# Bank Loan Case Study

## Project Description:

The **objective** of this Bank Loan Case Study is to use Exploratory Data Analysis (EDA) to identify key factors that influence loan default, enabling the finance company to make better decisions on loan approvals. Specifically, the goals are to:

1. **Reduce Financial Risk**: Identify customers who are likely to default on their loans so the company can either reject their applications, reduce loan amounts, or adjust the interest rates to mitigate risks.

2. **Maximize Business Opportunities**: Ensure that capable applicants are not wrongly rejected, allowing the company to approve loans for clients who are financially stable, thus maximizing profit.

3. **Understand Patterns**: Analyze customer and loan attributes to detect patterns related to loan default, such as income level, loan amount, credit history, and annuity payments.

4. **Improve Decision-Making**: Provide insights into which variables are the strongest predictors of loan default to enhance the company's loan approval process, reducing both approval risks and missed opportunities.

By achieving these objectives, the company aims to improve overall lending performance while maintaining financial health.

## Approach:

To achieve the objective of identifying patterns that influence loan default, the approach will involve structured steps focusing on data exploration, analysis, and insights extraction using Exploratory Data Analysis (EDA).

**Step 1: Data Understanding and Preparation**

1. **Data Collection**:
   o Obtain the loan application dataset, which includes customer and loan attributes such as income, loan amount, loan status (TARGET), and payment history.

2. **Data Inspection**:
   o Examine the dataset's structure, types of variables (categorical and numerical), and content.

- Ensure familiarity with key variables, particularly the TARGET variable that indicates loan default (1 for default, 0 for non-default).

3. **Data Cleaning**:
   - **Handle Missing Data**: Identify and deal with missing values using techniques like:
     - Removing rows or columns with excessive missing values.
     - Imputing missing data using averages, medians, or domain-specific values.
   - **Outlier Detection**: Identify outliers using statistical methods (e.g., interquartile range (IQR), z-scores) and assess whether they need to be removed or treated.

   **Tools**: Excel functions like COUNTIF, IF, AVERAGE, MEDIAN, QUARTILE, and conditional formatting.

**Step 2: Data Exploration (EDA)**

1. **Univariate Analysis**:
   - **Goal**: Understand the distribution of individual variables (e.g., income, loan amount).
   - **Techniques**:
     - Use descriptive statistics (mean, median, mode, standard deviation).
     - Create histograms and bar charts to visualize the distribution of numerical and categorical variables.

2. **Segmented Univariate Analysis**:
   - **Goal**: Compare variable distributions for clients with payment difficulties (TARGET = 1) and without payment difficulties (TARGET = 0).
   - **Techniques**:
     - Use pivot tables and filtering to create comparisons between the two segments.
     - Visualize the differences using stacked or grouped bar charts.

3. **Bivariate Analysis**:
   - **Goal**: Explore relationships between variables and their impact on loan default.
   - **Techniques**:
     - Correlation analysis for numerical variables (e.g., correlation between income and loan amount).
     - Scatter plots or heatmaps to visualize relationships between variables.
     - Analyze categorical variables using cross-tabulations.

**Step 3: Data Segmentation**

1. **Segment Data by Loan Status**:
   - o Create segments based on TARGET (loan default or non-default) and analyze each segment separately.
   - o Calculate the proportion of each class (TARGET = 0 and TARGET = 1) to identify potential data imbalance.

2. **Identify Correlations within Segments**:
   - o Perform correlation analysis within each segment to determine which variables (e.g., AMT_CREDIT, AMT_ANNUITY) have the strongest relationship with loan default.
   - o Rank the variables by correlation strength to identify the top indicators of default for each segment.

**Step 4: Analyze Data Imbalance**

1. **Goal**: Assess if there is an imbalance between the number of default and non-default cases.

2. **Techniques**:
   - o Calculate the ratio of default to non-default cases using COUNTIF.
   - o Visualize the class distribution using pie charts or bar charts.

**Step 5: Insights and Interpretation**

1. **Top Correlations**:
   - o Identify the variables that have the highest correlation with loan default in both segments (default and non-default).
   - o Focus on strong positive and negative correlations to highlight key indicators of financial risk (e.g., loan amount, income level).

2. **Outlier Investigation**:
   - o Investigate whether outliers (e.g., very high or low incomes) significantly affect default rates.

**Step 6: Visualization**

1. **Correlation Heatmaps**:

- Create heatmaps to visualize correlations across variables in both segments, highlighting strong correlations.

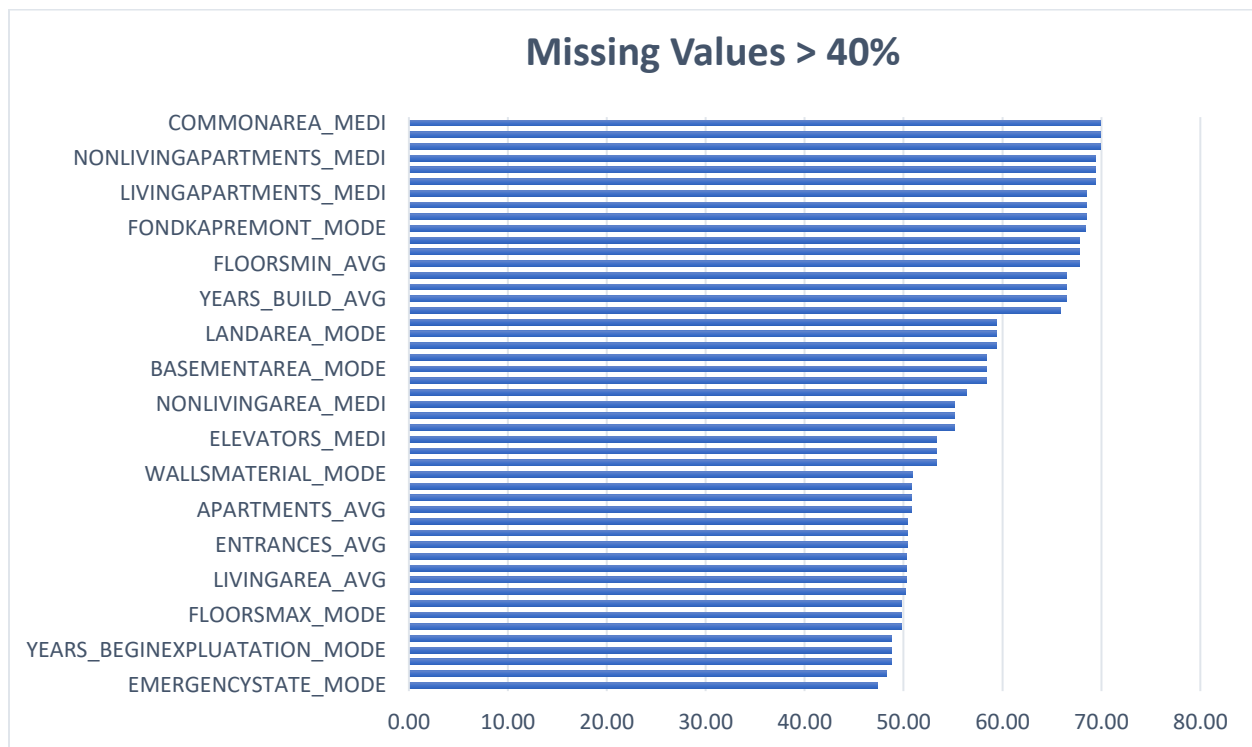2. **Bar Charts and Box Plots**:
   - Visualize the distribution of key variables like loan amount, income, and loan status using bar charts or box plots to gain a better understanding of their influence on default.

# Tech-Stack Used:

1. Microsoft excels

# Insights:

**A.     Identify Missing Data and Deal with it Appropriately:** Identified the missing data in the dataset and handle it with using excel function COUNT, ISBLANK, and IF. Also Perform imputation using excel function AVERAGE or MEDIAN.

| Column_Name | Blank_Percentage |
|---|---|
| EMERGENCYSTATE_MODE | 47.40 |
| TOTALAREA_MODE | 48.30 |
| YEARS_BEGINEXPLUATATION_AVG | 48.79 |
| YEARS_BEGINEXPLUATATION_MODE | 48.79 |
| YEARS_BEGINEXPLUATATION_MEDI | 48.79 |
| FLOORSMAX_AVG | 49.75 |
| FLOORSMAX_MODE | 49.75 |
| FLOORSMAX_MEDI | 49.75 |
| HOUSETYPE_MODE | 50.15 |
| LIVINGAREA_AVG | 50.28 |
| LIVINGAREA_MODE | 50.28 |
| LIVINGAREA_MEDI | 50.28 |
| ENTRANCES_AVG | 50.39 |
| ENTRANCES_MODE | 50.39 |
| ENTRANCES_MEDI | 50.39 |
| APARTMENTS_AVG | 50.77 |
| APARTMENTS_MODE | 50.77 |
| APARTMENTS_MEDI | 50.77 |
| WALLSMATERIAL_MODE | 50.92 |
| ELEVATORS_AVG | 53.30 |
| ELEVATORS_MODE | 53.30 |
| ELEVATORS_MEDI | 53.30 |
| NONLIVINGAREA_AVG | 55.15 |
| NONLIVINGAREA_MODE | 55.15 |
| NONLIVINGAREA_MEDI | 55.15 |
| EXT_SOURCE_1 | 56.35 |
| BASEMENTAREA_AVG | 58.40 |
| BASEMENTAREA_MODE | 58.40 |
| BASEMENTAREA_MEDI | 58.40 |
| LANDAREA_AVG | 59.44 |
| LANDAREA_MODE | 59.44 |
| LANDAREA_MEDI | 59.44 |
| OWN_CAR_AGE | 65.90 |
| YEARS_BUILD_AVG | 66.48 |
| YEARS_BUILD_MODE | 66.48 |
| YEARS_BUILD_MEDI | 66.48 |
| FLOORSMIN_AVG | 67.79 |
| FLOORSMIN_MODE | 67.79 |
| FLOORSMIN_MEDI | 67.79 |
| FONDKAPREMONT_MODE | 68.38 |
| LIVINGAPARTMENTS_AVG | 68.45 |
| LIVINGAPARTMENTS_MODE | 68.45 |
| LIVINGAPARTMENTS_MEDI | 68.45 |
| NONLIVINGAPARTMENTS_AVG | 69.43 |
| NONLIVINGAPARTMENTS_MODE | 69.43 |
| NONLIVINGAPARTMENTS_MEDI | 69.43 |
| COMMONAREA_AVG | 69.92 |
| COMMONAREA_MODE | 69.92 |
| COMMONAREA_MEDI | 69.92 |

**Fig. Missing Value >40**

| Column_Name | Blank_Percentage |
| --- | --- |
| SK_ID_CURR | 0.00 |
| TARGET | 0.00 |
| NAME_CONTRACT_TYPE | 0.00 |
| CODE_GENDER | 0.00 |
| FLAG_OWN_CAR | 0.00 |
| FLAG_OWN_REALTY | 0.00 |
| CNT_CHILDREN | 0.00 |
| AMT_INCOME_TOTAL | 0.00 |
| AMT_CREDIT | 0.00 |
| NAME_INCOME_TYPE | 0.00 |
| NAME_EDUCATION_TYPE | 0.00 |
| NAME_FAMILY_STATUS | 0.00 |
| NAME_HOUSING_TYPE | 0.00 |
| REGION_POPULATION_RELATIVE | 0.00 |
| DAYS_BIRTH | 0.00 |
| CUSTOMER_AGE | 0.00 |
| DAYS_EMPLOYED | 0.00 |
| EMPLOYMENT | 0.00 |
| DAYS_REGISTRATION | 0.00 |
| REGISTRATION_DAYS | 0.00 |
| DAYS_ID_PUBLISH | 0.00 |
| ID_PUBLISED_DAYS | 0.00 |
| FLAG_MOBIL | 0.00 |
| FLAG_EMP_PHONE | 0.00 |
| FLAG_WORK_PHONE | 0.00 |
| FLAG_CONT_MOBILE | 0.00 |
| FLAG_PHONE | 0.00 |
| FLAG_EMAIL | 0.00 |
| REGION_RATING_CLIENT | 0.00 |
| REGION_RATING_CLIENT_W_CITY | 0.00 |
| WEEKDAY_APPR_PROCESS_START | 0.00 |
| HOUR_APPR_PROCESS_START | 0.00 |
| REG_REGION_NOT_LIVE_REGION | 0.00 |
| REG_REGION_NOT_WORK_REGION | 0.00 |
| LIVE_REGION_NOT_WORK_REGION | 0.00 |
| REG_CITY_NOT_LIVE_CITY | 0.00 |
| REG_CITY_NOT_WORK_CITY | 0.00 |
| LIVE_CITY_NOT_WORK_CITY | 0.00 |
| ORGANIZATION_TYPE | 0.00 |
| FLAG_DOCUMENT_2 | 0.00 |
| FLAG_DOCUMENT_3 | 0.00 |
| FLAG_DOCUMENT_4 | 0.00 |
| FLAG_DOCUMENT_5 | 0.00 |
| FLAG_DOCUMENT_6 | 0.00 |
| FLAG_DOCUMENT_7 | 0.00 |
| FLAG_DOCUMENT_8 | 0.00 |
| FLAG_DOCUMENT_9 | 0.00 |
| FLAG_DOCUMENT_10 | 0.00 |
| FLAG_DOCUMENT_11 | 0.00 |
| FLAG_DOCUMENT_12 | 0.00 |
| FLAG_DOCUMENT_13 | 0.00 |
| FLAG_DOCUMENT_14 | 0.00 |
| FLAG_DOCUMENT_15 | 0.00 |
| FLAG_DOCUMENT_16 | 0.00 |
| FLAG_DOCUMENT_17 | 0.00 |
| FLAG_DOCUMENT_18 | 0.00 |
| FLAG_DOCUMENT_19 | 0.00 |
| FLAG_DOCUMENT_20 | 0.00 |
| FLAG_DOCUMENT_21 | 0.00 |
| AMT_ANNUITY | 0.00 |
| CNT_FAM_MEMBERS | 0.00 |
| DAYS_LAST_PHONE_CHANGE | 0.00 |
| AMT_GOODS_PRICE | 0.08 |
| EXT_SOURCE_2 | 0.25 |
| OBS_30_CNT_SOCIAL_CIRCLE | 0.34 |
| DEF_30_CNT_SOCIAL_CIRCLE | 0.34 |
| OBS_60_CNT_SOCIAL_CIRCLE | 0.34 |
| DEF_60_CNT_SOCIAL_CIRCLE | 0.34 |
| NAME_TYPE_SUITE | 0.38 |
| AMT_REQ_CREDIT_BUREAU_HOUR | 13.47 |
| AMT_REQ_CREDIT_BUREAU_DAY | 13.47 |
| AMT_REQ_CREDIT_BUREAU_WEEK | 13.47 |
| AMT_REQ_CREDIT_BUREAU_MON | 13.47 |
| AMT_REQ_CREDIT_BUREAU_QRT | 13.46826937 |
| AMT_REQ_CREDIT_BUREAU_YEAR | 13.46826937 |
| EXT_SOURCE_3 | 19.89 |
| OCCUPATION_TYPE | 31.31 |

**Fig. Missing Value<40%**

Fig. Missing value Handling

| Category | Count of missing value | Status Of Missing Data |
| --- | --- | --- |
| AMT_ANNUITY | 1 | Replace with median value 24939 |
| AMT_GOODS_PRICE | 38 | Calculate median and replace null values to that median |
| NAME_TYPE_SUITE | 192 | On the basis of majority "Unaccompanied", null value replaced with it |
| OWN_CAR_AGE | 32950 | Not found relevent data in previous loan sheet so it was deleted |
| OCCUPATION_TYPE | 15654 | replaced with "NA" |
| CNT_FAM_MEMBERS | 1 | Replaced with 0 (ZERO) |
| EXT_SOURCE_1 | 28172 | Calculate median and replace null values to that median |
| EXT_SOURCE_2 | 126 | Replace with median value |
| EXT_SOURCE_3 | 9944 | Replace with median value |
| APARTMENTS_AVG | 25385 | Replace with median value |
| BASEMENTAREA_AVG | 29199 | Replace with median value |
| YEARS_BEGINEXPLUATATION_AVG | 24394 | Replace with median value |
| YEARS_BUILD_AVG | 33239 | Replace with median value |
| COMMONAREA_AVG | 34960 | Replace with median value |
| ELEVATORS_AVG | 26651 | Replace with median value |
| ENTRANCES_AVG | 25195 | Replace with median value |
| FLOORSMAX_AVG | 24875 | Replace with median value |
| FLOORSMIN_AVG | 33894 | Replace with median value |
| LANDAREA_AVG | 29721 | Replace with median value |
| LIVINGAPARTMENTS_AVG | 34226 | Replace with median value |
| LIVINGAREA_AVG | 25137 | Replace with median value |
| NONLIVINGAPARTMENTS_AVG | 34714 | Replace with median value |
| NONLIVINGAREA_AVG | 27572 | Replace with median value |
| APARTMENTS_MODE | 25385 | Replace with median value |
| BASEMENTAREA_MODE | 29199 | Replace with median value |
| YEARS_BEGINEXPLUATATION_MODE | 24394 | Replace with median value |
| YEARS_BUILD_MODE | 33239 | Replace with median value |
| COMMONAREA_MODE | 34960 | Replace with median value |
| ELEVATORS_MODE | 26651 | Replace with median value |
| ENTRANCES_MODE | 25195 | Replace with median value |
| FLOORSMAX_MODE | 24875 | Replace with median value |
| FLOORSMIN_MODE | 33894 | Replace with median value |
| LANDAREA_MODE | 29721 | Replace with median value |
| LIVINGAPARTMENTS_MODE | 34226 | Replace with median value |
| LIVINGAREA_MODE | 25137 | Replace with median value |
| NONLIVINGAPARTMENTS_MODE | 34714 | Replace with median value |
| NONLIVINGAREA_MODE | 27572 | Replace with median value |
| APARTMENTS_MEDI | 25385 | Replace with median value |
| BASEMENTAREA_MEDI | 29199 | Replace with median value |
| YEARS_BEGINEXPLUATATION_MEDI | 24394 | Replace with median value |
| YEARS_BUILD_MEDI | 33239 | Replace with median value |
| COMMONAREA_MEDI | 34960 | Replace with median value |
| ELEVATORS_MEDI | 26651 | Replace with median value |
| ENTRANCES_MEDI | 25195 | Replace with median value |
| FLOORSMAX_MEDI | 24875 | Replace with median value |
| FLOORSMIN_MEDI | 33894 | Replace with median value |
| LANDAREA_MEDI | 29721 | Replace with median value |
| LIVINGAPARTMENTS_MEDI | 34226 | Replace with median value |
| LIVINGAREA_MEDI | 25137 | Replace with median value |
| NONLIVINGAPARTMENTS_MEDI | 34714 | Replace with median value |
| NONLIVINGAREA_MEDI | 27572 | Replace with median value |
| FONDKAPREMONT_MODE | 34191 | Column deleted |
| HOUSETYPE_MODE | 25075 | Column deleted |

**Fig. Missing value status**

**B.      Identify Outliers in the Dataset:** Identified the outliers using the IQR method, which is a common statistical technique. The IQR represents the middle 50% of the data, and outliers are typically values that fall below Q1 - 1.5 * IQR or above Q3 + 1.5 * IQR.
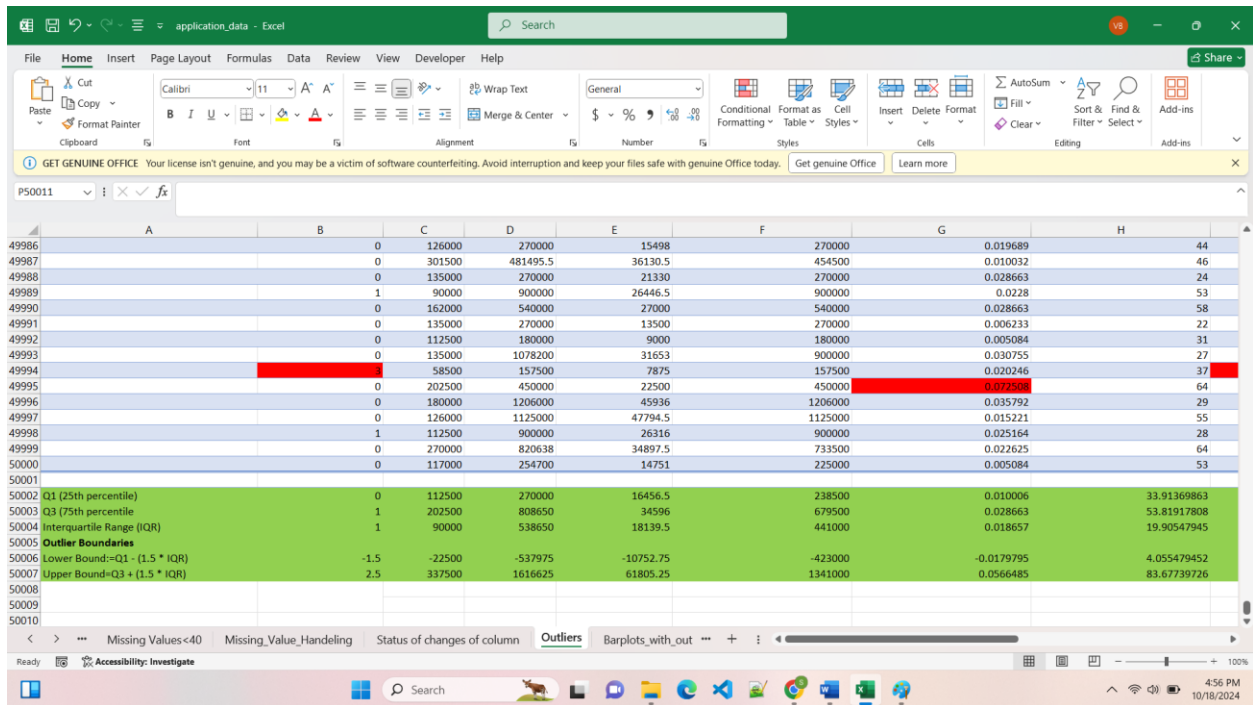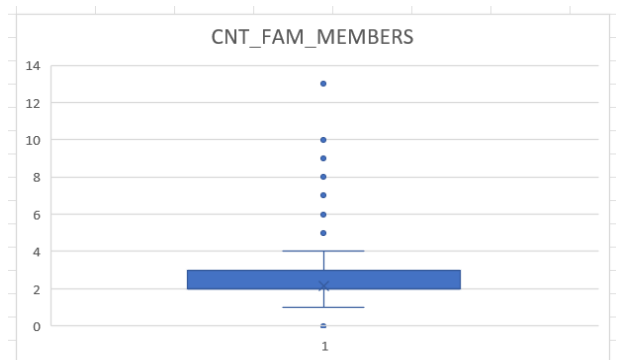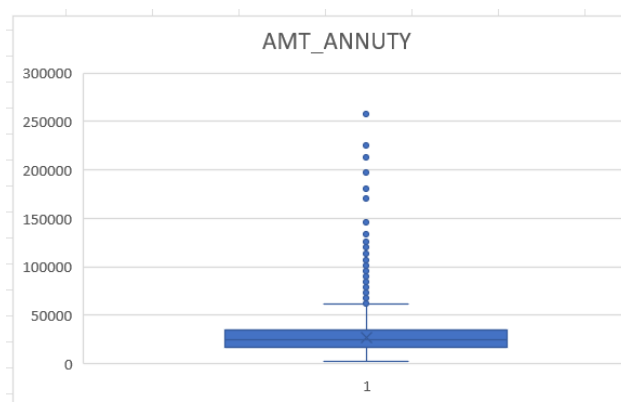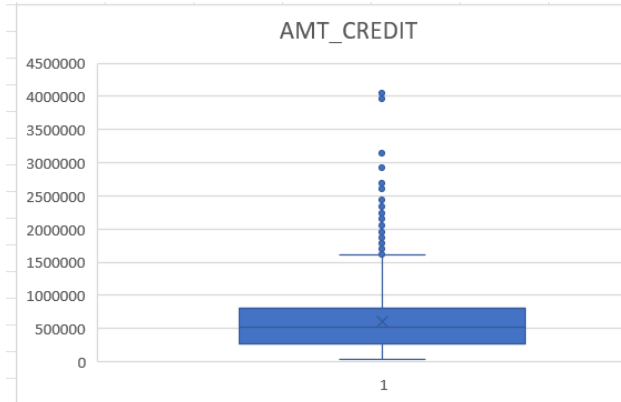


**Fig. Outlier Detection**

EXT_SOURCE_1

BASEMENTAREA_AVG

EXT_SOURCE_2

YEARS_BEGINXPOPULATION_AVG

EXT_SOURCE_3

YEARS_BUIL_AVG

APPARTMENT_AVG

LANDAREA_AVG
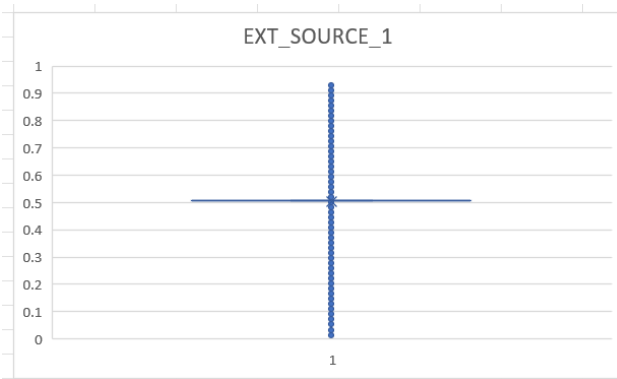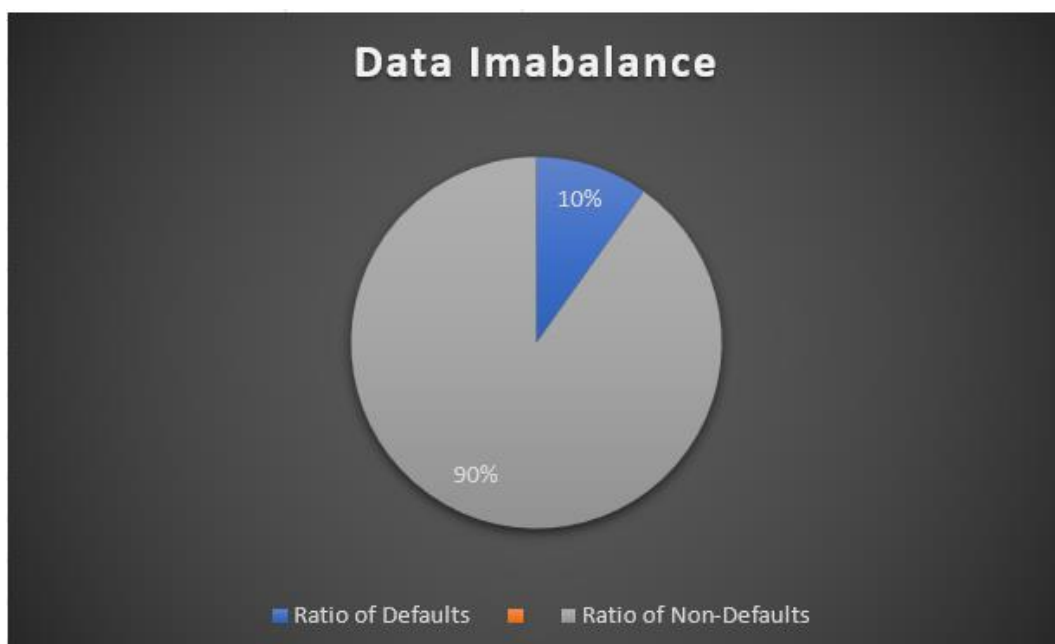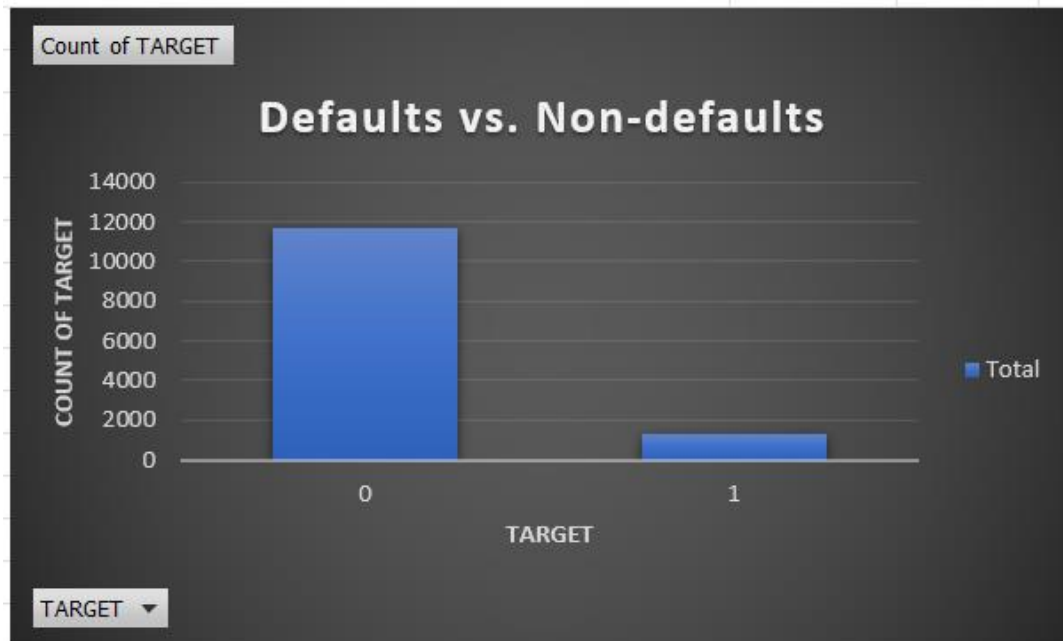
**C.** **Analyze Data Imbalance:** Determined data imbalance in the loan application dataset and calculated the ratio of data imbalance using Excel function COUNTIF and SUM to calculate the proportions of each class.

## Class Proportions

| Row Labels | Count of TARGET |
|---|---|
| 0 | 11698 |
| 1 | 1277 |
| Grand Total | 12975 |

Note: 1 represents customers with repayment difficulties (defaulted)
0 represents customers without repayment difficulties (non-defaulted)

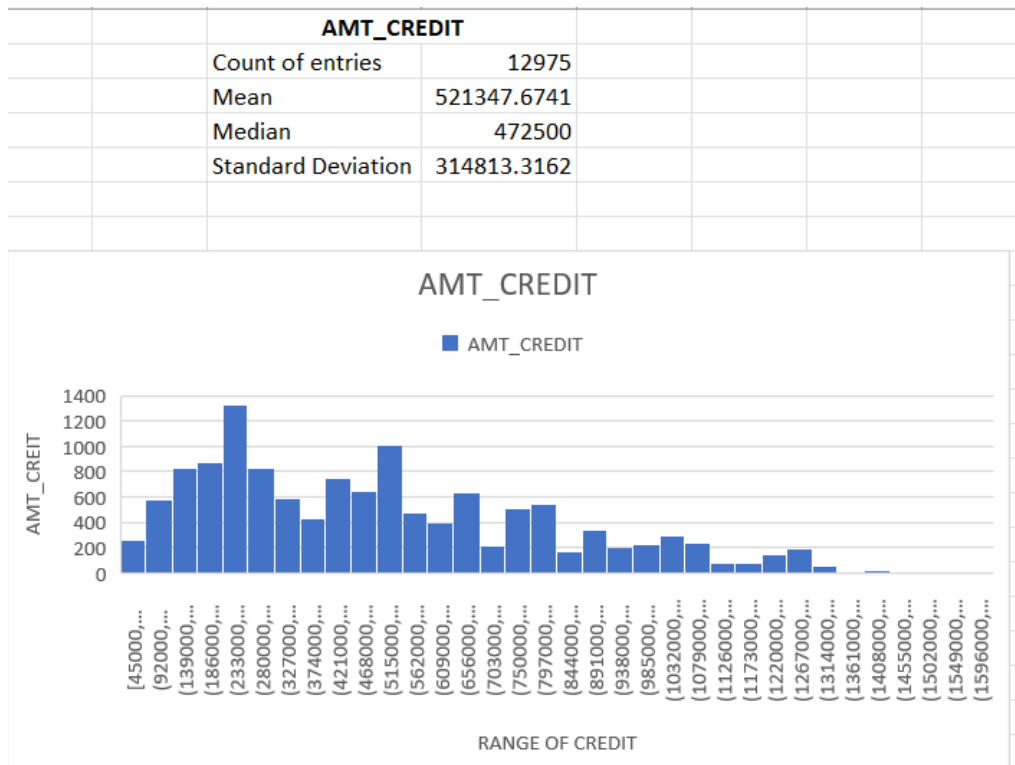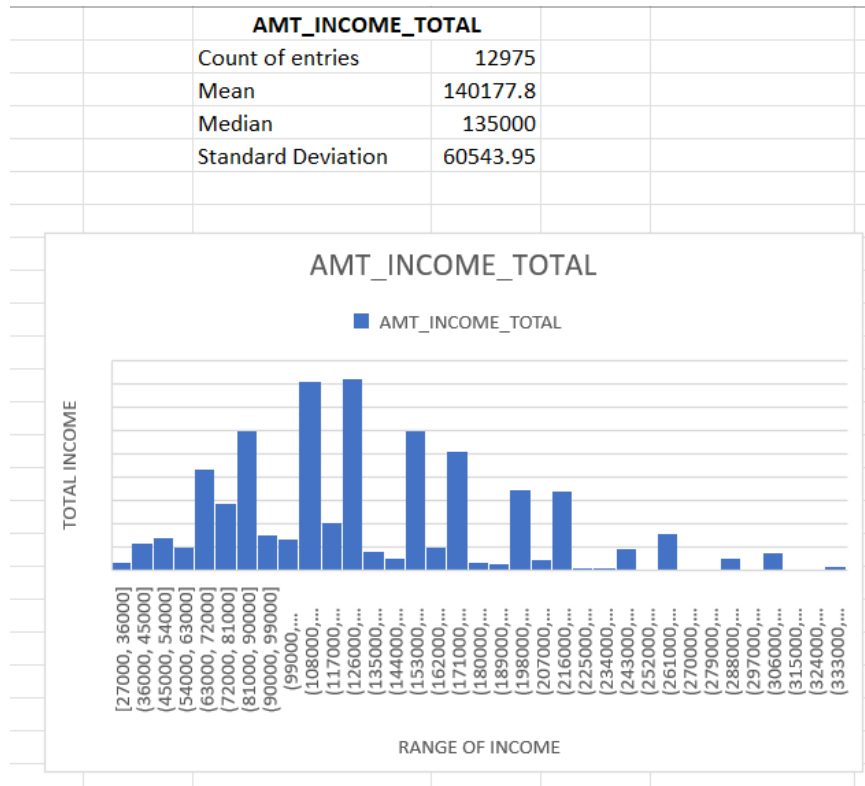| Target | Ratio | Interpret Data Imbalance |
|---|---|---|
| Ratio of Defaults | 10% | the ratio is significantly skewed (e.g., 90% non-defaults and 10% defaults), the dataset is imbalanced. |
| | | |
| Ratio of Non-Defaults | 90% | |
| | | |



Data Imabalance

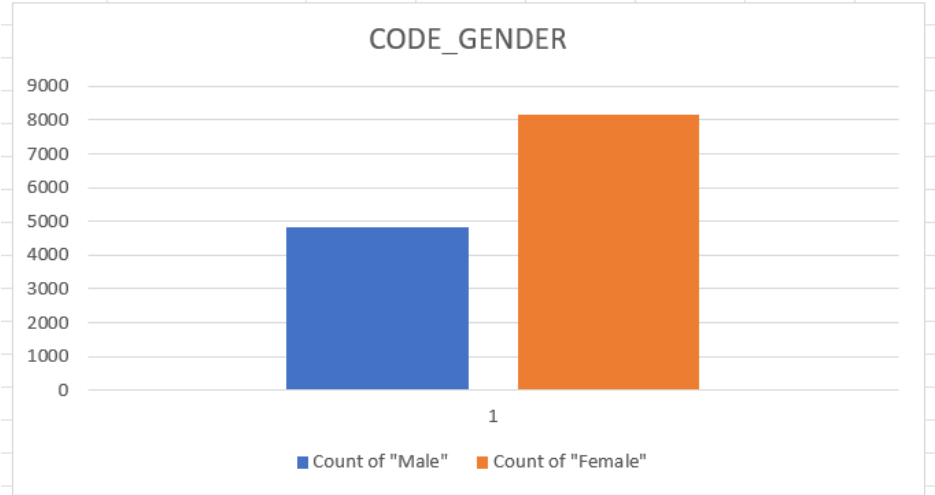**D.** **Perform Univariate, Segmented Univariate, and Bivariate Analysis:**

1. **Univariate Analysis:** Univariate analysis involves examining the distribution of a single variable at a time to understand its characteristics, such as central tendency, variability, and shape.

| AMT_INCOME_TOTAL | AMT_CREDIT | COAD_GENDER | NAME_CONTRACT_TYPE |
|---|---|---|---|
| 67500 | 135000 | M | Revolving loans |
| 135000 | 312682.5 | F | Cash loans |
| 121500 | 513000 | M | Cash loans |
| 99000 | 490495.5 | M | Cash loans |
| 135000 | 405000 | M | Revolving loans |
| 81000 | 270000 | F | Revolving loans |
| 90000 | 544491 | F | Cash loans |
| 112500 | 327024 | M | Cash loans |
| 202500 | 604152 | F | Cash loans |
| 202500 | 661702.5 | M | Cash loans |
| 90000 | 180000 | F | Revolving loans |
| 202500 | 305221.5 | F | Cash loans |
| 99000 | 260640 | F | Cash loans |
| 67500 | 298728 | F | Cash loans |
| 157500 | 755190 | M | Cash loans |
| 135000 | 675000 | F | Cash loans |
| 202500 | 1288350 | F | Cash loans |
| 112500 | 135000 | F | Revolving loans |
| 81000 | 252000 | F | Cash loans |
| 157500 | 760225.5 | M | Cash loans |
| 225000 | 270000 | M | Revolving loans |
| 72000 | 450000 | F | Cash loans |
| 126000 | 263686.5 | F | Cash loans |
| 135000 | 391194 | M | Cash loans |

**Fig. Categorial and numerical data for Univariate analysis**

| AMT_INCOME_TOTAL | |
|---|---|
| Count of entries | 12975 |
| Mean | 140177.8 |
| Median | 135000 |
| Standard Deviation | 60543.95 |



| AMT_CREDIT | |
|---|---|
| Count of entries | 12975 |
| Mean | 521347.6741 |
| Median | 472500 |
| Standard Deviation | 314813.3162 |

**COAD_GENDER**

| | |
|---|---|
| Count of "Male" | 4833 |
| Count of "Female" | 8141 |

## CODE_GENDER



**NAME_CONTRACT_TYPE**

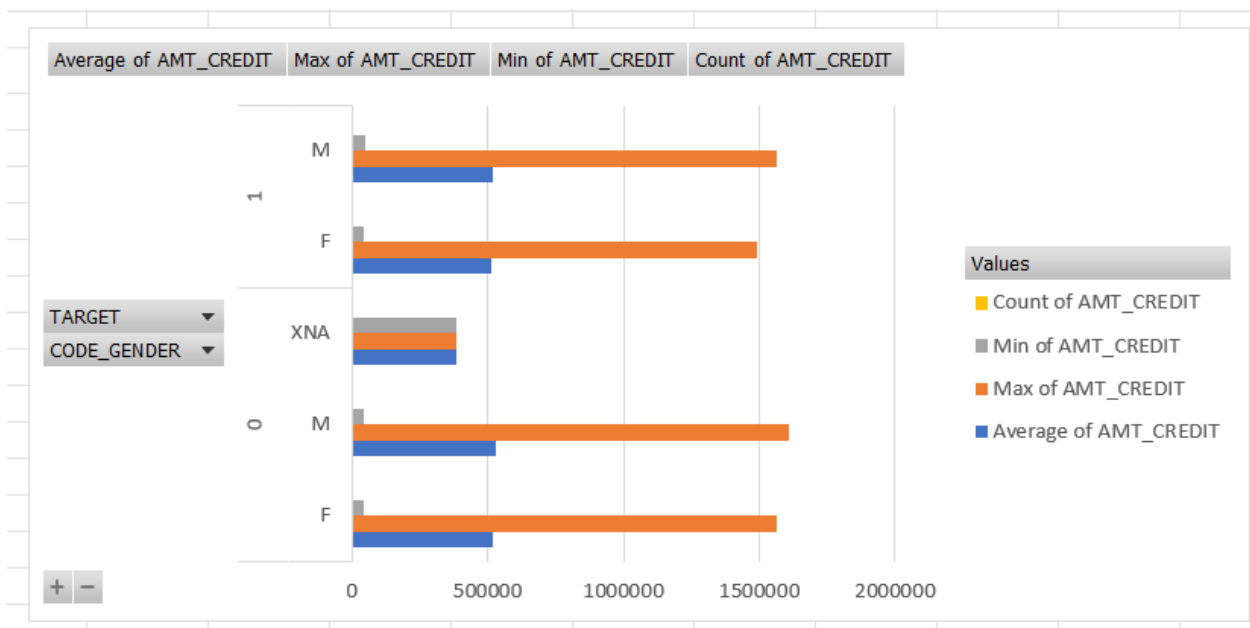| | |
|---|---|
| Count of "Cash loans" | 11936 |
| Count of "Revolving loans" | 1038 |

## NAME_CONTRACT_TYPE

2. **Segmented Univariate Analysis:** Segmented univariate analysis compares the distribution of a variable across different groups or scenarios, such as comparing AMT_INCOME_TOTAL for loan defaults vs. non-defaults.

| Comparison of Total Income with Gender and Default and Non Default Loan Payer | | | | |
|---|---|---|---|---|
| Row Labels | Average of AMT_INCOME_TOTAL | Max of AMT_INCOME_TOTAL | Min of AMT_INCOME_TOTAL | Count of AMT_INCOME_TOTAL |
| ⊟0 | | | | |
| F | 130259.7776 | 337500 | 27000 | 7480 |
| M | 157826.719 | 337500 | 27000 | 4217 |
| XNA | 207000 | 207000 | 207000 | 1 |
| ⊟1 | | | | |
| F | 127000.8722 | 337500 | 31500 | 661 |
| M | 153822.2143 | 337500 | 36000 | 616 |
| Grand Total | 140177.8398 | 337500 | 27000 | 12975 |



| Comparison of Loan Amount with Gender and Default and Non Default Loan Payer | | | | |
|---|---|---|---|---|
| Row Labels | Average of AMT_CREDIT | Max of AMT_CREDIT | Min of AMT_CREDIT | Count of AMT_CREDIT |
| ⊟0 | 522012.3406 | 1609272 | 45000 | 11698 |
| F | 517249.8475 | 1566909 | 45000 | 7480 |
| M | 530493.0046 | 1609272 | 45000 | 4217 |
| XNA | 382500 | 382500 | 382500 | 1 |
| ⊟1 | 515258.9753 | 1563291 | 45000 | 1277 |
| F | 512187.1611 | 1494486 | 45000 | 661 |
| M | 518555.1916 | 1563291 | 50940 | 616 |
| Grand Total | 521347.6741 | 1609272 | 45000 | 12975 |

3. **Bivariate Analysis:** Bivariate analysis examines the relationship between two variables. For this, you'll explore how various customer and loan attributes relate to the target variable (loan default).

| | Target(Default and Non-Default) | | | |
|---|---|---|---|---|
| | 0 | | 1 | |
| CODE_GENDER | Count of AMT_INCOME_TOTAL | Count of AMT_CREDIT | Count of AMT_INCOME_TOTAL | Count of AMT_CREDIT |
| F | 7480 | 7480 | 661 | 661 |
| M | 4217 | 4217 | 616 | 616 |
| XNA | 1 | 1 | | |

Count of AMT_INCOME_TOTAL    Count of AMT_CREDIT

TARGET

Values

- 1 - Count of AMT_CREDIT
- 1 - Count of AMT_INCOME_TOTAL
- 0 - Count of AMT_CREDIT
- 0 - Count of AMT_INCOME_TOTAL

CODE_GENDER

XNA

M

F

0    2000    4000    6000    8000



Correlation Analysis

**E.    Identify Top Correlations for Different Scenarios:** Identified top correlations for different scenarios in the Bank Loan Case Study, segment the dataset based on different groups (e.g., clients with payment difficulties and clients without payment difficulties) and analyze which variables have the strongest correlations with the target variable (TARGET, where 1 = default, 0 = no default)
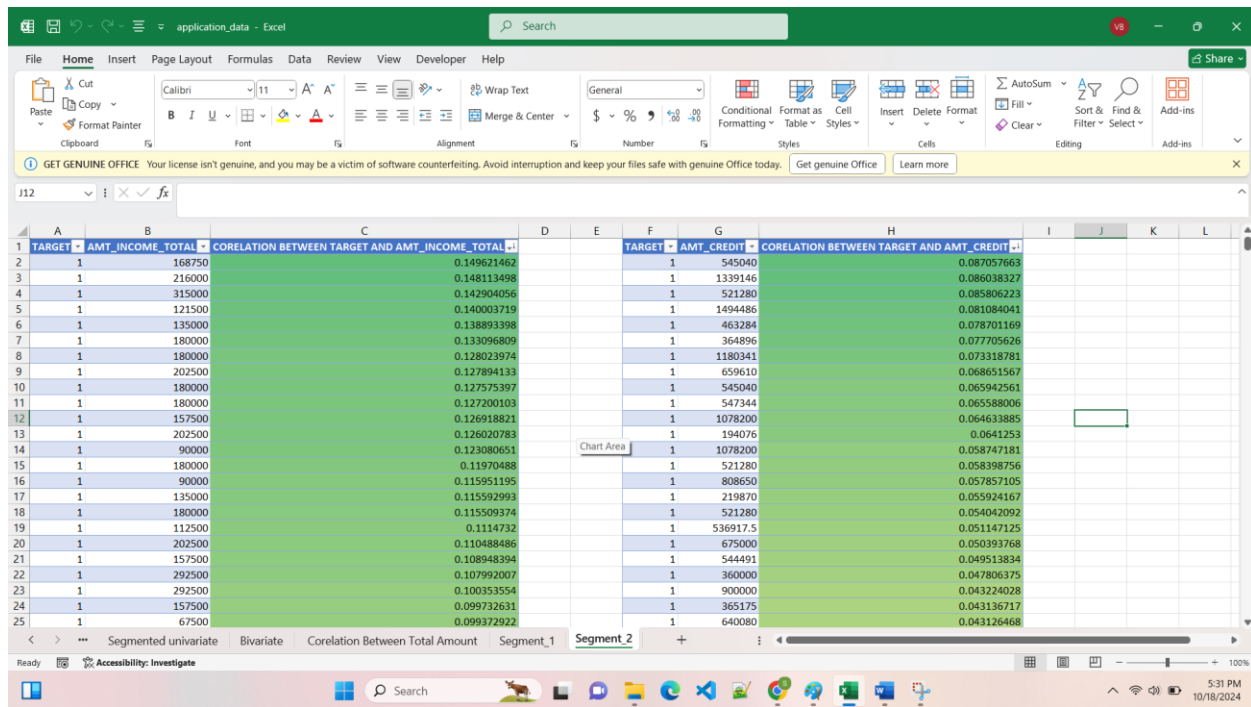
1.  **Segment_1:**

## 2. Segment_2:



## Result:

The findings from this analysis can greatly improve the company's loan approval process by helping identify risky applicants before approval. By focusing on key variables such as income, annuity, and loan amounts, the company can better balance financial risk while capturing new business opportunities. Going forward, advanced modeling techniques can be implemented to create more precise predictions, further enhancing the company's ability to manage its loan portfolio effectively.