

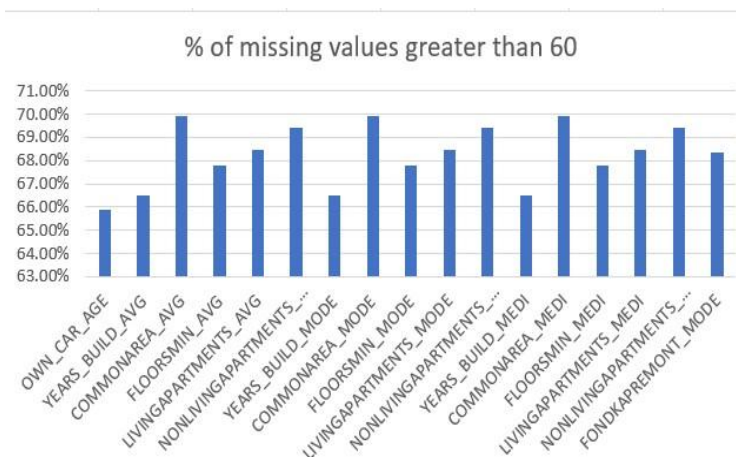
BANK LOAN CASE STUDY

Project Description: This project analyses loan application data to identify patterns that influence the likelihood of loan default. It also focuses on understanding the key factors and customer behaviour behind loan default so it can make better decisions about loan approval and rejection, reduce financial risks and ensure capable applicants are not rejected.

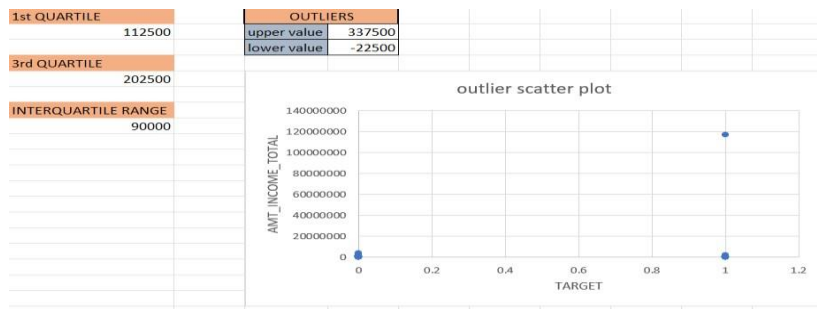
Approach: To perform the analysis of the Bank Loan Case Study datasets, I performed Exploratory Data Analysis focusing on handling missing data, filtering and filling missing values, deleting rows and deleting unnecessary columns especially the ones with more than 50% missing values. After that further Data Analysis tasks like identifying outliers, addressing data imbalance, performing univariate, segmented univariate, bivariate analysis and top correlations are carried out to uncover key factors driving loan default.

Data Analytics Tasks:

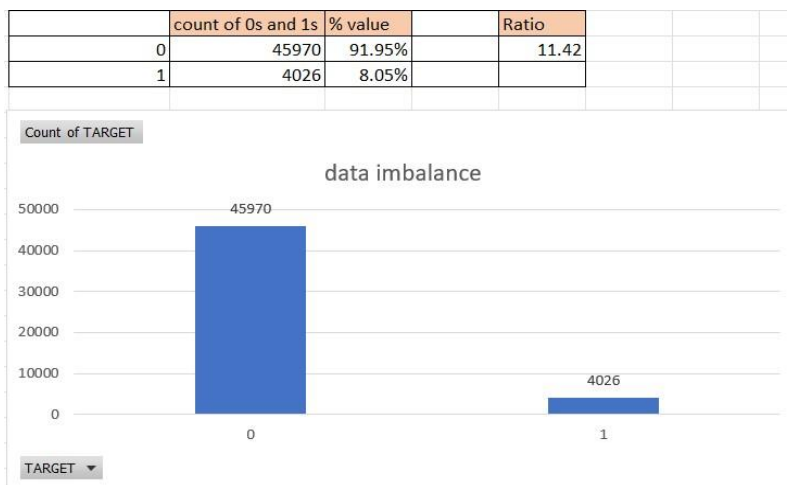
- A. Identify Missing Data:** The task involves identifying missing data in the dataset using countblank function and highlighting missing data for more than 60% using conditional formatting. For columns having missing values less than 40%, I imputed the numerical values using average formula and for categorical columns I simply put unknown.



- B. Identify Outliers in the Dataset:** The task involves detecting outliers in the dataset using excel functions like quartile and Interquartile Range on numerical columns.



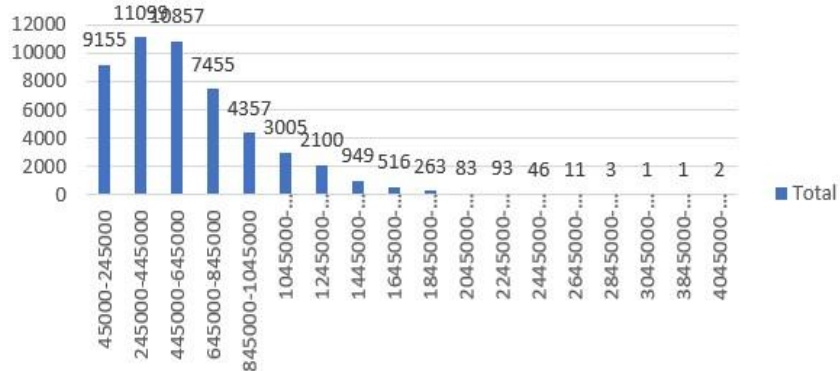
C. Analyze Data Imbalance: The task involves analysing data imbalance in the dataset and calculating the ratio of data imbalance using pivot tables and excel formulas.



D. Perform Univariate, Segmented Univariate and Bivariate Analysis: The task involves analysis on consumer and loan attributes. For univariate analysis, distribution of individual variables like income range, loan amount range and age range are conducted. For segmented univariate analysis, variables like income range, loan amount range and age are compared with the target variable. For bivariate analysis, the income range is compared with the average credit amount using pivot tables and pivot charts.

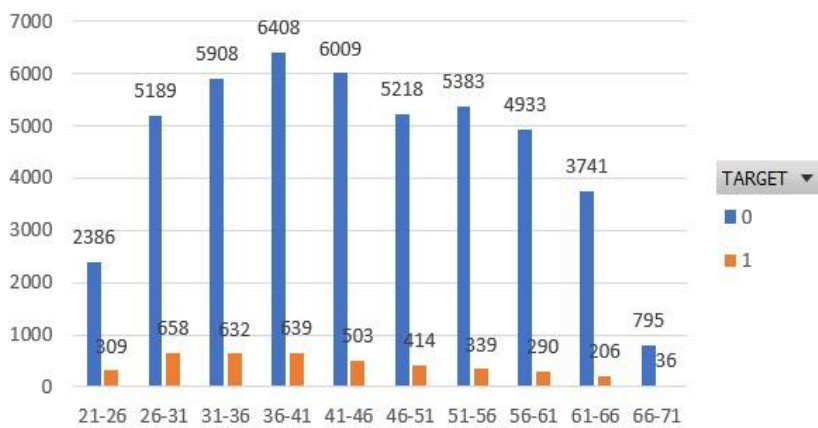
Count of SK_ID_CURR

count of applicants for loan amount range



AMT_CREDIT

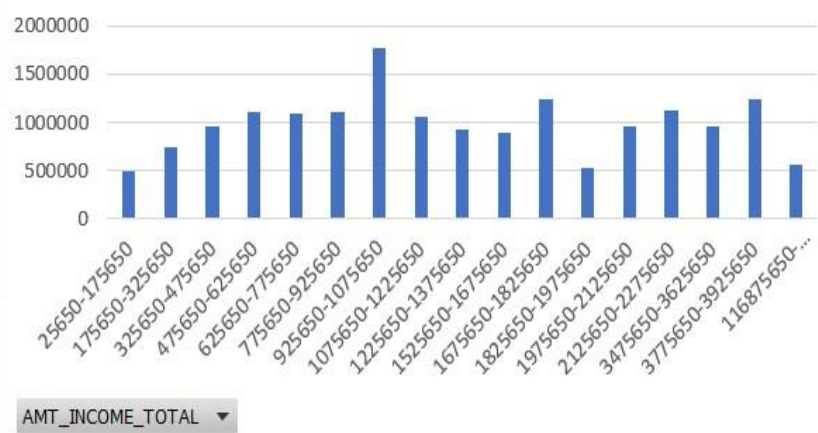
Count of TARGET



DAYS_BIRTH(YRS)

Average of AMT_CREDIT

average credit amount for income range



AMT_INCOME_TOTAL

- E. Identify Top Correlations for Different Scenarios:** The task involves identifying the top correlations for numerical variables in the dataset using the Excel correl function and visualizing the correlations using conditional formatting.

TOP CORRELATIONS	
CNT_CHILDREN	0.026354
AMT_INCOME_TOTAL	0.010897
AMT_CREDIT	-0.03243
REGION_POPULATION_RELATIVE	-0.0408
AGE	-0.07681

Tech-Stack Used: For this project, the software used is Microsoft Excel 2021 as it is a very popular tool for both beginners and advanced professionals. It is widely used for data analysis because of its accessibility, powerful functionality, versatility and ease of use.

Insights:

- Certain less significant variables having missing values were also removed to ensure data quality.
- Outliers were detected in variables like amount of income and number of children.
- The dataset was imbalanced with a higher proportion of non-default cases compared to default cases.
- As the age of applicants increases, the chances of defaulting decreases.
- Lowest income range has highest number of defaults.

Result: This project provided comprehensive insights into the Bank Loan Case Study dataset and helped me understand how key factors like customer attributes and loan attributes help us understand loan defaults and how we can make better decisions about loan approvals and loan rejections so that capable applicants are not rejected as to not lose business and incapable applicants do not get approved mistakenly to avoid financial loss.