# Analysis of fraud detection in bank data using EDA

# Problem statement

- A certain bank wants to find a way to provide loan for the applicants to run its business .

- But they facing difficulty in finding the client as :
  1. providing loan to a client who is not likely to repay loan will lead to financial loss.
  2. Denying loan to a client who is likely to repay loan will lead to business loss.

- To find a solution, we are given a three datasets which contains the current application and previous application data of the clients and description of the columns.

- We are required to perform the EDA and use the data visualization techniques on these data to find the clients who are most likely to repay the loan and who are not.

# Assumptions

- Since some the data are missing in the data , we are removing columns which has higher percentage missing data as imputing values to them may affect the analysis.

- For the data where the missing values are less, we are imputing them with median value of the remaining data for continuous data and mode value of the remaining data for categorical data.

# Approach and methodology

**UNDERSTANDING THE DATA**

- The primary step involves understanding the data, here we are given a dataset containing column description to understand the meaning of the data and their importance

**IMPORTING THE DATA**

- Now we load the data into the platform to check its shape, datatypes and other information.

**HANDING MISSING DATA**

- Generally the data(column) containing more than 30% (depending on total number of columns) of the missing values are removed.

- For the columns containing lesser missing values, are replaced median value for continuous data and mode value for categorical data.
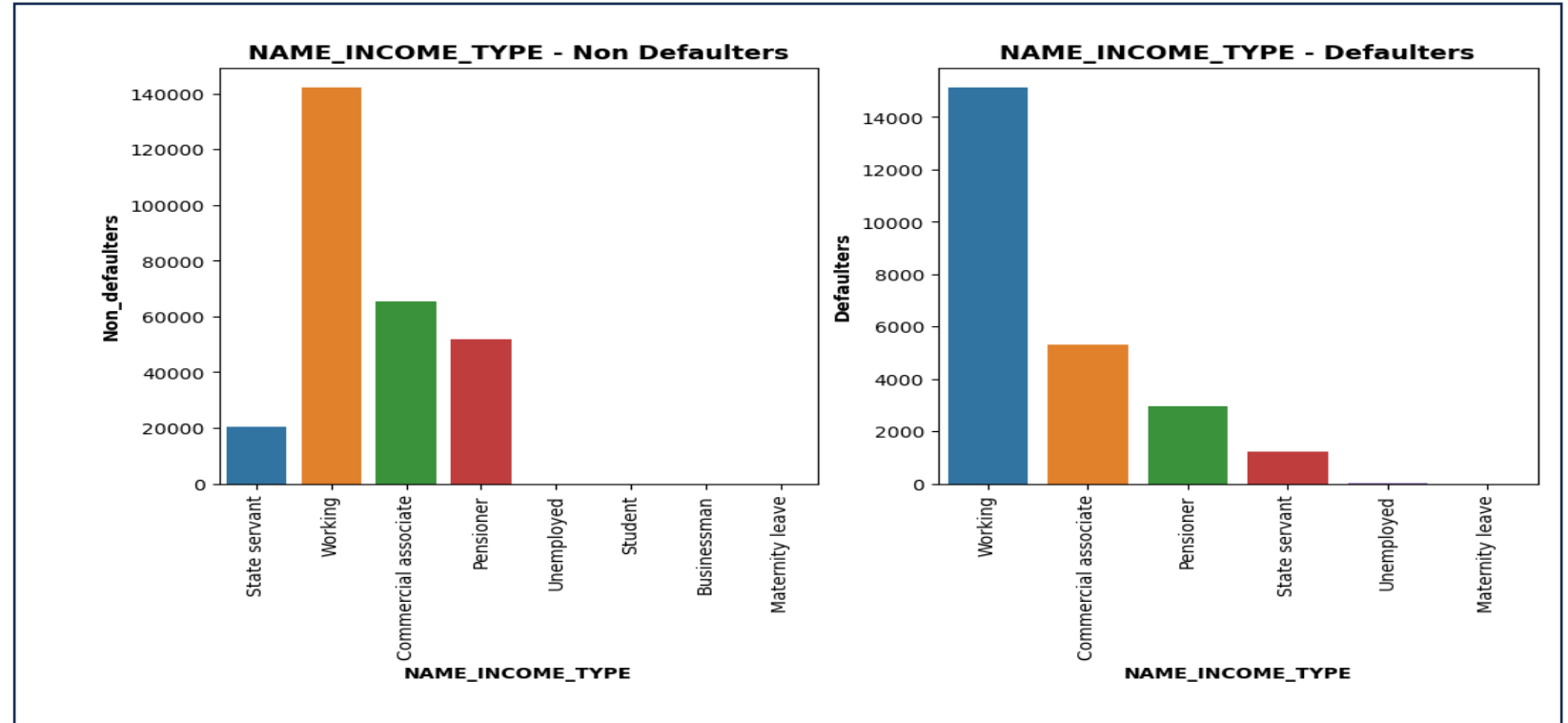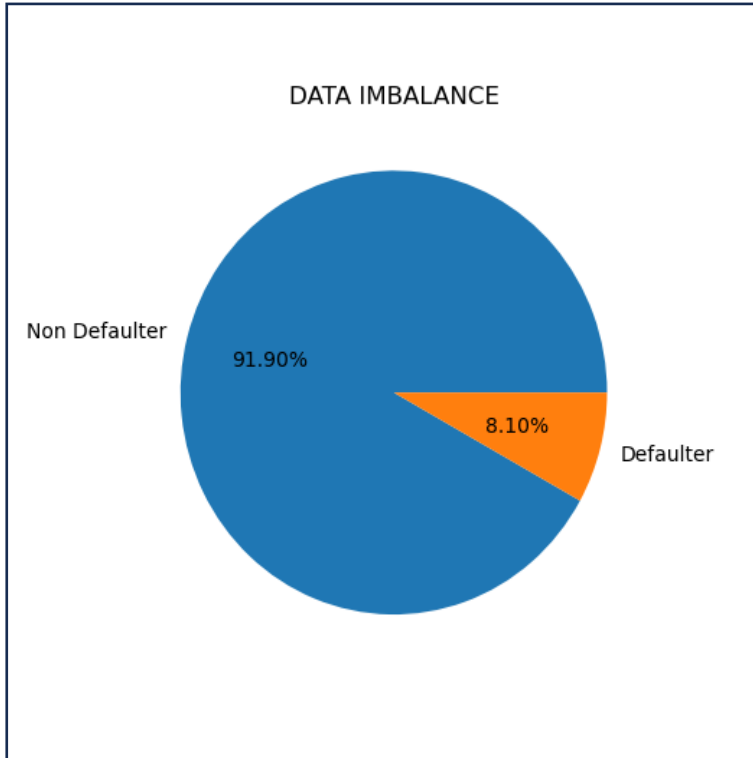
# Approach and methodology

**OUTLIERS**

- After handling the missing data, we use describe method or boxplot to visually identify the potential outliers which are later verified further for data correctness.
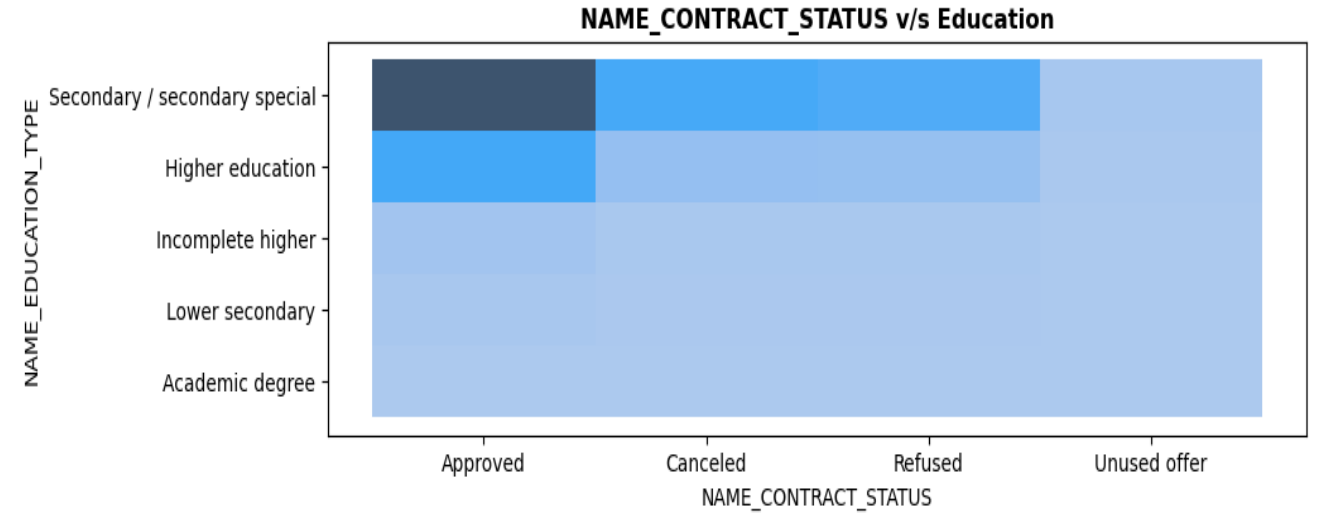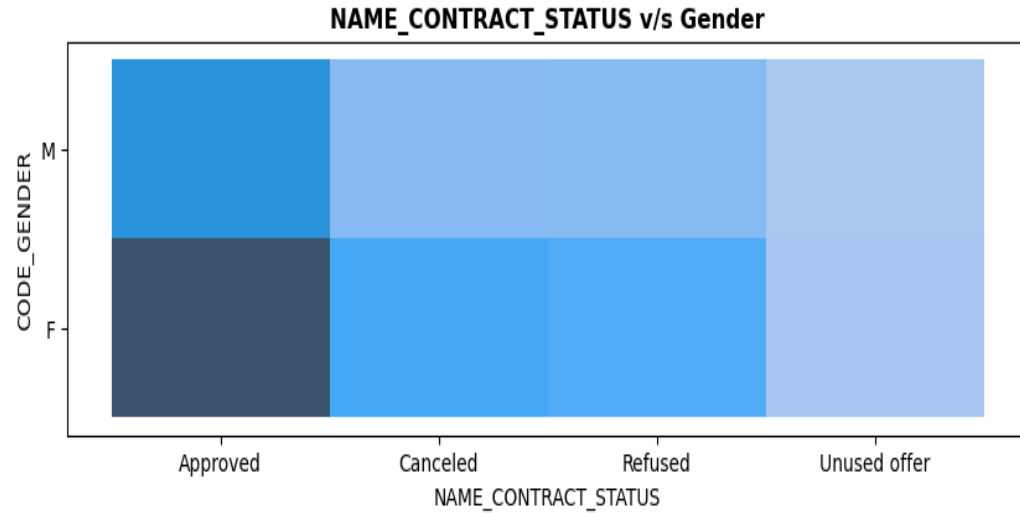
**VISUAL ANALYSIS**

- Here we use various sophisticated visual techniques such as UNIVARIATE, BIVARIATE AND MULTIVARIATE ANALYSIS for identifying the important feature, trends and correlation among the data

- The above information is used to find the valuable insights and applied in favour of the target variable

- Such examples can be seen in next slide.
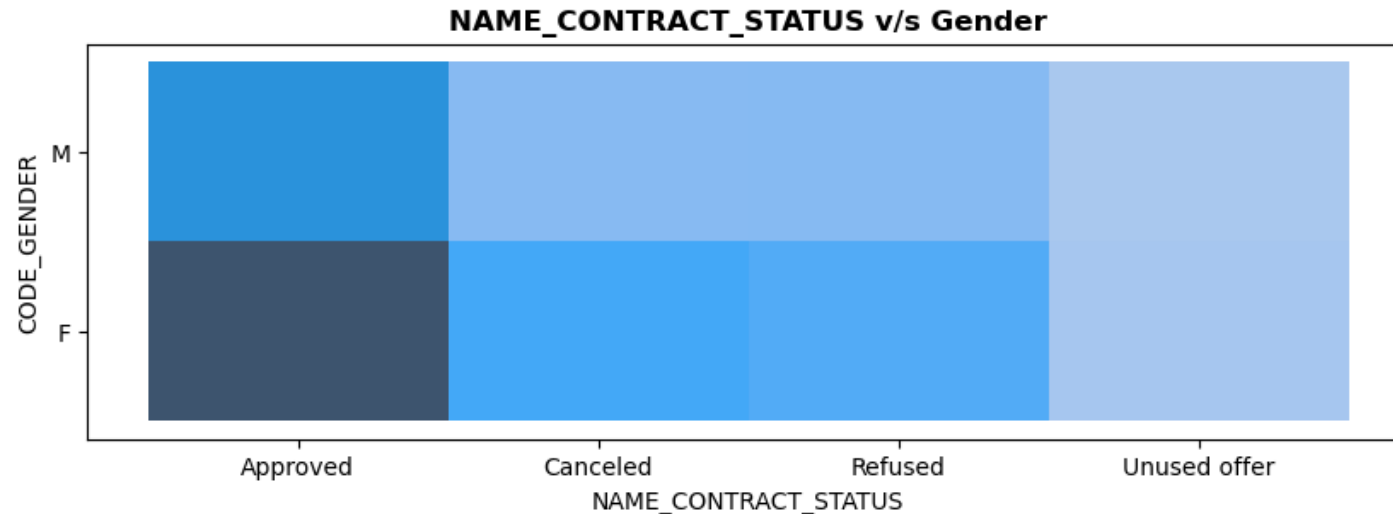
# Graphs and Insights - Univariate



- Here we can see that the data imbalance ratio is 11.35, which is a favorable condition for the bank as it indicated that the defaulter are less compared to non defaulters.
- In the second graph we can see that the working class people are applying more for the loan.

# Graphs and Insights - Bivariate

### NAME_CONTRACT_STATUS v/s Gender



### NAME_CONTRACT_STATUS v/s Education



Here we can see that the approval rate is high repeater applicant and in applicant with secondary education

### NAME_CONTRACT_STATUS v/s Gender



- Approval rate is higher in females than males

# Graphs and Insights - Multivariate

- AMT_GOODS_PRICE and AMT_CREDIT have the highest correlation

- AMT CREDIT and AMT_ANNUITY also has high correlation

# Recommendations and Insights

- We have observed that the applicant who have secondary education are most likely to repay the loan, hence the bank can lend loan to them.

- Banks can offer the loans to repeater as they already have a good credit history and likely to repay the loan

- Students and businessmen have almost negligible defaulters, hence bank can offer the loan for them as well

- Banks can avoid clients with academic degree and medium level income as the defaulters are more than non defaulters in these class

# THANK YOU