

Name: Naman Rawal

Email: rawalnaman49@gmail.com

Bank Loan Case Study - Final Project

1. Project Description

This project involves analyzing the bank loan dataset to identify key factors that influence loan default. As a data analyst, the objective is to use Exploratory Data Analysis (EDA) to analyze patterns in customer and loan attributes and help the bank make informed loan approval decisions.

The two main risks faced by the company are:

- Losing business by rejecting capable applicants.
 - Suffering financial loss by approving applicants who cannot repay loans.
-

2. Approach

We started by cleaning the dataset, identifying missing values, and handling outliers. We conducted various analyses, including univariate, bivariate, and correlation analysis, to explore the relationships between customer attributes and loan defaults. All work was done using Microsoft Excel.

3. Tech Stack Used

- **Microsoft Excel 2007:** Used for data cleaning, analysis, and visualization.
 - **Pivot Tables:** Used to segment data for deeper insights.
 - **Excel Functions:** Functions like COUNTIF, MEDIAN, QUARTILE, CORREL, and more were utilized.
-

4. Data Analytics Tasks

A. Missing Data Identification

- **Task:** Identify and handle missing data.
- **Approach:** Used ISBLANK and COUNT functions to detect missing data. Imputation was done using AVERAGE for numerical values and MODE for categorical values where appropriate.

These are the columns which has null values more than or equal to 50%. These columns need to be dropped.

Column name	no_of null values	Percentage_of null_values
COMMONAREA_AVG	34960	70%
COMMONAREA_MODE	34960	70%
COMMONAREA_MEDI	34960	70%
NONLIVINGAPARTMENTS_AVG	34714	69%
NONLIVINGAPARTMENTS_MODE	34714	69%
NONLIVINGAPARTMENTS_MEDI	34714	69%
LIVINGAPARTMENTS_AVG	34226	68%
LIVINGAPARTMENTS_MODE	34226	68%
LIVINGAPARTMENTS_MEDI	34226	68%
FONDKAPREMONT_MODE	34191	68%
FLOORSMIN_AVG	33894	68%
FLOORSMIN_MODE	33894	68%
FLOORSMIN_MEDI	33894	68%
YEARS_BUILD_AVG	33239	66%
YEARS_BUILD_MODE	33239	66%
YEARS_BUILD_MEDI	33239	66%
OWN_CAR_AGE	32949	66%
LANDAREA_AVG	29721	59%
LANDAREA_MODE	29721	59%
LANDAREA_MEDI	29721	59%
BASEMENTAREA_AVG	29199	58%
BASEMENTAREA_MODE	29199	58%
BASEMENTAREA_MEDI	29199	58%
EXT_SOURCE_1	28172	56%
NONLIVINGAREA_AVG	27572	55%
NONLIVINGAREA_MODE	27572	55%
NONLIVINGAREA_MEDI	27572	55%
ELEVATORS_AVG	26651	53%
ELEVATORS_MODE	26651	53%
ELEVATORS_MEDI	26651	53%
WALLSMATERIAL_MODE	25459	51%
APARTMENTS_AVG	25385	51%
APARTMENTS_MODE	25385	51%
APARTMENTS_MEDI	25385	51%
ENTRANCES_AVG	25195	50%
ENTRANCES_MODE	25195	50%
ENTRANCES_MEDI	25195	50%
LIVINGAREA_AVG	25137	50%
LIVINGAREA_MODE	25137	50%
LIVINGAREA_MEDI	25137	50%
HOUSETYPE_MODE	25075	50%
FLOORSMAX_AVG	24875	50%
FLOORSMAX_MODE	24875	50%
FLOORSMAX_MEDI	24875	50%

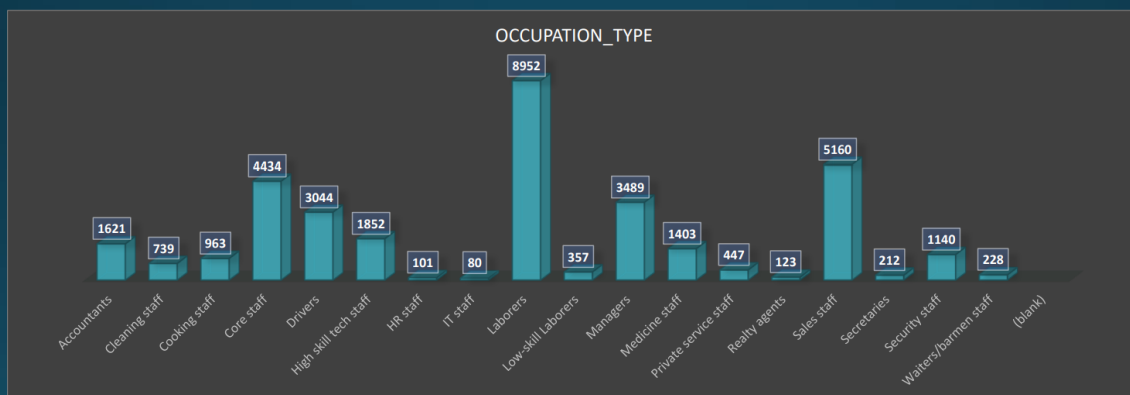
These are the columns which have irrelevant data for analysis. These columns need to be dropped.

Column name	no_of_null_values	Percentage_of_null_values
FLAG_MOBIL	0	0%
FLAG_EMP_PHONE	0	0%
FLAG_WORK_PHONE	0	0%
FLAG_CONT_MOBILE	0	0%
FLAG_PHONE	0	0%
FLAG_EMAIL	0	0%
CNT_FAM_MEMBERS	1	0%
REGION_RATING_CLIENT	0	0%
REGION_RATING_CLIENT_W_CITY	0	0%
EXT_SOURCE_2	126	0%
EXT_SOURCE_3	9944	20%
YEARS_BEGINEXPLUATATION_AVG	24394	49%
YEARS_BEGINEXPLUATATION_MODE	24394	49%
YEARS_BEGINEXPLUATATION_MEDI	24394	49%
TOTALAREA_MODE	24148	48%
EMERGENCYSTATE_MODE	23698	47%
DAYS_LAST_PHONE_CHANGE	1	0%
FLAG_DOCUMENT_2	0	0%
FLAG_DOCUMENT_3	0	0%
FLAG_DOCUMENT_4	0	0%
FLAG_DOCUMENT_5	0	0%
FLAG_DOCUMENT_6	0	0%
FLAG_DOCUMENT_7	0	0%
FLAG_DOCUMENT_8	0	0%
FLAG_DOCUMENT_9	0	0%
FLAG_DOCUMENT_10	0	0%
FLAG_DOCUMENT_11	0	0%
FLAG_DOCUMENT_12	0	0%
FLAG_DOCUMENT_13	0	0%
FLAG_DOCUMENT_14	0	0%
FLAG_DOCUMENT_15	0	0%
FLAG_DOCUMENT_16	0	0%
FLAG_DOCUMENT_17	0	0%
FLAG_DOCUMENT_18	0	0%
FLAG_DOCUMENT_19	0	0%
FLAG_DOCUMENT_20	0	0%
FLAG_DOCUMENT_21	0	0%



Mode Imputations:-

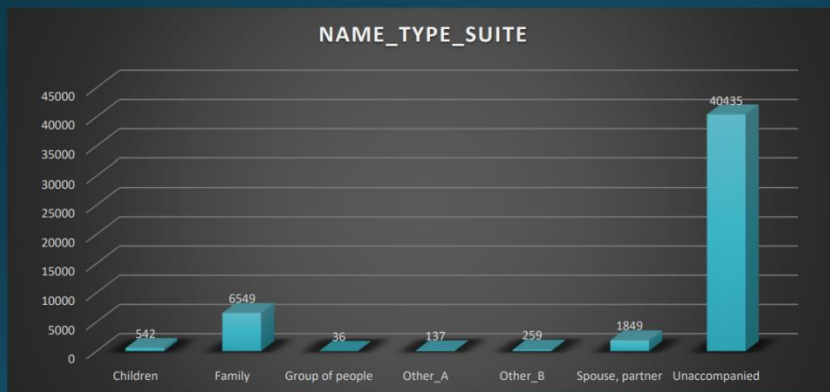
1. OCCUPATION_TYPE



Most Occurring Variable is Laborers. We will replace blanks with 8952.

Mode Imputations:-

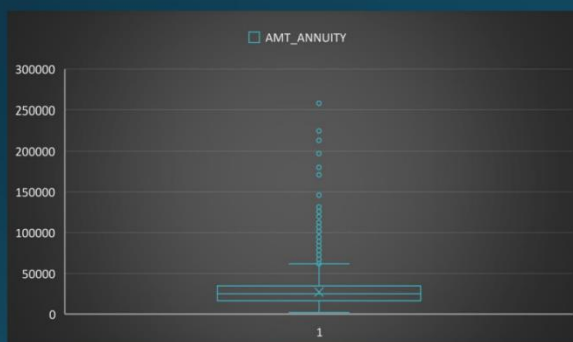
2. NAME_TYPE_SUITE



Most Occurring Variable is Unaccompanied.

Median Imputations:-

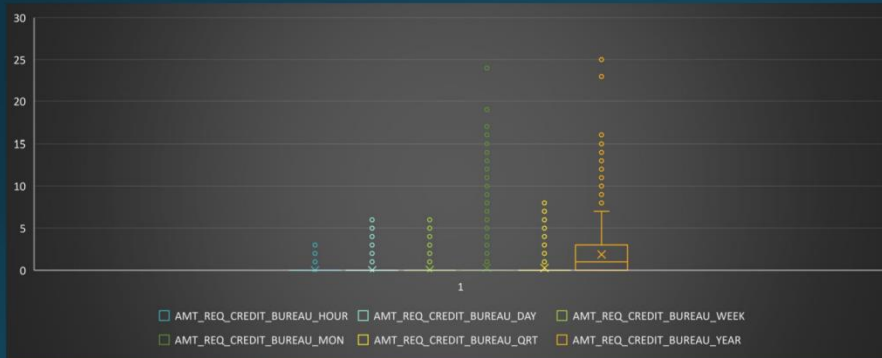
1. AMT_ANNUIITY



2. AMT_GOODS_PRICE

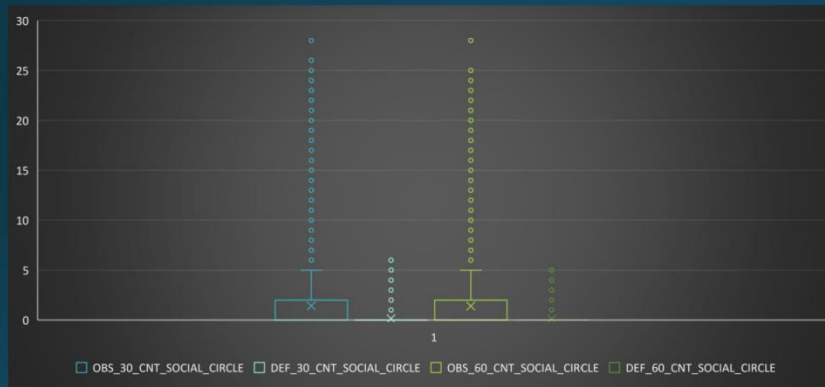


Median Imputations:-



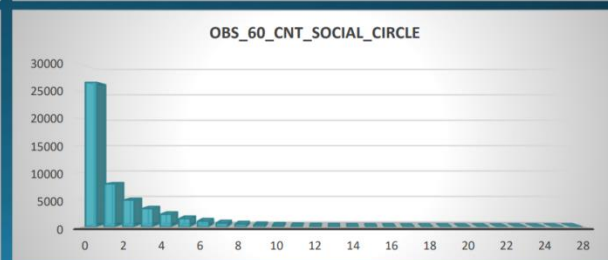
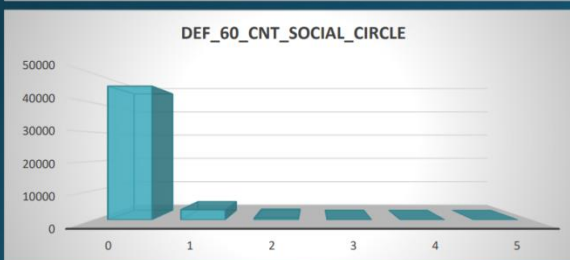
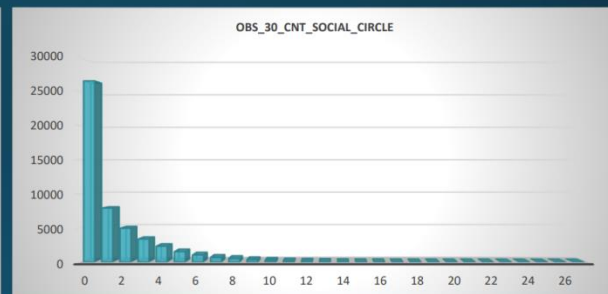
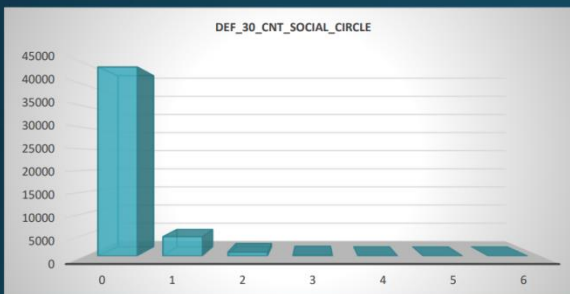
3. AMT_REQ_CREDIT_BUREAU_HOUR
4. AMT_REQ_CREDIT_BUREAU_DAY
5. AMT_REQ_CREDIT_BUREAU_WEEK
6. AMT_REQ_CREDIT_BUREAU_MON
7. AMT_REQ_CREDIT_BUREAU_QRT
8. AMT_REQ_CREDIT_BUREAU_YEAR

Median/Mode Imputations:-



1. DEF_30_CNT_SOCIAL_CIRCLE
2. OBS_30_CNT_SOCIAL_CIRCLE
3. DEF_60_CNT_SOCIAL_CIRCLE
4. OBS_60_CNT_SOCIAL_CIRCLE

Median/Mode Imputations:-



Previous_application datasets

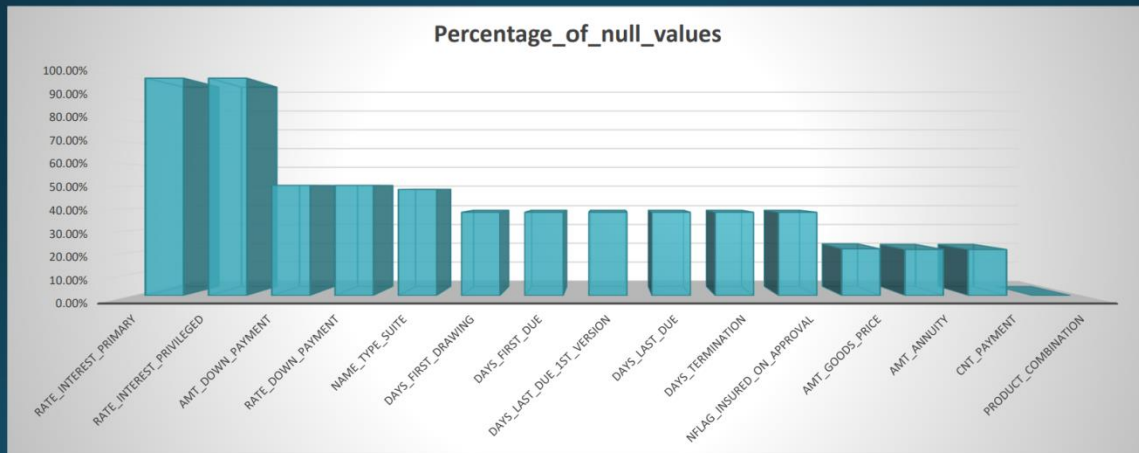
Column name	no_of_null_values	Percentage_of_null_values
RATE_INTEREST_PRIMARY	49833	99.67%
RATE_INTEREST_PRIVILEGED	49833	99.67%
AMT_DOWN_PAYMENT	25197	50.40%
RATE_DOWN_PAYMENT	25197	50.40%

These are the columns which has null values more than or equal to 50%. These columns need to be dropped.

Column name	no_of_null_values	Percentage_of_null_values
NAME_TYPE_SUITE	24243	48.49%
PRODUCT_COMBINATION	8	0.02%
WEEKDAY_APPR_PROCESS_START	0	0.00%
HOUR_APPR_PROCESS_START	0	0.00%
FLAG_LAST_APPL_PER_CONTRACT	0	0.00%
NFLAG_LAST_APPL_IN_DAY	0	0.00%

These are the columns which have irrelevant data for analysis. These columns need to be dropped.

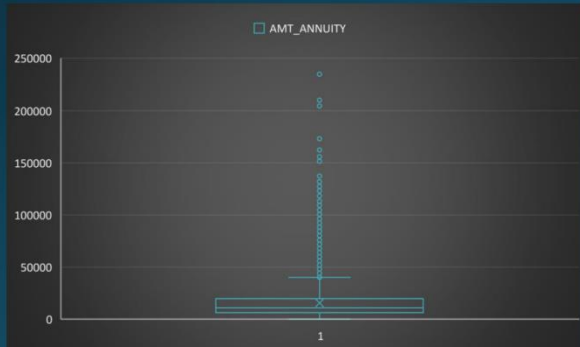
Previous_application datasets



Previous_application datasets

Median Imputations:-

1. AMT_ANNUITY



2. AMT_GOODS_PRICE

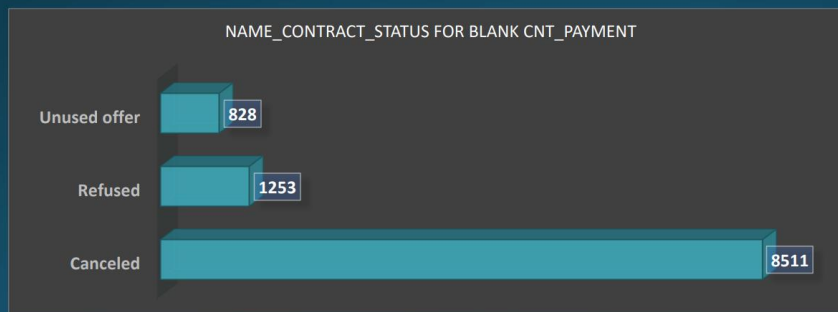


Previous_application datasets

Custom Imputations:-

1. CNT_PAYMENT

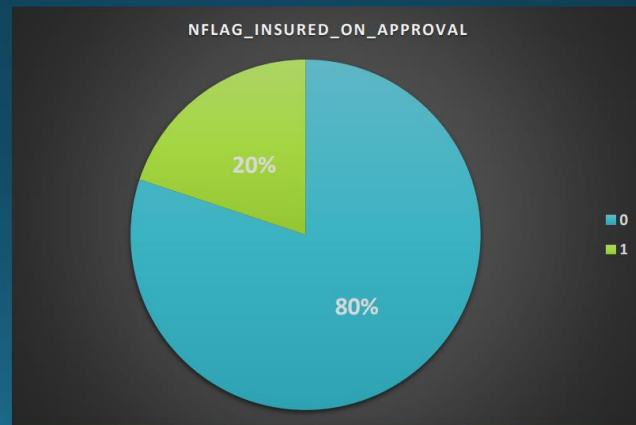
Most of Blank cells of cnt_payments have contract_status as canceled, refused, unused offer. So it makes more sense replacing them with 0 rather than Mean or Median.



Previous_application datasets

Mode Imputations:-

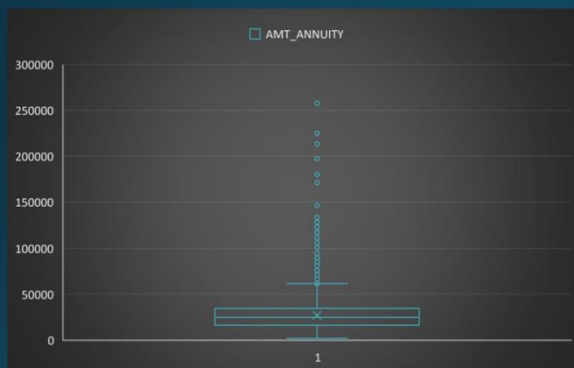
1. NFLAG_INSURED_ON_APPROVAL



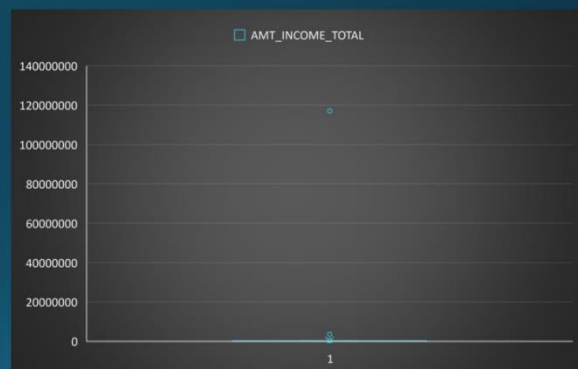
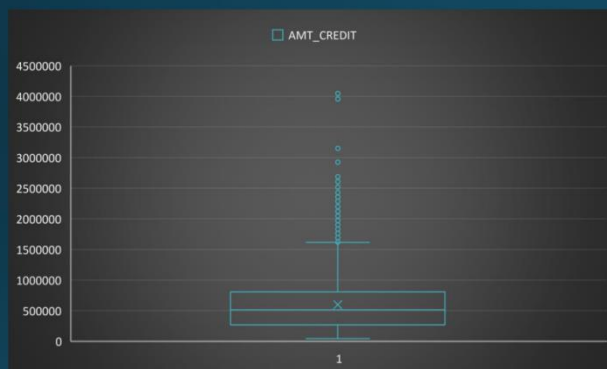
B. Identifying Outliers

- **Task:** Detect outliers in numerical variables.
- **Approach:** Used QUARTILE and IQR methods to detect outliers. Data points outside $1.5 * IQR$ were considered outliers.

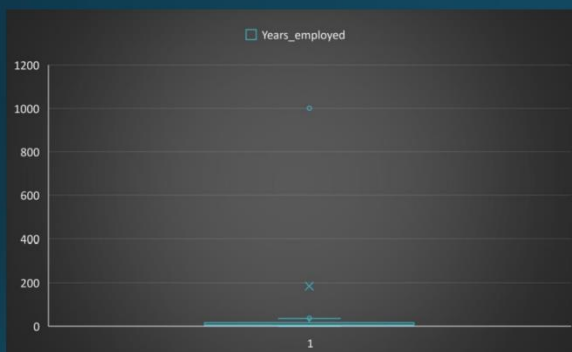
AMT_INCOME_TOTAL	
Quartile 1	112500
Quartile 3	202500
IQR	90000
Upper Limit	337500
Lower Limit	-22500



In the chart we can see there are few outliers in columns like AMT_ANNUITY and AMT_GOODS_PRICE.



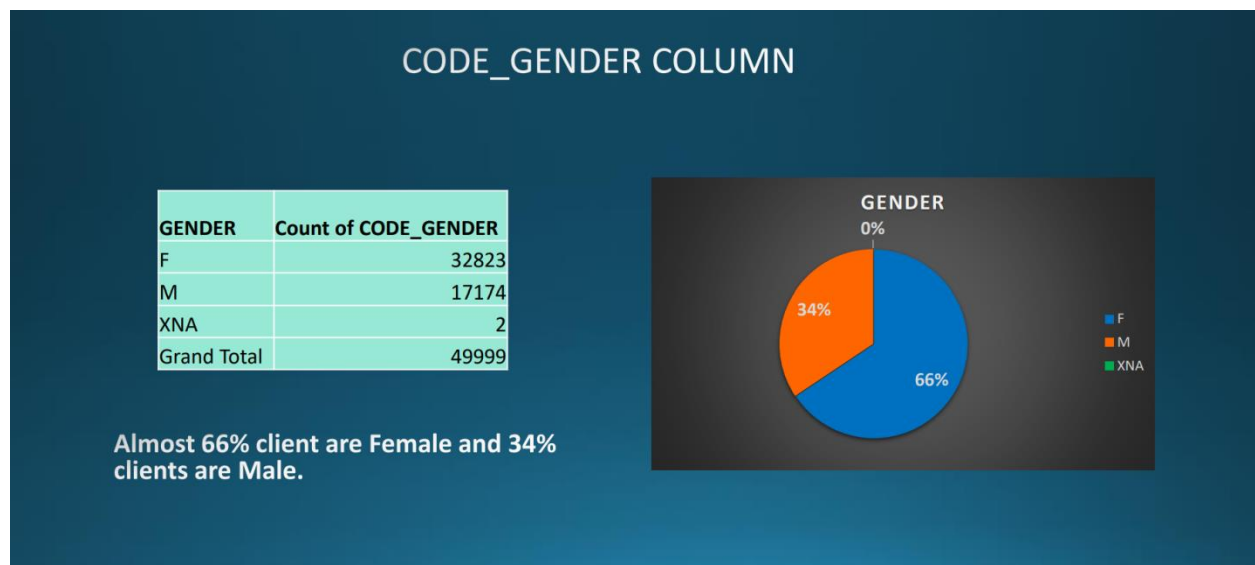
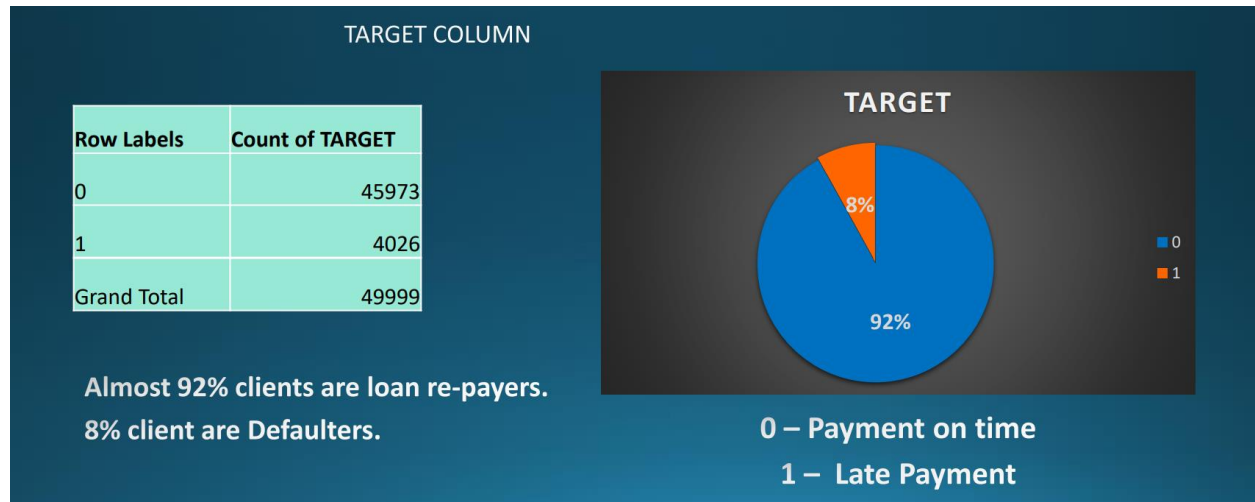
There are few outliers in columns like AMT_CREDIT and AMT_INCOME_TOTAL where amount is higher than normal. In AMT_INCOME_TOTAL one of extreme outlier is 1170000000 but we will not remove because income of person varies. We will not remove outlier from AMT_CREDIT too.



In Column Years_employed we can see people being employed for 1001 yrs which is not possible. Column CNT_CHILDREN shows people are having 11 children which is impossible in today's age

C. Data Imbalance

- **Task:** Analyze the distribution of the target variable (loan default) to check for class imbalance.
- **Approach:** Used COUNTIF to calculate class proportions. Found an imbalance in the data, with a higher proportion of successful loans compared to defaults.



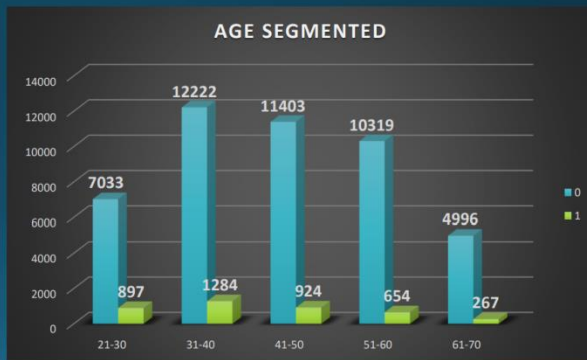
D. Univariate, Segmented Univariate, and Bivariate Analysis

- **Task:** Analyze individual variables and their relationships with loan default.

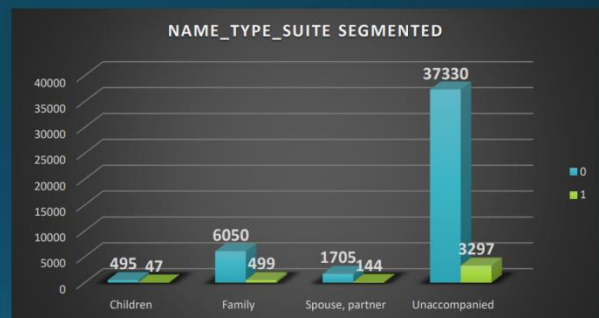
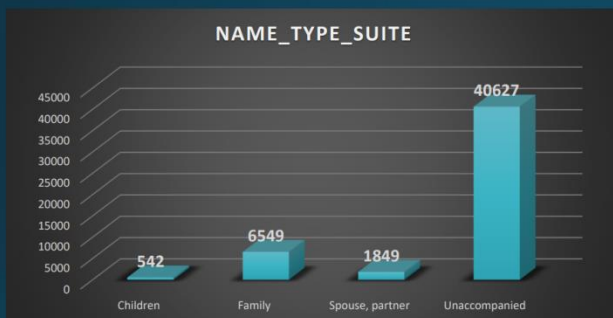
- **Univariate Analysis:** Created histograms to understand the distribution of numerical variables such as loan amount, income, and payment duration.
- **Segmented Univariate Analysis:** Used pivot tables to segment data by customer type (defaulted vs non-defaulted).
- **Bivariate Analysis:** Analyzed the relationships between loan amount, income, and default status using scatter plots.

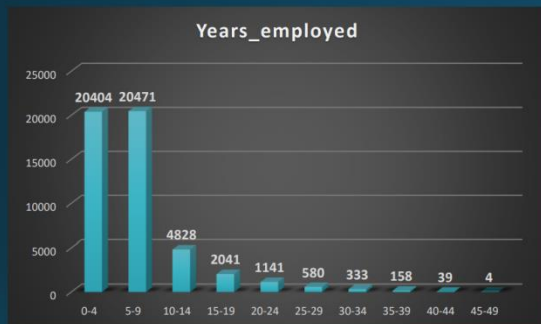


Majority of the Clients are in the age group 31-40.



we can see as age increases , chances of defaulter decreases.

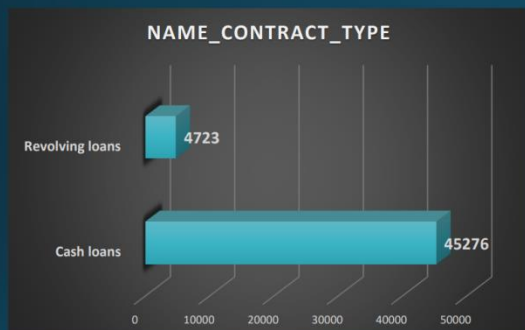




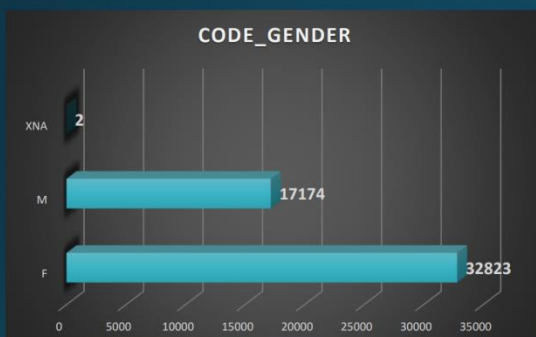
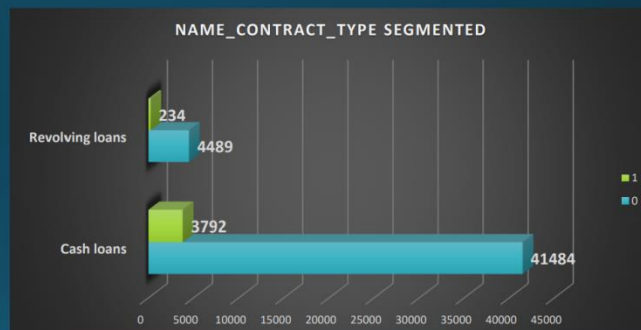
Majority of the Clients are having 0-9 years of experience.



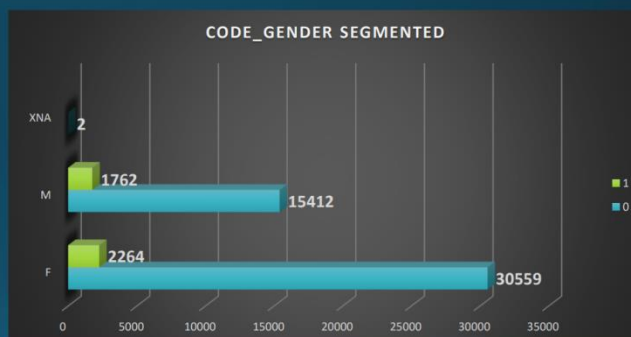
we can see as experience increases , chances of defaulting decreases.

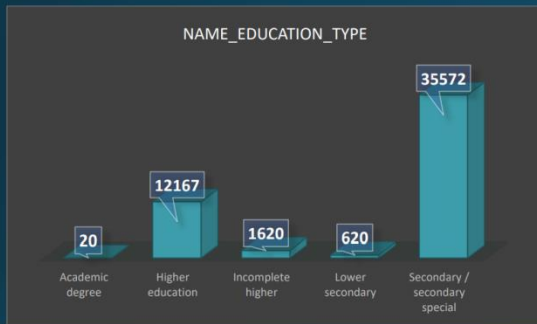


Majority of the Clients are taking Cash loans.

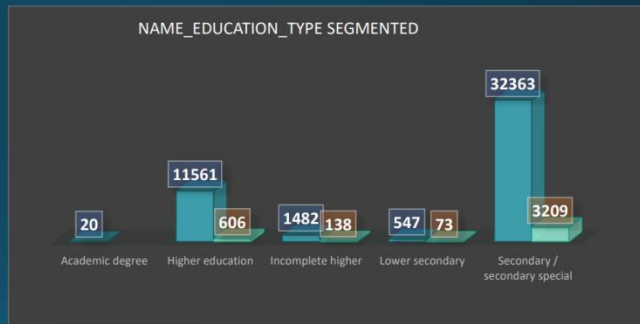


Male are less defaulters compared to Female.

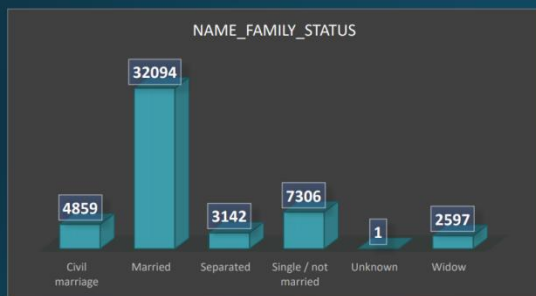




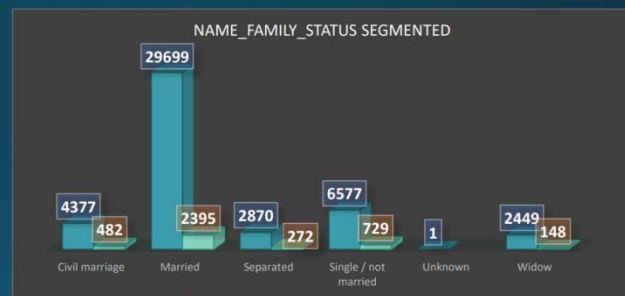
The numbers of loans taken by Clients with Secondary special Education is the highest and Academic degree is the lowest



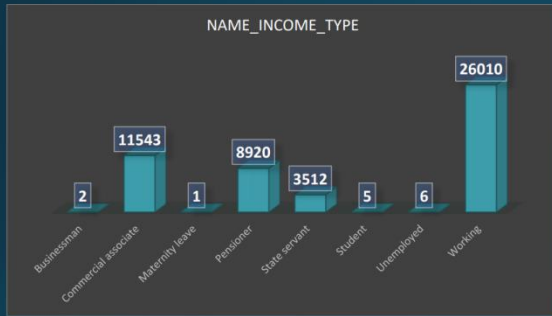
Least default: Academic degree
Highest default: Secondary special



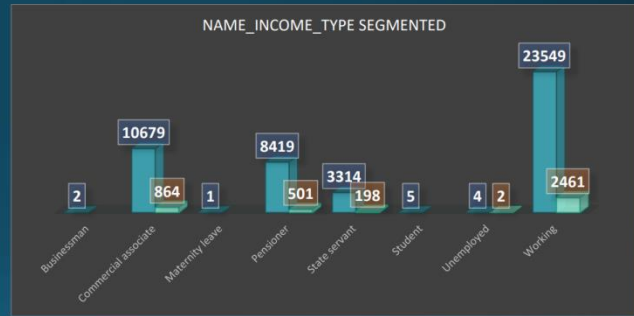
The number of loans taken by Married clients are the highest and clients who are widows are the least if we ignore unknown.



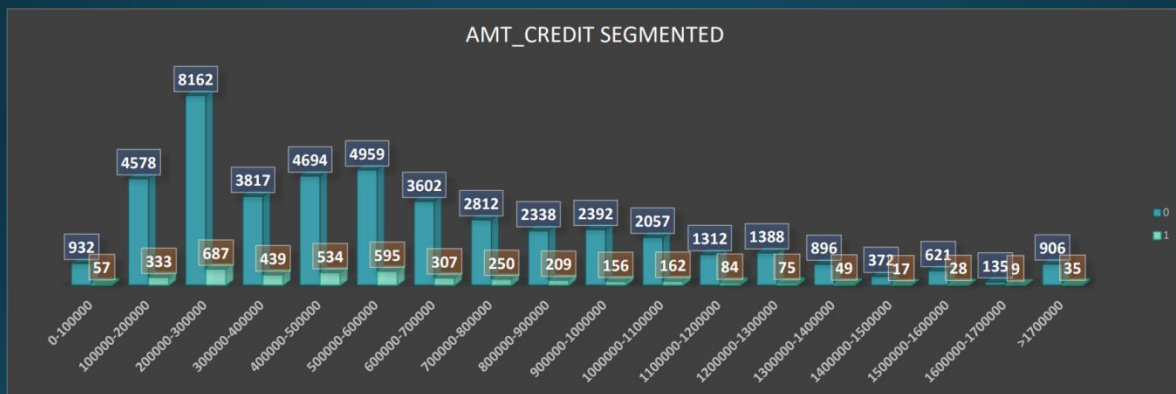
Least Defaulter: Widow
Highest Defaulter : Married



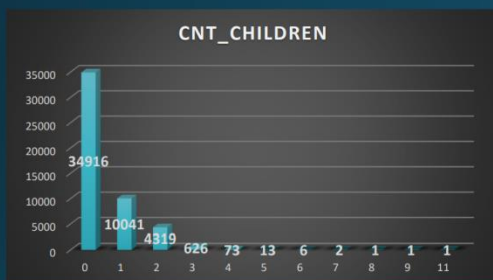
Bank target those groups whose income type is working.



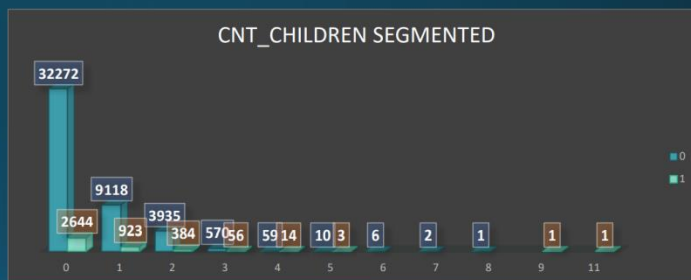
Least default: Client who is Businessman or student or at Maternity leave.
Highest default: Client who is working



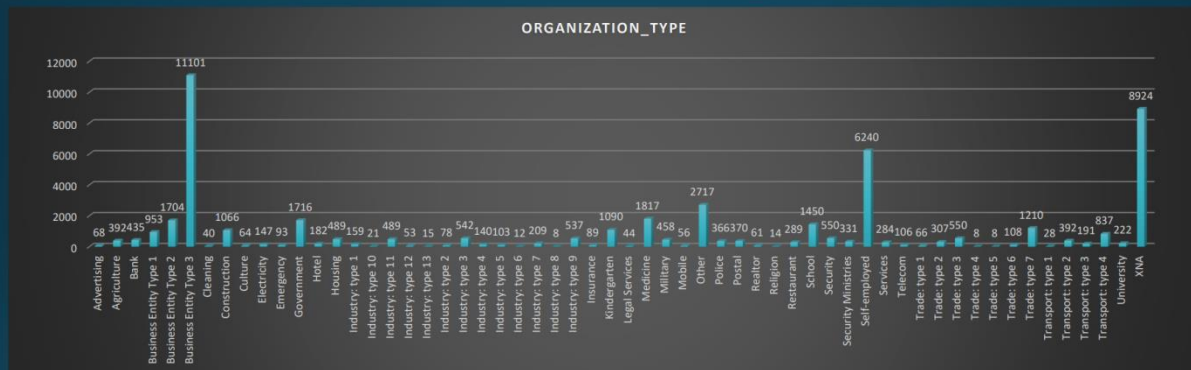
Majority of the Clients took the loan between 2L – 3L.



The highest number of loans are taken by Clients who does not have a child

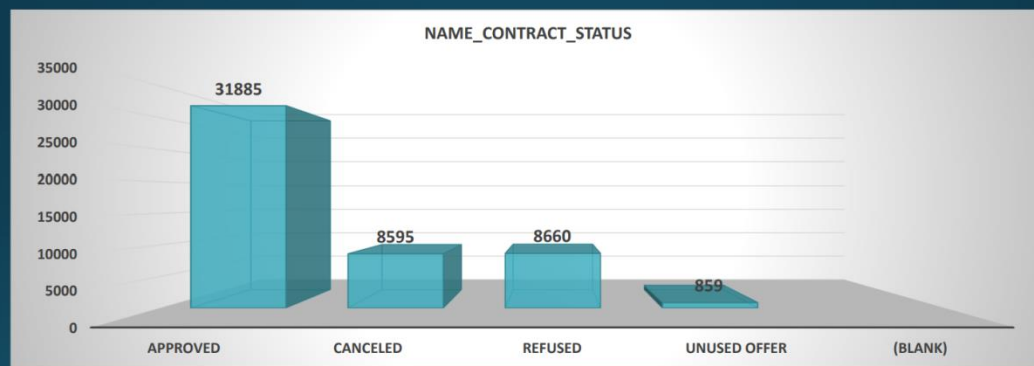


As number of children increases, number of client who took loan decreases.

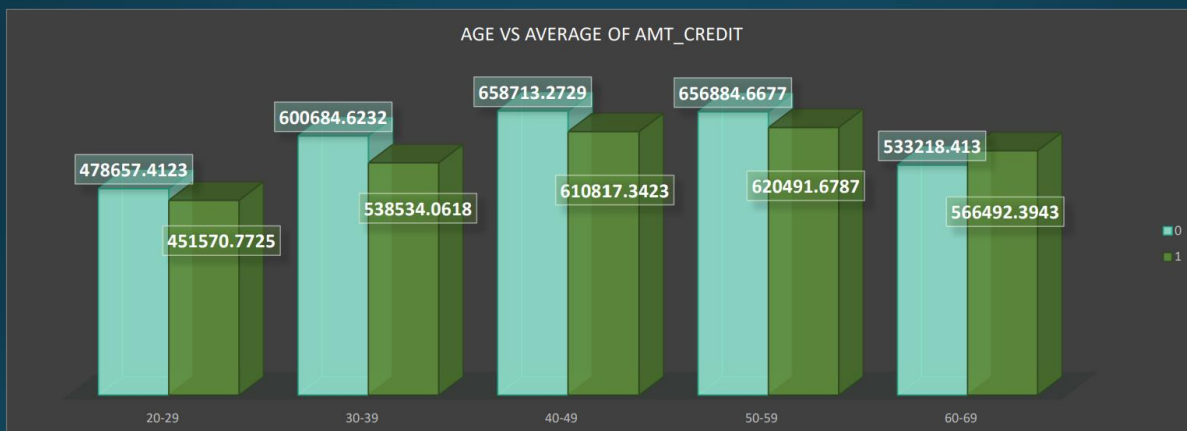


Clients who are working in business Entity type of Organization took the highest number of loans.

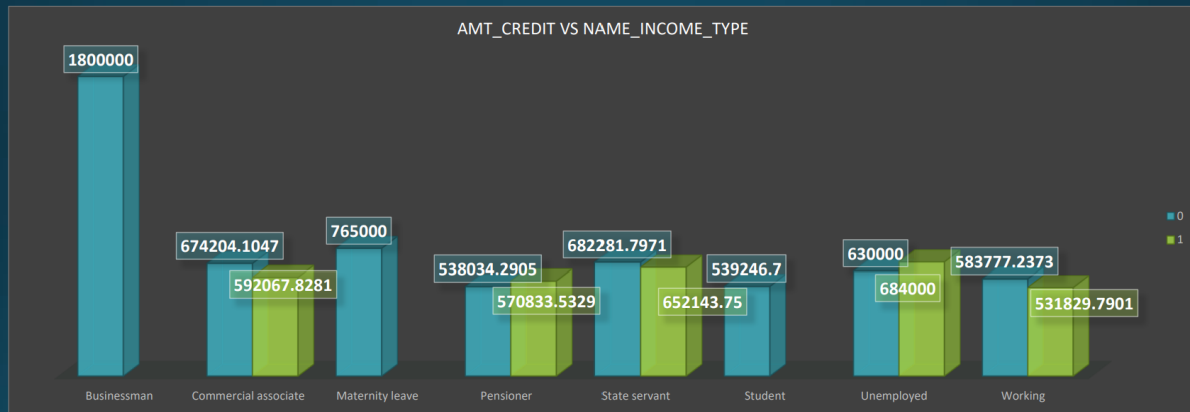
Previous_application datasets



More number of Clients were approved for loans previously.



Age group 40-49 took the highest amount of loan but age group 50-59 are defaulter with highest amount of loan.



As we see Businessman took the highest amount of loan and did the payment on time. Clients who are unemployed have highest amount of loan which they didn't repay on time.

E. Top Correlations for Different Scenarios

- **Task:** Identify the strongest correlations between variables and loan default for customers with and without payment difficulties.
- **Approach:** Used CORREL to calculate correlations between variables such as income, loan amount, and repayment status.

Top Correlation Coefficients for Payment difficulties are:-

Correlation between Columns	Value
AMT_CREDIT - AMT_GOODS_PRICE	0.982267963
OBS_60_CNT_SOCIAL_CIRCLE - OBS_30_CNT_SOCIAL_CIRCLE	0.998065853
DEF_60_CNT_SOCIAL_CIRCLE - DEF_30_CNT_SOCIAL_CIRCLE	0.89051161
REG_REGION_NOT_WORK_REGION - LIVE_REGION_NOT_WORK_REGION	0.806743886
REG_CITY_NOT_WORK_CITY - LIVE_CITY_NOT_WORK_CITY	0.783754676
AMT_CREDIT - AMT_ANNUITY	0.749665201
AMT_GOODS_PRICE - AMT_ANNUITY	0.74950403

Correlation between Columns	Value
OBS_60_CNT_SOCIAL_CIRCLE - OBS_30_CNT_SOCIAL_CIRCLE	0.998357563
AMT_GOODS_PRICE - AMT_CREDIT	0.986051701
LIVE_REGION_NOT_WORK_REGION - REG_REGION_NOT_WORK_REGION	0.861374946
DEF_60_CNT_SOCIAL_CIRCLE - DEF_30_CNT_SOCIAL_CIRCLE	0.850995792
REG_CITY_NOT_WORK_CITY - LIVE_CITY_NOT_WORK_CITY	0.825358079
AMT_ANNUITY - AMT_GOODS_PRICE	0.774006842
AMT_ANNUITY - AMT_CREDIT	0.770772818

Correlation between Columns	Value
OBS_60_CNT_SOCIAL_CIRCLE - OBS_30_CNT_SOCIAL_CIRCLE	0.998357563
AMT_GOODS_PRICE - AMT_CREDIT	0.986051701
LIVE_REGION_NOT_WORK_REGION - REG_REGION_NOT_WORK_REGION	0.861374946
DEF_60_CNT_SOCIAL_CIRCLE - DEF_30_CNT_SOCIAL_CIRCLE	0.850995792
REG_CITY_NOT_WORK_CITY - LIVE_CITY_NOT_WORK_CITY	0.825358079
AMT_ANNUITY - AMT_GOODS_PRICE	0.774006842
AMT_ANNUITY - AMT_CREDIT	0.770772818

[illegible]

- Majority of clients are responsible loan payers.
- The bank tends to grant more loans to women, though men have a lower default rate compared to women.
- As both age and work experience increase, the likelihood of defaulting on a loan decreases.
- Cash loans are the most common type of loan among clients.
- Clients with higher education levels tend to default less frequently than those with lower levels of education, such as those with secondary or specialized education. The bank should prioritize lending to more educated individuals.
- The number of loan applicants tends to decrease as the number of children they have increases.
- The bank needs to exercise greater caution when providing loans to unemployed clients, as they represent the highest default rates with the largest loan amounts.
- Older clients typically take out larger loans, but they also have a lower default rate, making them a safer and more profitable demographic for the bank.

The analysis successfully identified key factors that influence loan defaults. These insights can help the bank implement better risk management practices, such as offering higher interest rates to riskier customers or denying loans to those likely to default.

7. Hyperlinks

- **Excel Sheet:**

https://docs.google.com/spreadsheets/d/1r50pZlWwjCJCFqI6h13qX_M0kG05CTbG/edit?usp=sharing&ouid=107712337603641298783&rtpof=true&sd=true
