

# BANK LOAN CASE STUDY

Data Analytics Project by Soumya Singhal

# PROJECT DESCRIPTION

- In this project, the goal is to analyze loan application data to identify patterns that predict the likelihood of loan default.
- The project involves Exploratory Data Analysis (EDA) to explore how different customer attributes and loan characteristics influence the probability of defaulting on a loan.
- The key business objectives are to reduce risk by identifying potentially problematic applicants and improving the loan approval process.

# APPROACH

To achieve the project objectives, a systematic and structured approach was followed, leveraging Excel's advanced features and statistical techniques to derive meaningful insights.

## **1.Data Preparation:**

- Analyzed missing data using functions like COUNTBLANK and IF.
- Filled missing values with median or mode where applicable, and removed columns with over 40% missing data.

## **2.Outlier Detection:**

- Used statistical measures like IQR to identify and visualize outliers.
- Created box plots to display the data distribution and pinpoint anomalies.

### 3.Data Imbalance Analysis:

- Evaluated the distribution of the target variable using COUNTIF.
- Visualized the class imbalance through pie charts and bar charts to highlight potential issues.

### 4.Exploratory Data Analysis:

- Performed **Univariate Analysis** to examine individual variable distributions using histograms and bar charts.
- Conducted **Segmented Univariate Analysis** by comparing variable distributions for defaulters and non-defaulters.
- Executed **Bivariate Analysis** to study relationships between variables, employing scatter plots and heatmaps for insights.

### 5.Correlation Analysis:

- Generated correlation matrices for the overall dataset, defaulters, and non-defaulters separately.
- Highlighted strong correlations ( $>0.5$ ) using color coding for easy interpretation.
- Identified key correlated variables and summarized them in a table for further analysis.

This comprehensive approach ensured a deep understanding of the dataset and provided actionable insights into loan default patterns.

## TECH-STACK USED

- MS Excel
- MS PowerPoint

## TASK A: IDENTIFY AND HANDLE MISSING DATA

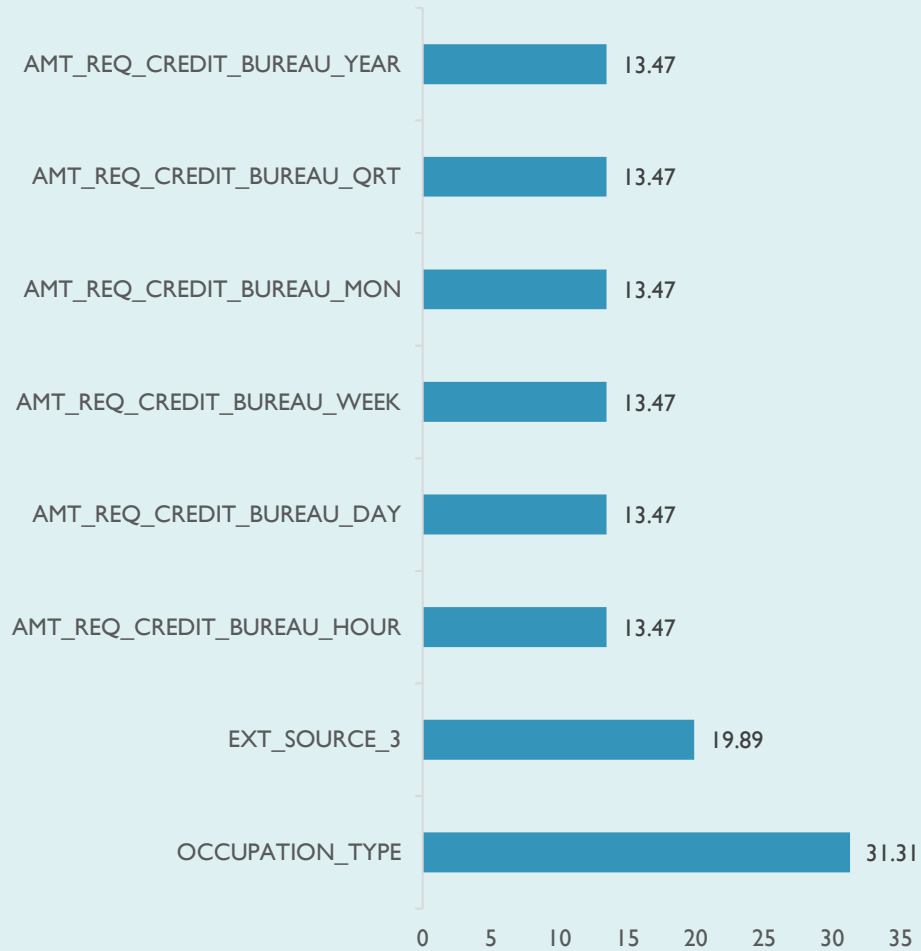
- Counted the number of blank rows in each column using the COUNTBLANK function and calculated the percentage of missing values.
- Created a bar chart and column chart to visualise the proportion of missing values for each variable.
- Removed columns(49) that contained more than 40% blank cells.
- Filled missing values in the `NAME\_TYPE\_SUITE` column using the mode i.e. Unaccompanied.
- Filled missing values in other columns, where less than 40% of the data was missing, using the median or mode.

## NULL PERCENT >40%

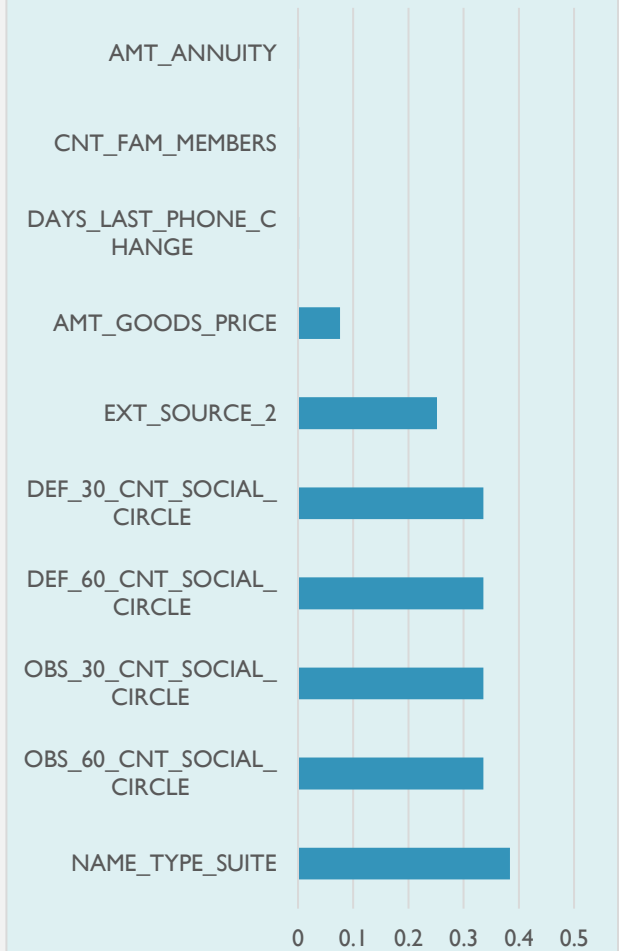
80  
70  
60  
50  
40  
30  
20  
10  
0

COMMONAREA\_AVG  
COMMONAREA\_MODE  
COMMONAREA\_MEDI  
NONLIVINGAPARTMENTS\_AVG  
NONLIVINGAPARTMENTS\_MODE  
NONLIVINGAPARTMENTS\_MEDI  
LIVINGAPARTMENTS\_AVG  
LIVINGAPARTMENTS\_MODE  
LIVINGAPARTMENTS\_MEDI  
FONDKAPREMENT\_MODE  
FLOORSMIN\_AVG  
FLOORSMIN\_MODE  
FLOORSMIN\_MEDI  
YEARS\_BUILD\_AVG  
YEARS\_BUILD\_MODE  
YEARS\_BUILD\_MEDI  
OWN\_CAR\_AGE  
LANDAREA\_AVG  
LANDAREA\_MODE  
LANDAREA\_MEDI  
BASEMENTAREA\_AVG  
BASEMENTAREA\_MODE  
BASEMENTAREA\_MEDI  
EXT\_SOURCE\_I  
NONLIVINGAREA\_AVG  
NONLIVINGAREA\_MODE  
NONLIVINGAREA\_MEDI  
ELEVATORS\_AVG  
ELEVATORS\_MODE  
ELEVATORS\_MEDI  
WALLSMATERIAL\_MODE  
APARTMENTS\_AVG  
APARTMENTS\_MODE  
APARTMENTS\_MEDI  
ENTRANCES\_AVG  
ENTRANCES\_MODE  
ENTRANCES\_MEDI  
LIVINGAREA\_AVG  
LIVINGAREA\_MODE  
LIVINGAREA\_MEDI  
HOUSETYPE\_MODE  
FLOORSMAX\_AVG  
FLOORSMAX\_MODE  
FLOORSMAX\_MEDI  
YEARS\_BEGINEXPLUATATION\_AVG  
YEARS\_BEGINEXPLUATATION\_MODE  
YEARS\_BEGINEXPLUATATION\_MEDI  
TOTALAREA\_MODE  
EMERGENCYSTATE\_MODE

### NULL PERCENT



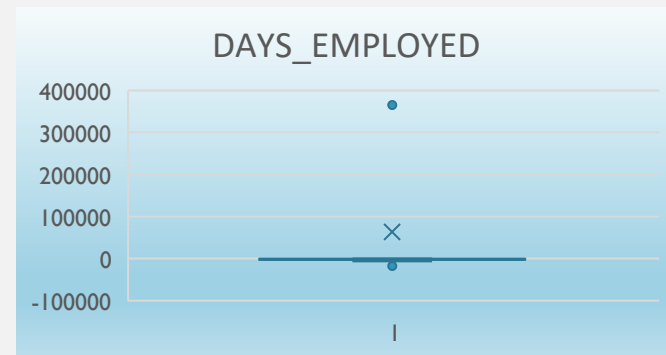
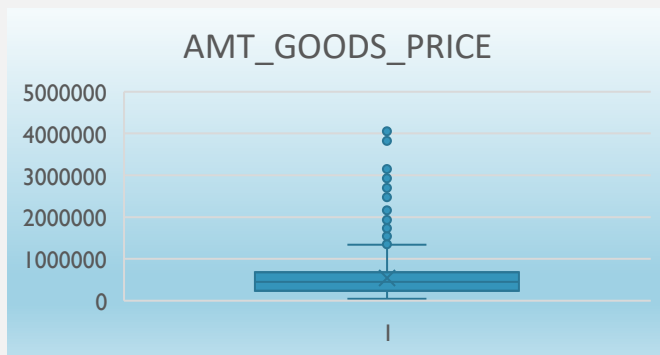
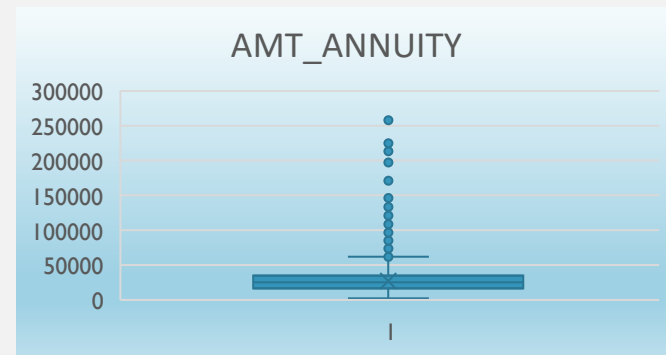
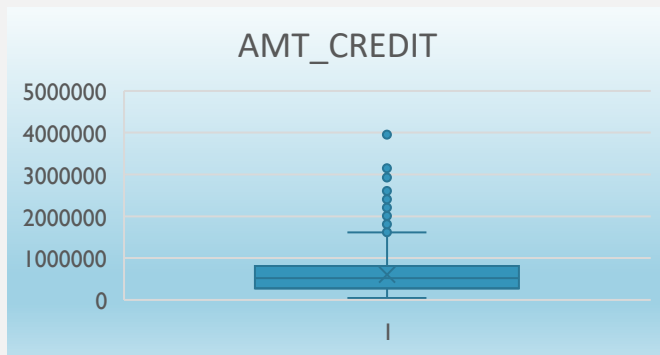
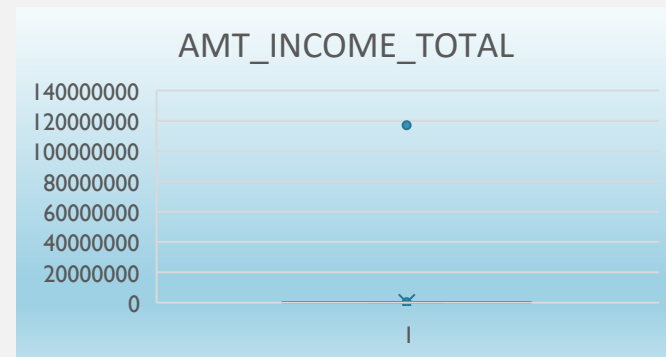
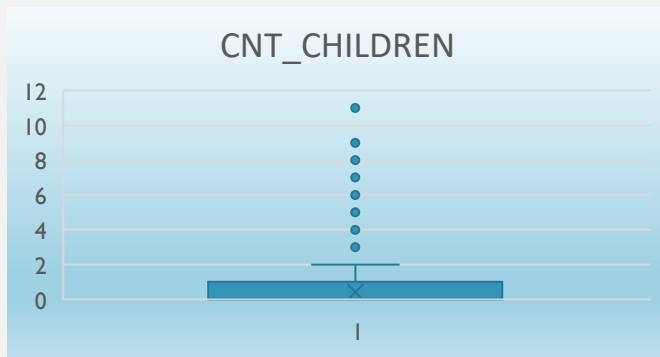
### Blanks <1%

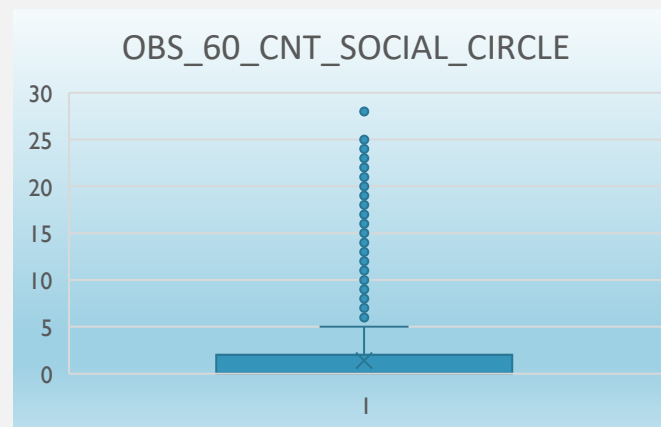
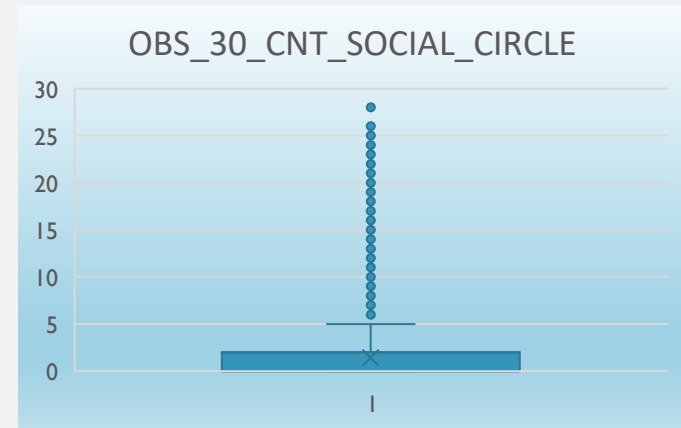
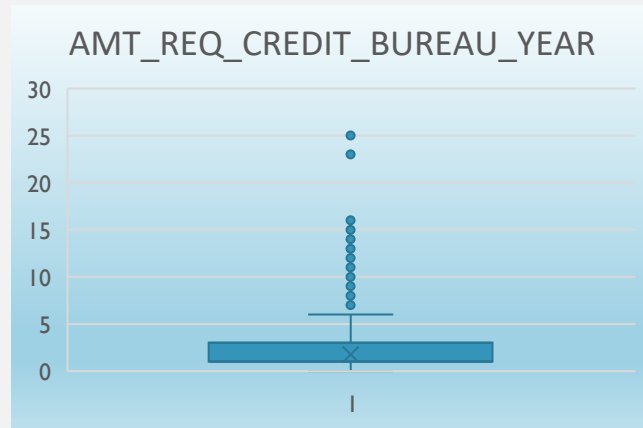




## TASK B: IDENTIFY OUTLIERS

- Detect outliers using statistical methods such as QUARTILE.EXC and IQR in Excel.
- Visualize outliers using box plots or scatter plots.

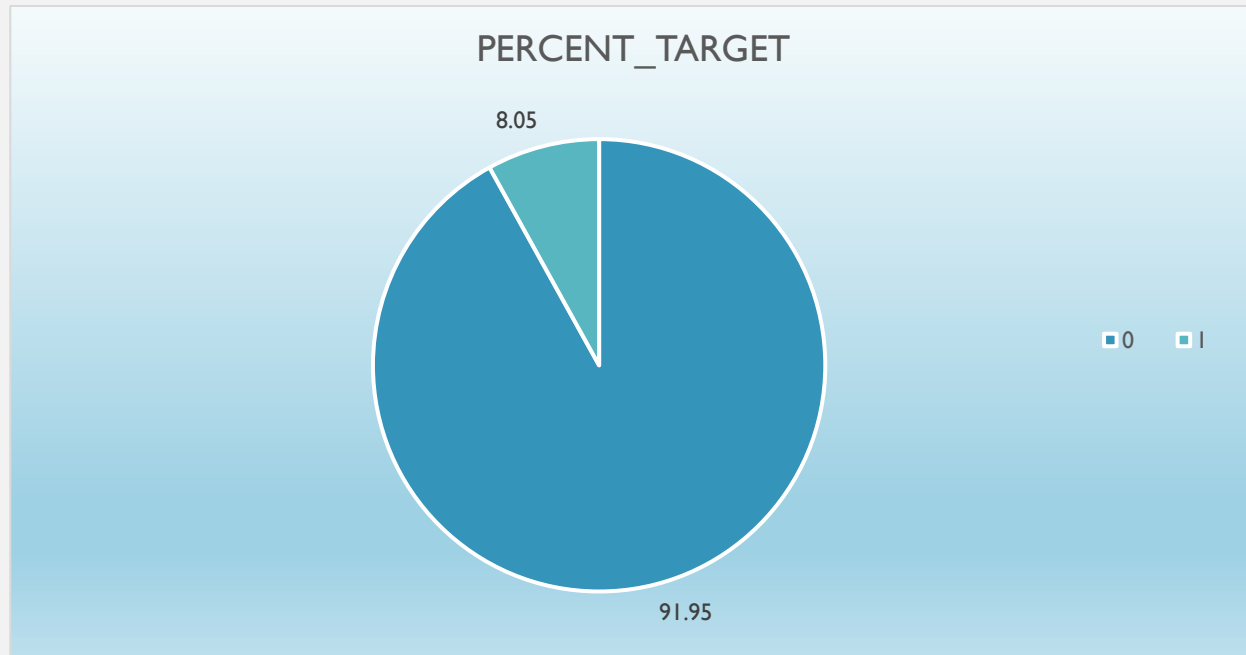




## TASK C: ANALYZE DATA IMBALANCE

- Determine if data imbalance exists using COUNTIF and SUM functions in Excel.
- Visualize the distribution of the target variable with pie charts or bar charts.
- The data is significantly imbalanced for non-defaulters (Target = 0).
- I have created a bar chart to clearly illustrate this imbalance.

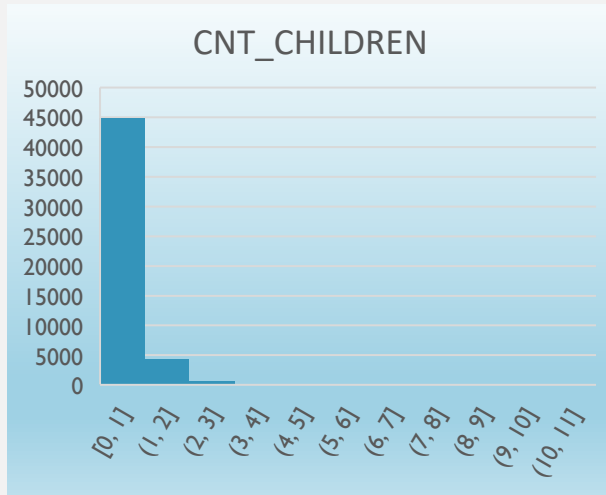
TARGET	COUNT_TARGET	PERCENT_TARGET
0	45973	91.94783896
1	4026	8.052161043



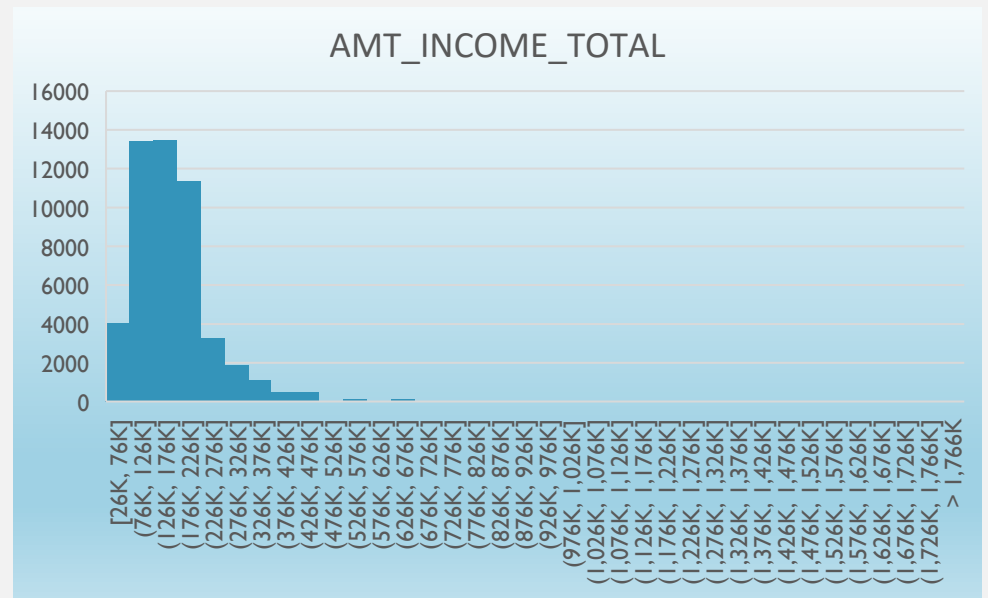
## TASK D: UNIVARIATE AND BIVARIATE ANALYSIS

- Perform univariate analysis to understand individual variable distributions.
- Conduct bivariate analysis to explore relationships between variables and the target variable.
- Visualize the relationships using histograms, bar charts, scatter plots, and heatmaps.

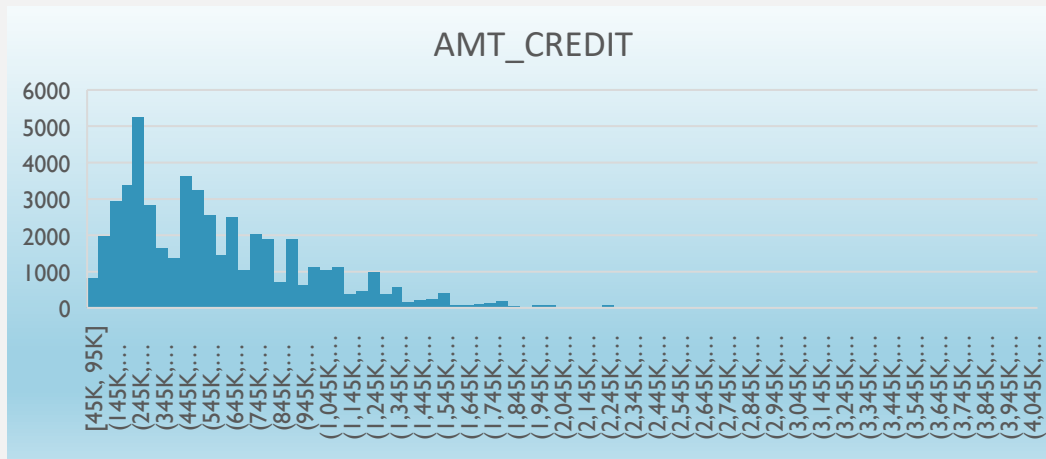
Coloumns	AVERAGE	MEDIAN	STD DV	MINIMUM	MAXIMUM
CNT_CHILDREN	0.419848397	0	0.724031307	0	11
AMT_INCOME_TOTAL	170767.5905	145800	531813.7768	25650	117000000
AMT_CREDIT	599700.5815	514777.5	402411.4096	45000	4050000
AMT_ANNUITY	27107.33399	24939	14562.6564	2052	258025.5
AMT_GOODS_PRICE	538992.3491	450000	369717.1252	45000	4050000
REGION_RATING_CLI ENT	2.051661033	2	0.507972786	1	3
YEAR_BIRTH	43.8960057	43.09863014	11.94892234	21.04109589	68.99726027



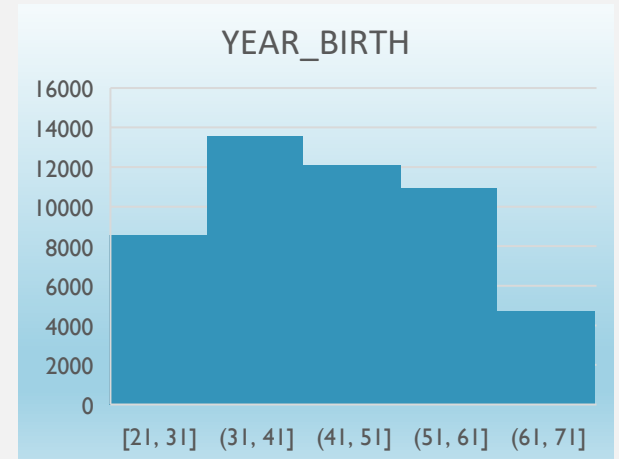
Most people who take out loans have either no children or just one child.



Individuals with incomes between 70k and 2.5 lakh tend to take out more loans, whereas those with higher incomes may require fewer loans.



Most people require loans ranging from 250,000 to 300,000 rupees.

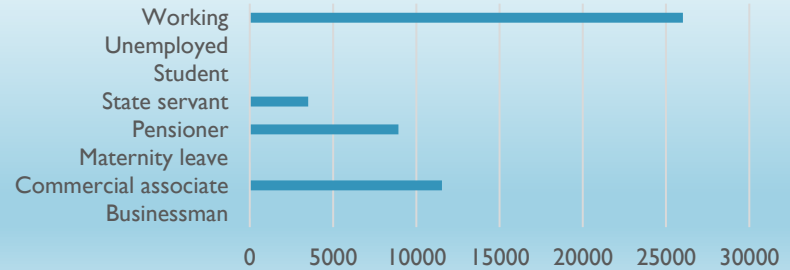


People aged 31 to 41, along with those in middle age, are more likely to take out loans.

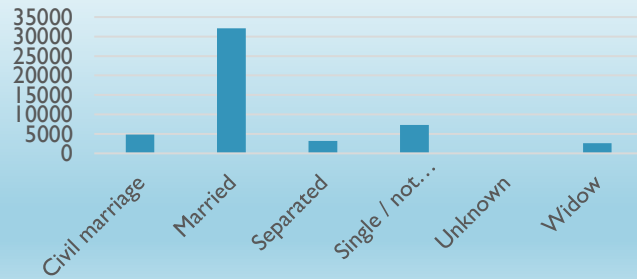
NAME\_CONTRACT\_TYPE



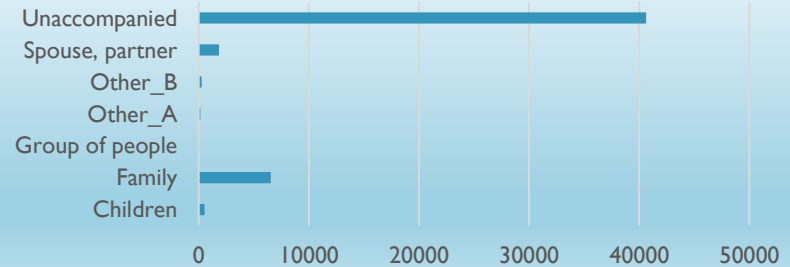
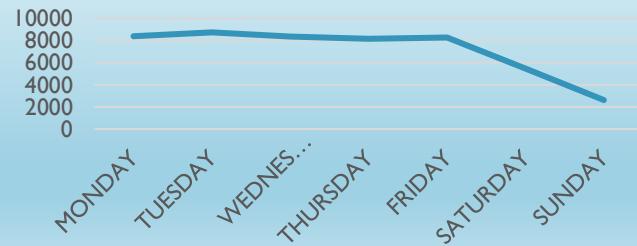
NAME\_INCOME\_TYPE



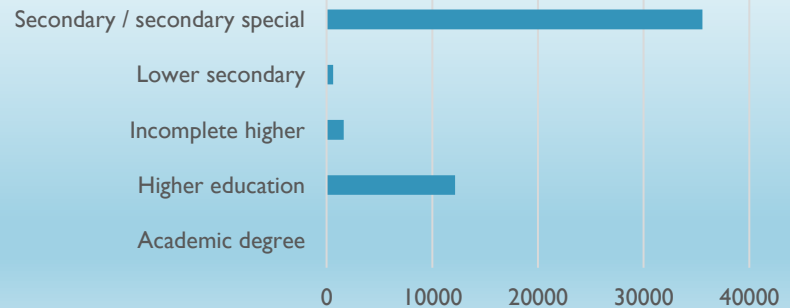
NAME\_FAMILY\_STATUS



NAME\_TYPE\_SUITE

Count of  
WEEKDAY\_APPR\_PROCESS\_ST  
ART

NAME\_EDUCATION\_TYPE

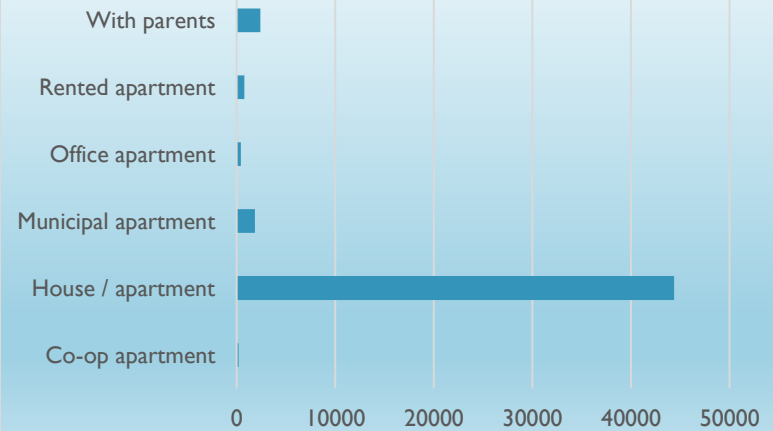




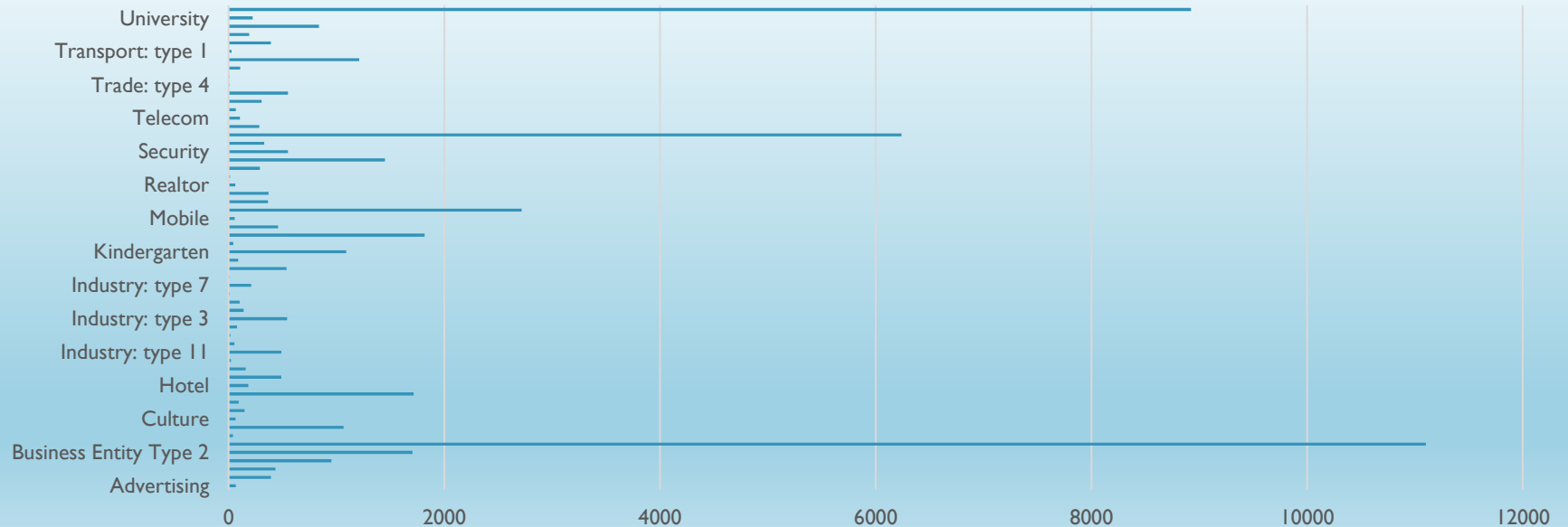
### Count of OCCUPATION\_TYPE

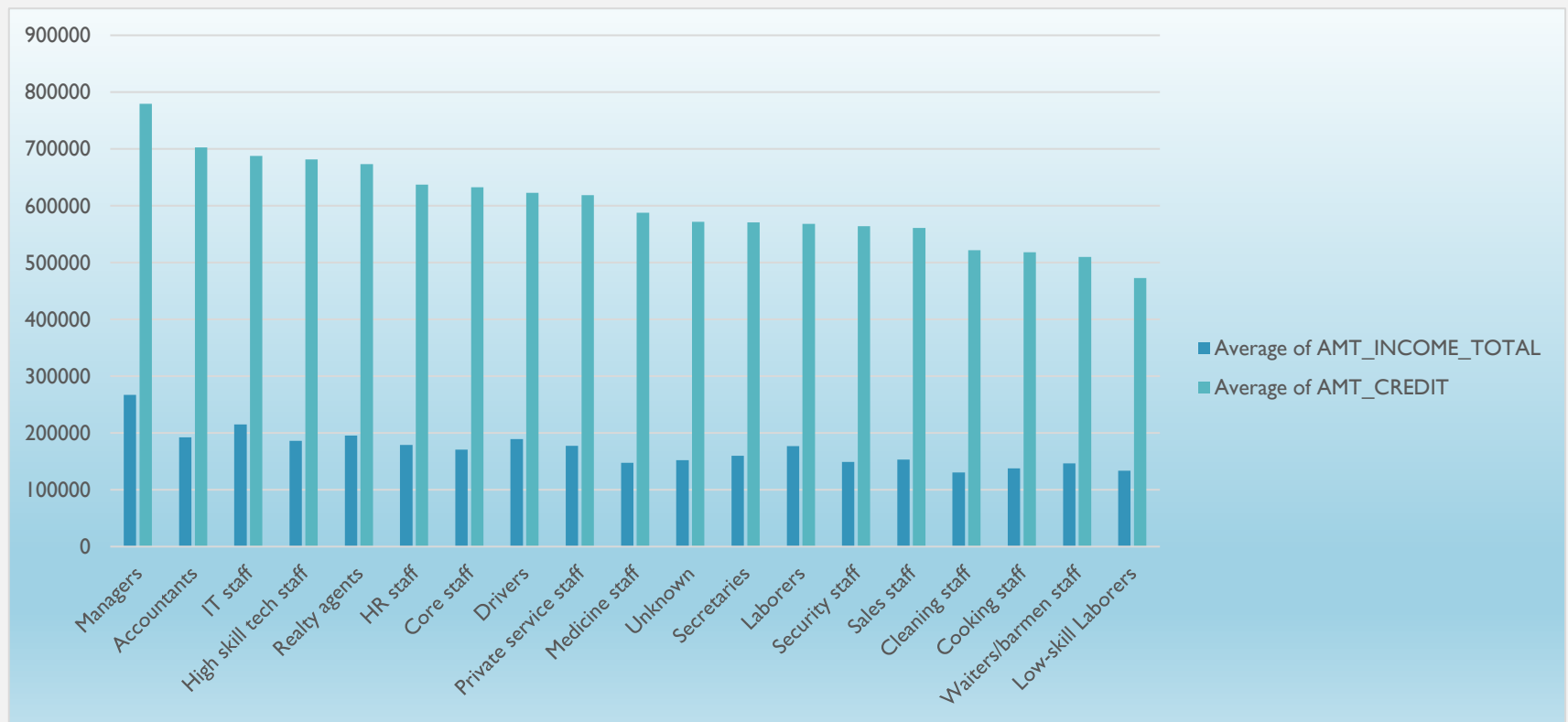
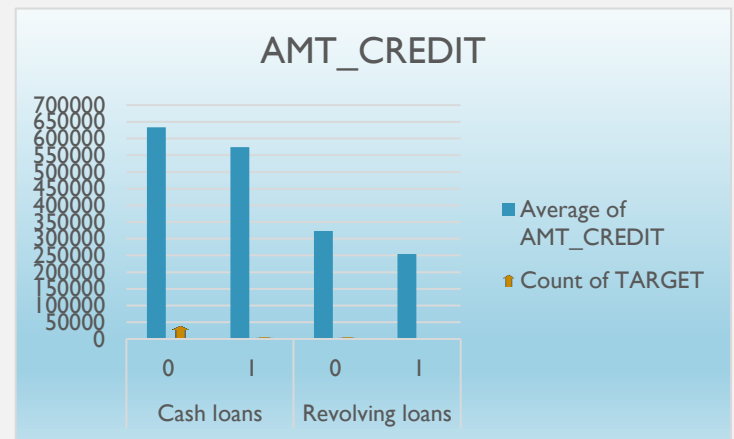
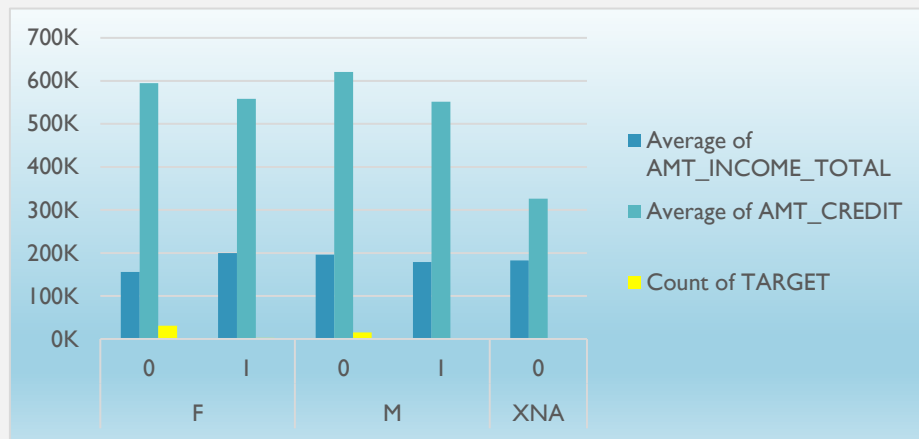


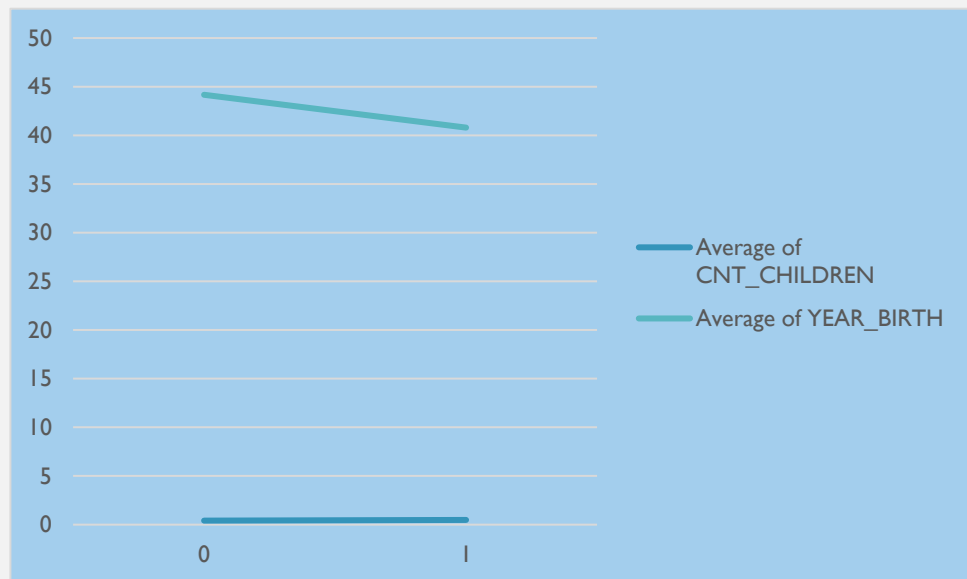
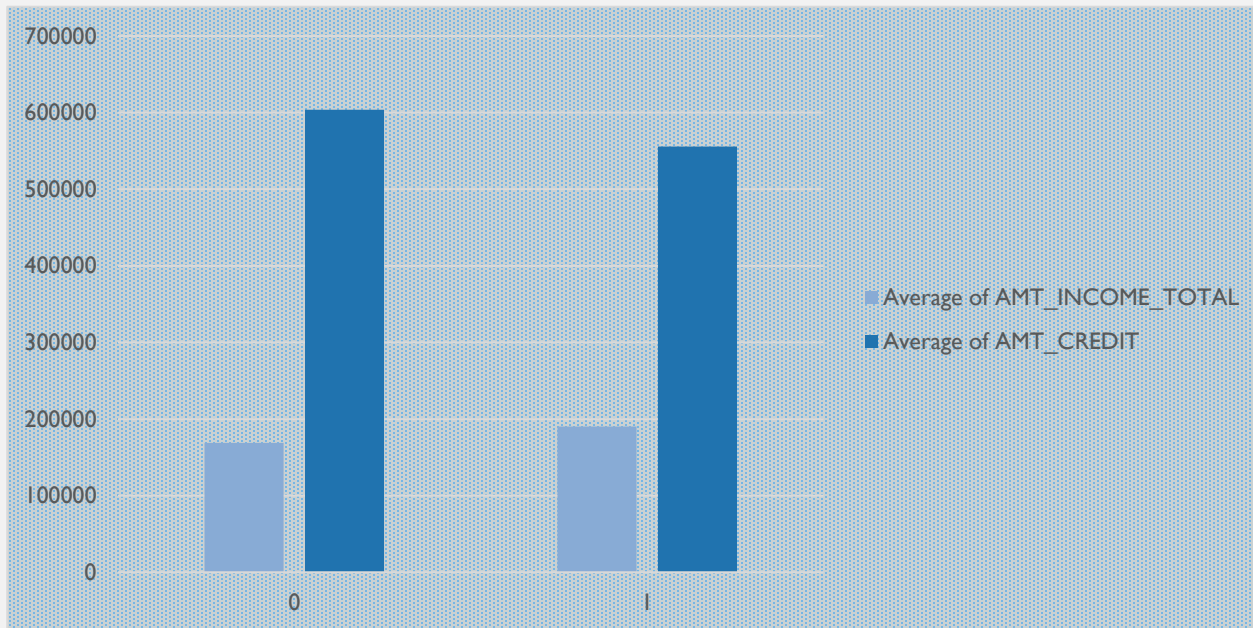
### Count of NAME\_HOUSING\_TYPE

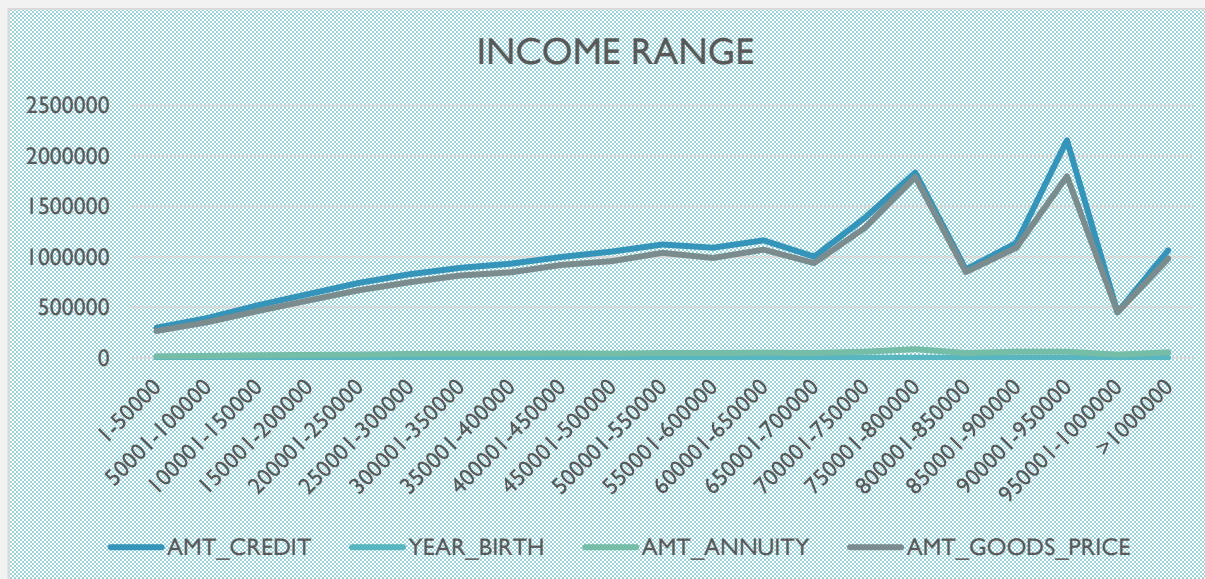
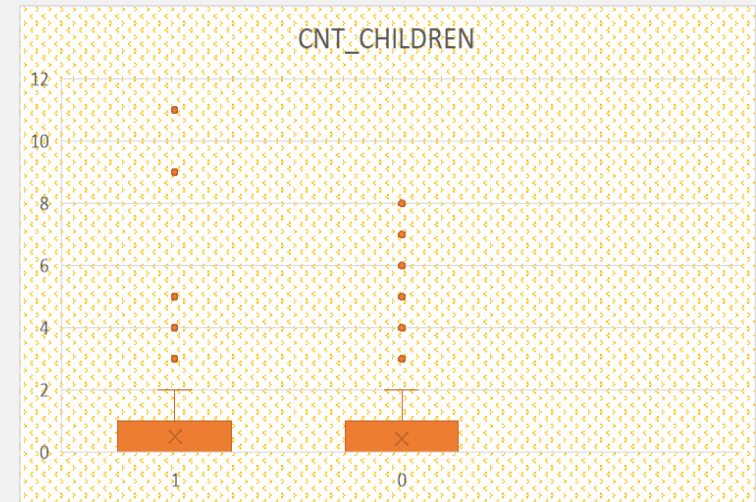
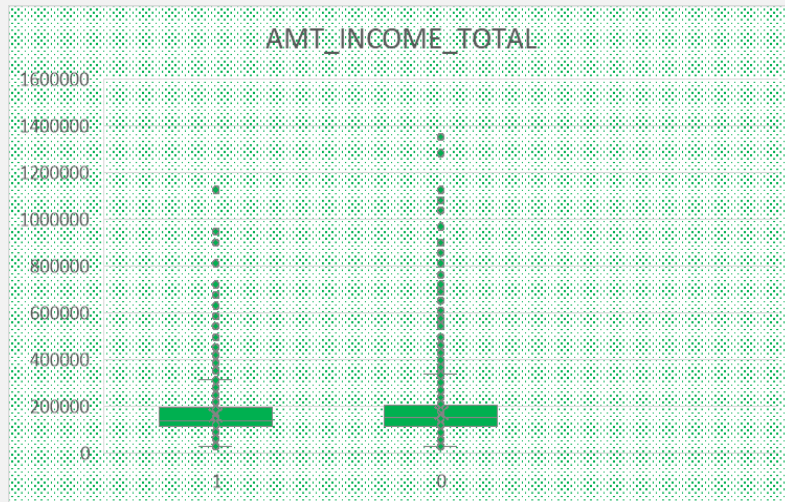


### Count of ORGANIZATION\_TYPE









## TASK E: IDENTIFY TOP CORRELATIONS

- Created correlation matrices for the entire dataset, as well as for defaulters and non-defaulters.
- Organized the data into tables for improved clarity.
- Used deep red to highlight correlations greater than 0.5 and light red for correlations between 0 and 0.5.
- Summarized the highest correlations in a table of key insights.
- Separated the analysis for defaulters and non-defaulters to identify specific patterns within each group.

# COMBINED CORRELATION

Variable 1	Variable 2	CORRELATION
OBS_60_CNT_SOCIAL_CIRCLE	OBS_30_CNT_SOCIAL_CIRCLE	0.998332863
AMT_GOODS_PRICE	AMT_CREDIT	0.986704386
REGION_RATING_CLIENT_W_CITY	REGION_RATING_CLIENT	0.950710179
CNT_FAM_MEMBERS	CNT_CHILDREN	0.880453292
REG_REGION_NOT_WORK_REGION	LIVE_REGION_NOT_WORK_REGION	0.857141677
DEF_30_CNT_SOCIAL_CIRCLE	DEF_60_CNT_SOCIAL_CIRCLE	0.856262392
LIVE_CITY_NOT_WORK_CITY	REG_CITY_NOT_WORK_CITY	0.821583789
AMT_GOODS_PRICE	AMT_ANNUITY	0.774134041
FLAG_EMP_PHONE	DAYS_BIRTH	0.617702887
DAYS_EMPLOYED	FLAG_DOCUMENT_6	0.591672763

## INSIGHTS AND RESULTS

- The dataset needed to be cleaned for empty cells, but it was free of duplicate rows.
- Many important columns contained outliers, which significantly affected the average.
- The data was highly imbalanced with respect to non-defaulters.
- During the univariate and bivariate analysis, we created numerous graphs that provided valuable insights into how individuals' loan-seeking behavior depended on age, gender, occupation, income, and education.
- We also discovered the correlation between many variables.

# LINKS

- [Excel](#)
- [Video Presentaion](#)