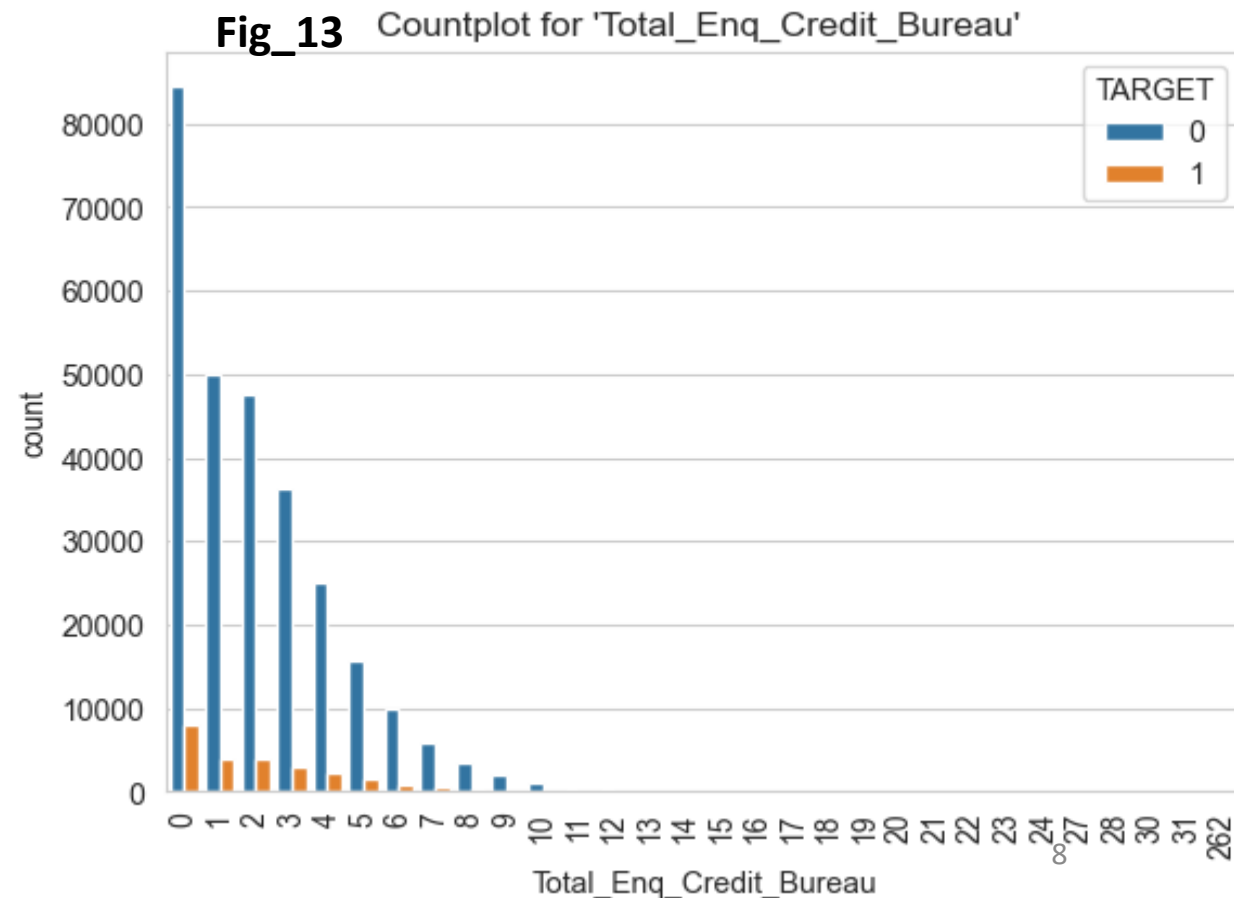
**Fig\_12**

In fig\_12, we can observe that correlation is there between the following:

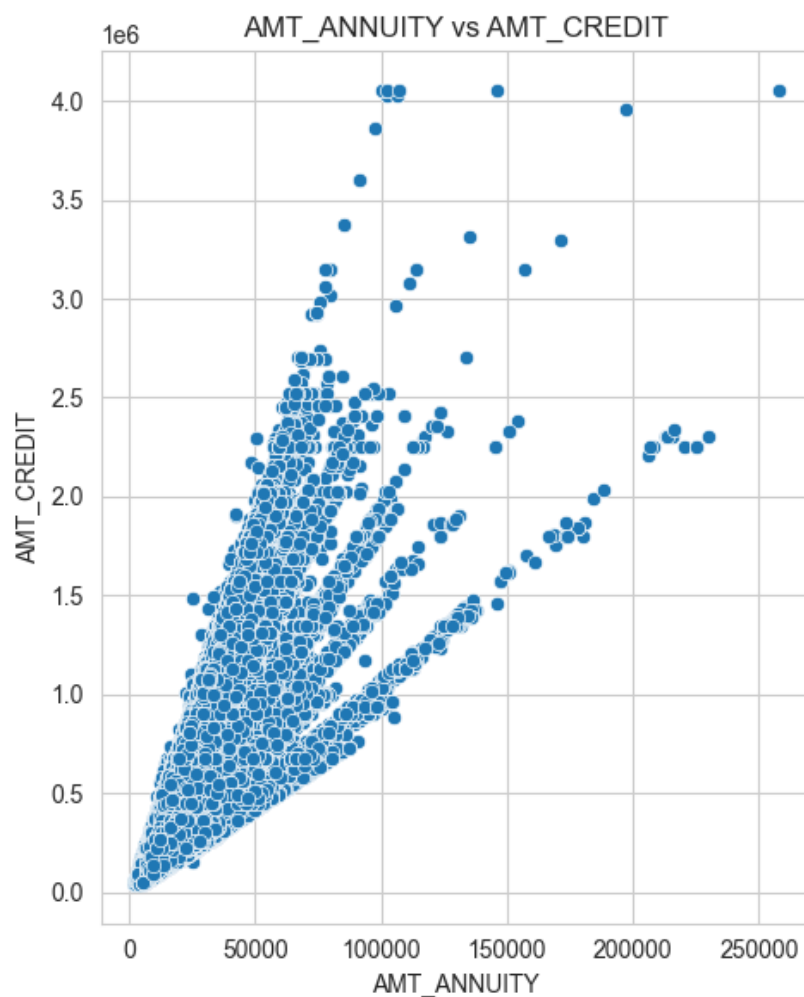
strong correlation: CNT\_FAM\_MEMBERS vs CNT\_CHILDREN, OBS\_60\_CNT\_SOCIAL\_CIRCLE vs OBS\_30\_CNT\_SOCIAL\_CIRCLE, DEF\_60\_CNT\_SOCIAL\_CIRCLE vs DEF\_30\_CNT\_SOCIAL\_CIRCLE, DAYS\_BIRTH vs DAYS\_EMPLOYED, AMT\_ANNUITY vs AMT\_CREDIT, AMT\_GOODS\_PRICE vs AMT\_CREDIT, AMT\_GOODS\_PRICE vs AMT\_ANNUITY

weak correlation: DAYS\_BIRTH vs CNT\_CHILDREN, TARGET vs EXT\_SOURCE\_2

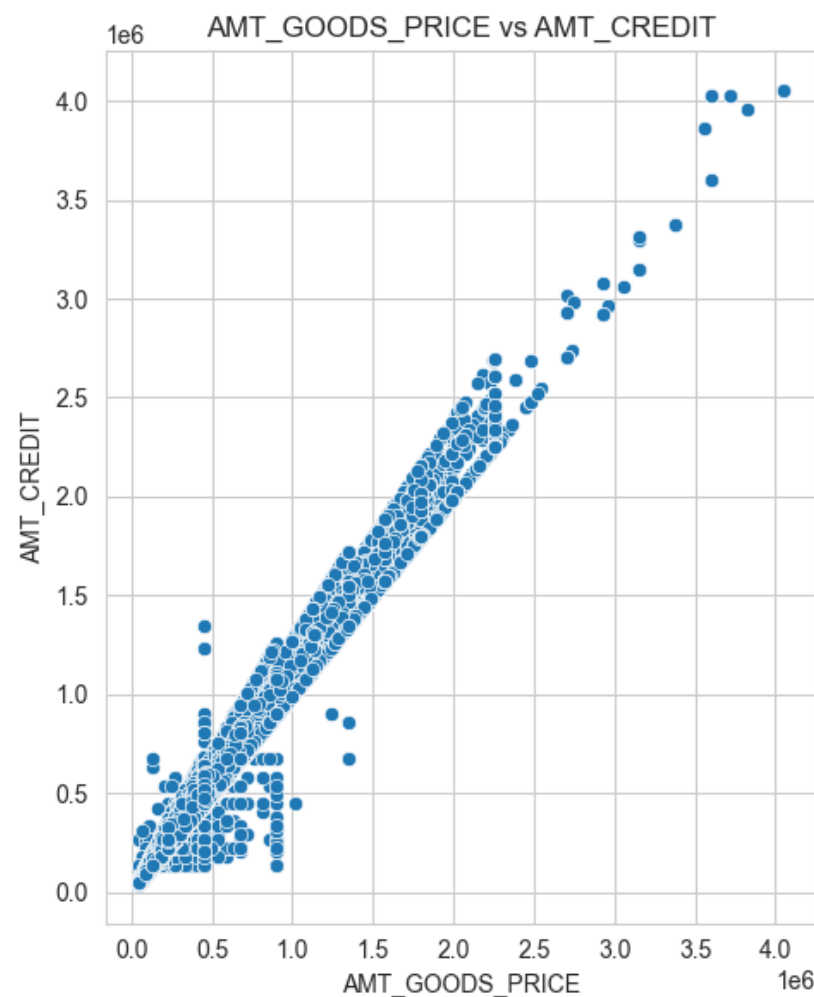
In fig\_13, we observe that most of the applicants have lesser number of credit bureau enquiries



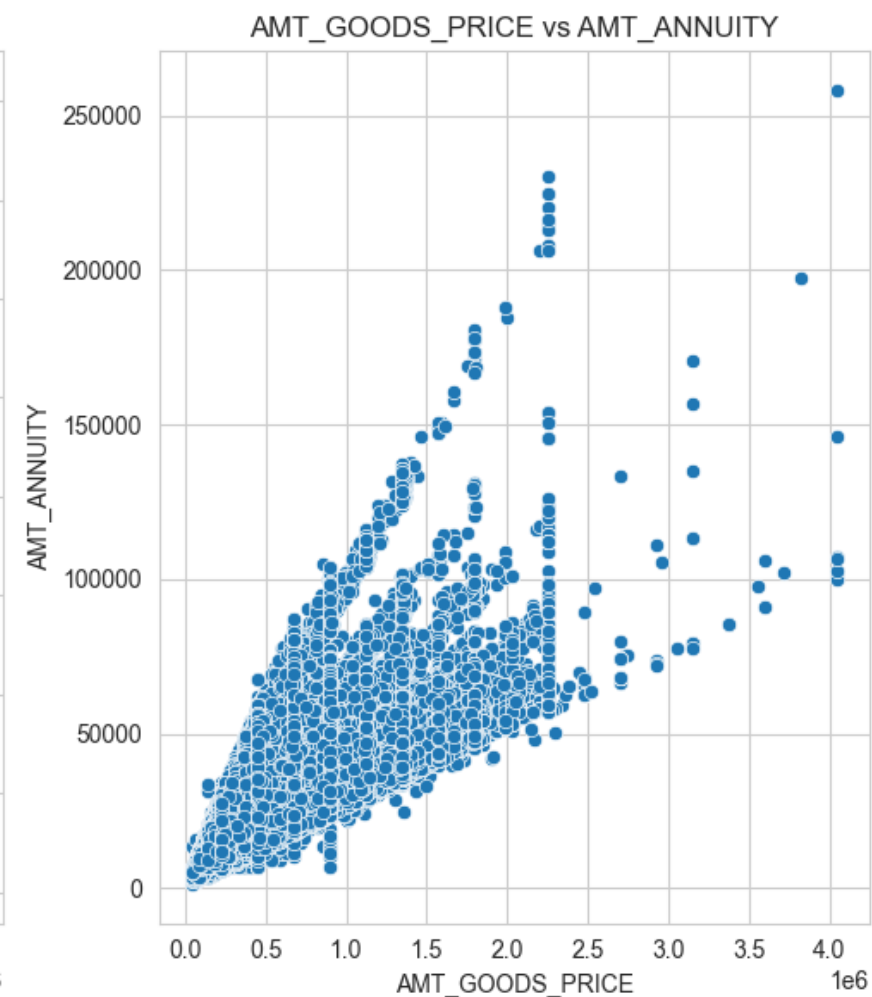




**Fig\_14**

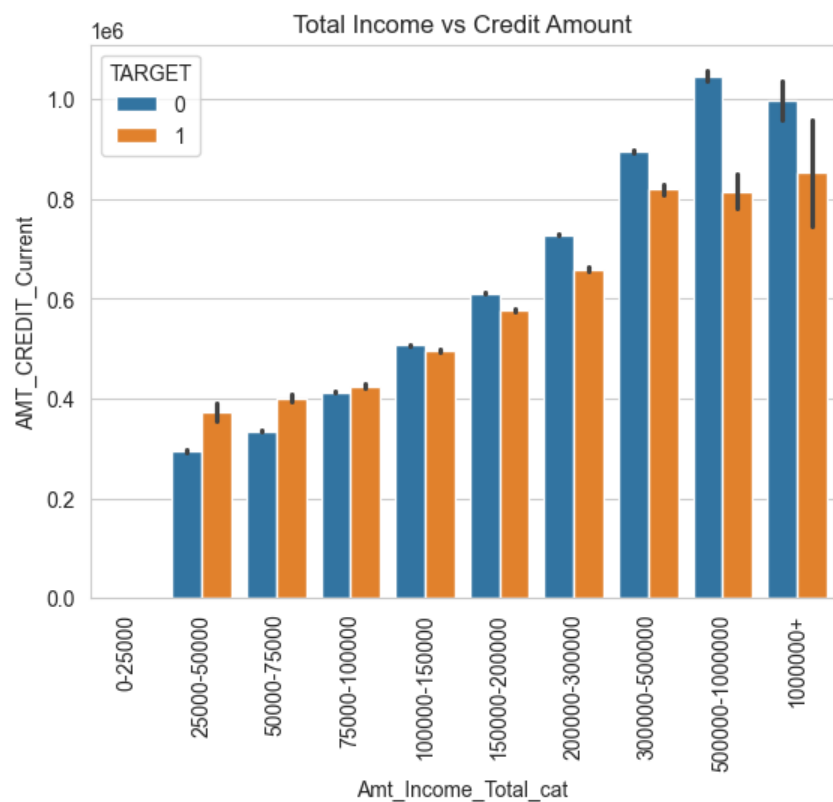


**Fig\_15**



**Fig\_16**

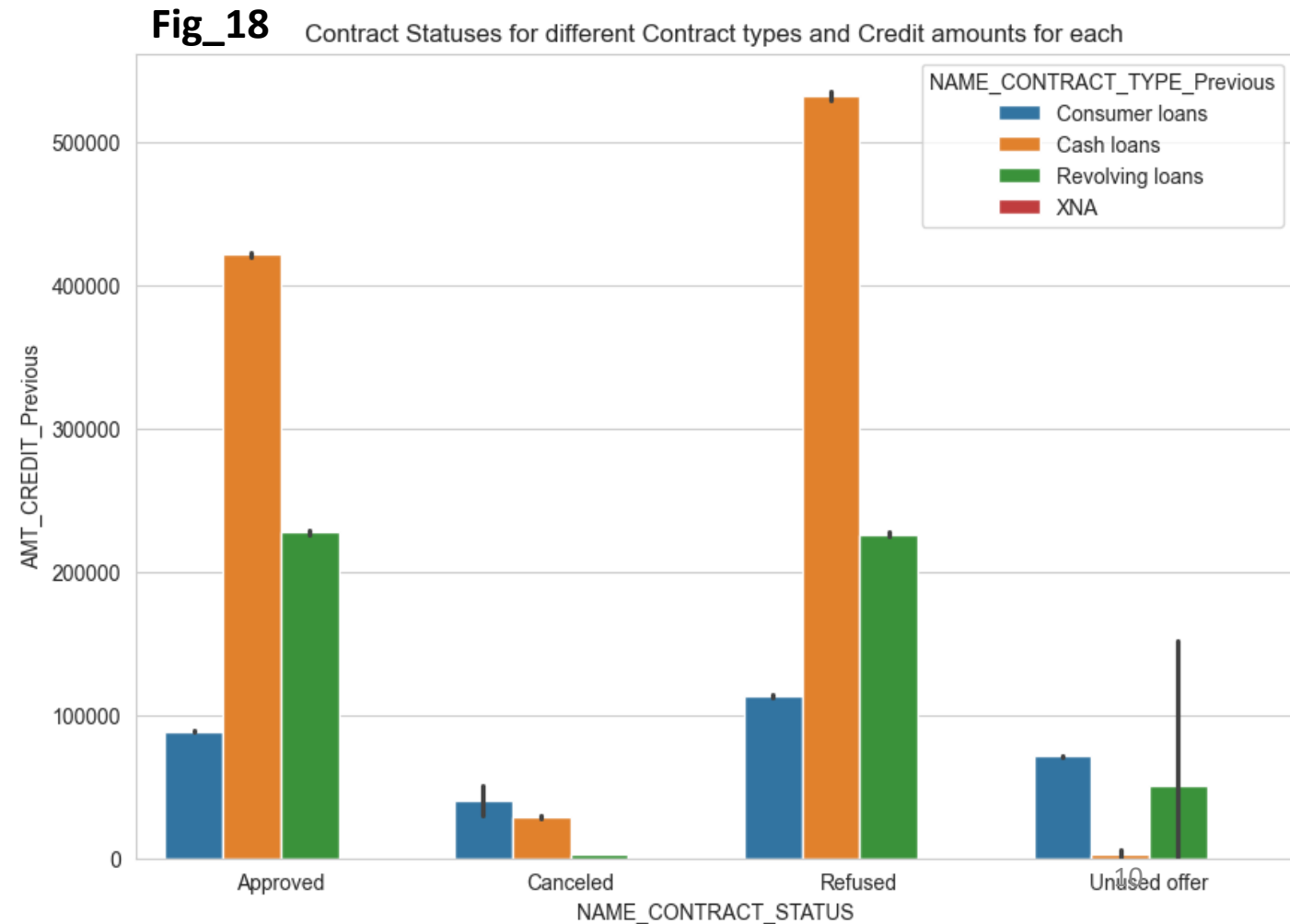
Fig\_14, Fig\_15 and Fig\_16 show that the variables: 'AMT\_ANNUITY', 'AMT\_CREDIT', 'AMT\_GOODS\_PRICE' are hugely correlated as one of them increases other two also show similar increase in values.



**Fig\_17**

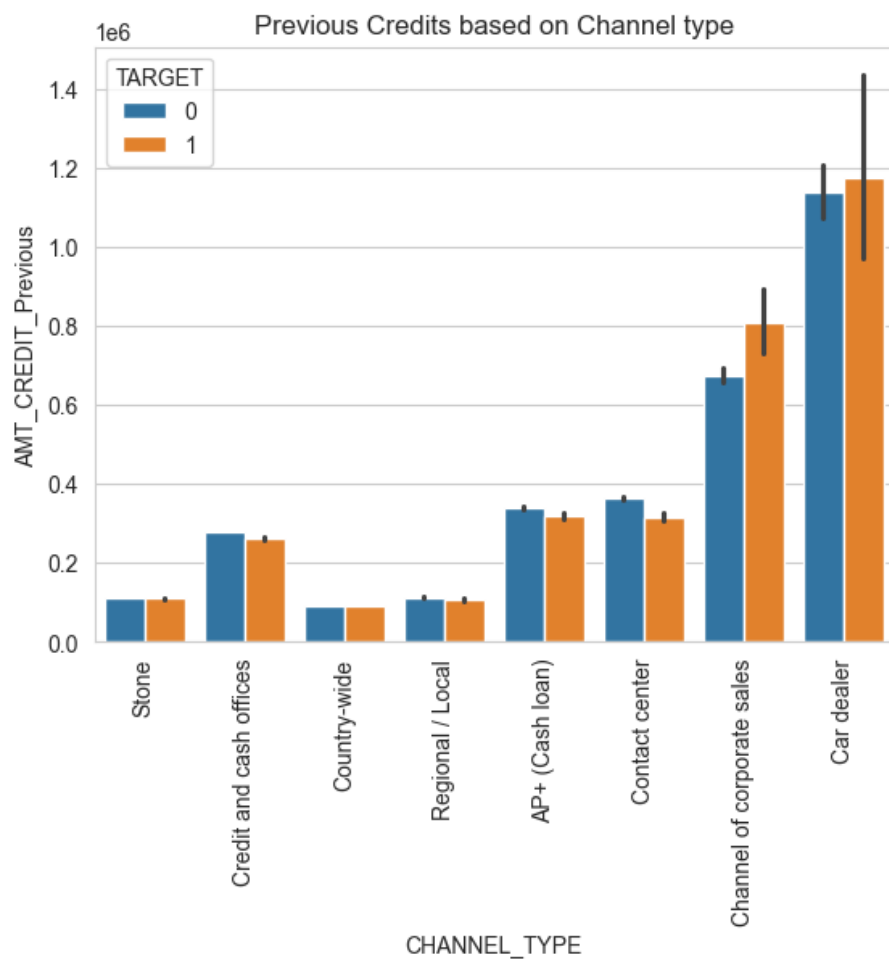
In Fig\_17, we clearly see that a higher total income increases the probability of having higher credit amount

In fig\_18, we observe that more cash loans are more refused than approved. Also Revolving loans have almost 50% chance of getting approved or refused. Revolving loan applications are rarely cancelled. Cash loans are engaged with most amount of credit



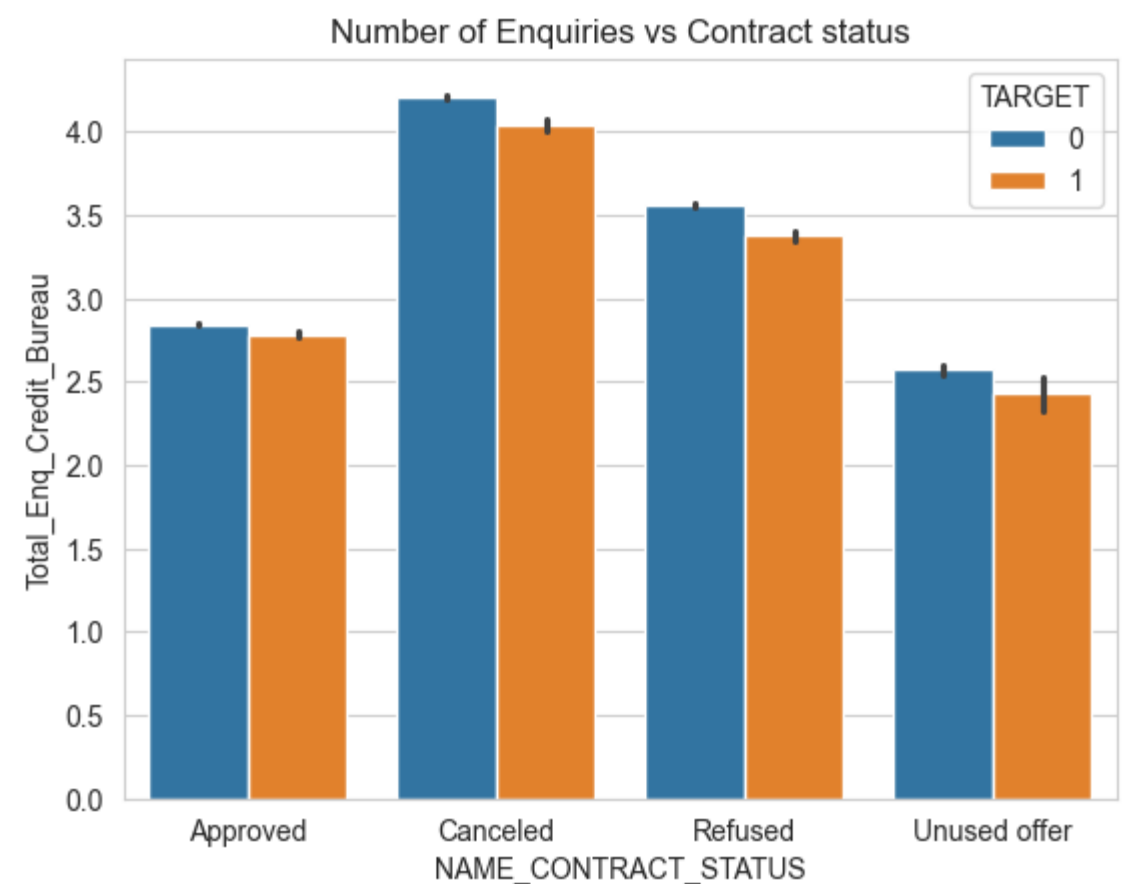
**Fig\_18**

Contract Statuses for different Contract types and Credit amounts for each



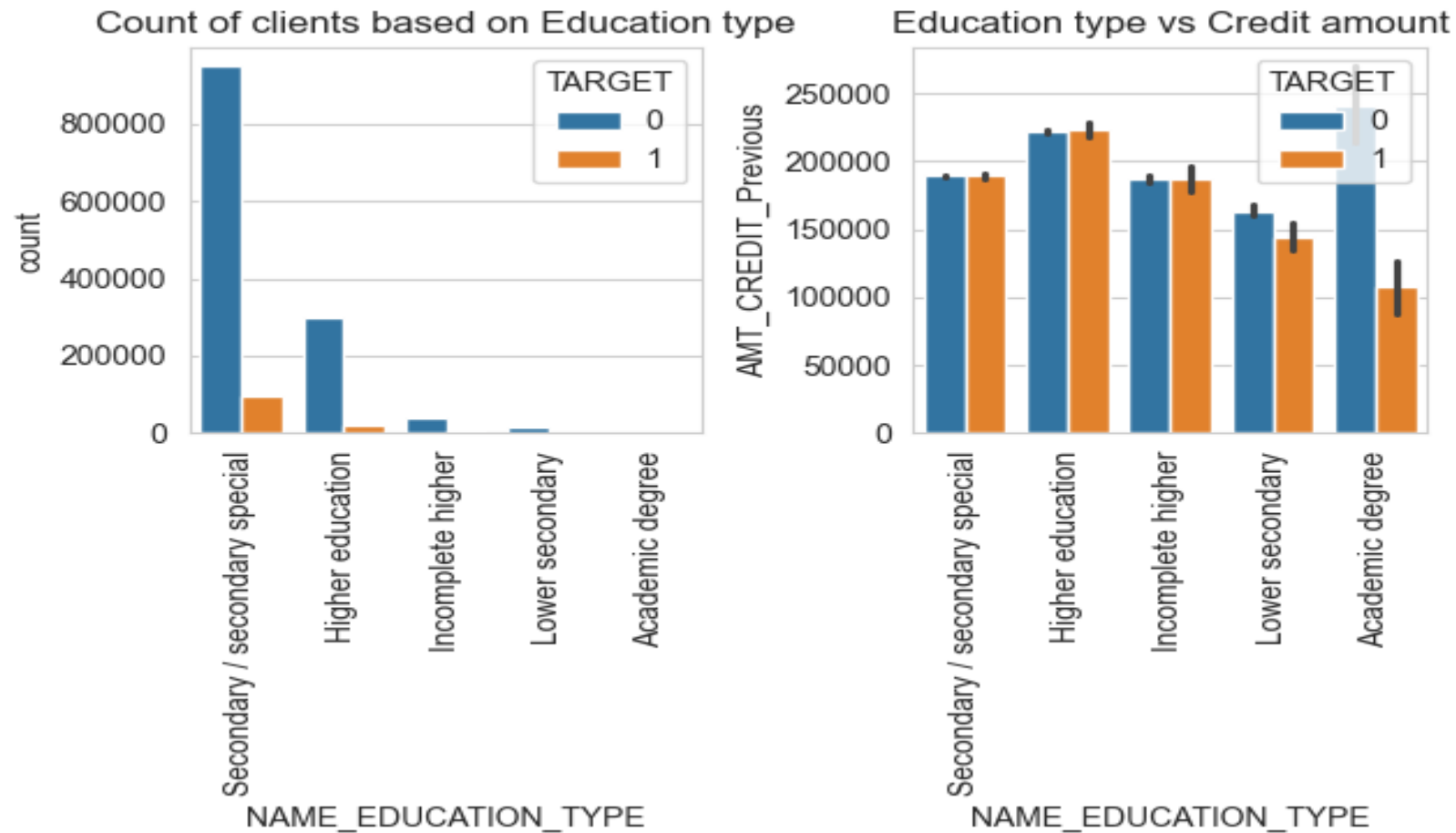
**Fig\_19**

In fig\_19, we observe that car dealers and Corporate Sales Channel have the highest amount of credit, while local channel, stone and credit and cash offices had lower credit amounts in historical data



**Fig\_20**

In fig\_20, we see that more number of Credit Bureau enquiries correspond to either Cancelling or Refusing of the contracts. While Lesser number of Enquiries increases likelihood of Approval or getting (Unused) offer for the loans. However, we cannot tell much about possibility of Defaulting based on these two variables.



**Fig\_21**

In Fig\_21, we observe that people with higher education and secondary special education apply for most number of loans. We also see that Academic degree clients are very few but they default when credit amount is lesser. People with Higher Education are most likely to get a higher credit amount

# Final Recommendations after Analysis

*The following variables exhibit significant potential in predicting the likelihood of a client defaulting on payments (preference for non-defaulting clients):*

- **AMT\_INCOME\_TOTAL** (higher=>non-defaulter)
- **NAME\_EDUCATION\_TYPE** (prefer: higher education)
- **AMT\_ANNUITY** (higher => non-defaulter)
- **NAME\_INCOME\_TYPE** (prefer: working, commercial associate, pensioner, state servant,  
avoid: Unemployed, Maternity leave)
- **DAYS\_EMPLOYED** (the more the better)
- **DAYS\_BIRTH** (the more the better)
- **Total\_Enq\_Credit\_Bureau** (custom column) (lesser the better)
- **CHANNEL\_TYPE** (contact center has lease defaulter to non defaulters ratio)
- **OCCUPATION\_TYPE** (prefer: private service, medicine, high skill teck staff, realty agents, secretaries;  
avoid: low-skilled laborers, laborers, cleaning staff, sales staff, drivers)
- **NAME\_HOUSING\_TYPE** (prefer: office apartment,house/apartment, co-op apt  
avoid: with parents and rented apartment)

This concludes our analysis. However, data analysis is inherently an ongoing process, and there are always additional insights that can be discovered. Therefore, the journey of exploration and understanding never truly ends.

# The End