# Bank Loan Case Study

**Data Understanding:**

**Download the Dataset using the link given under dataset section on the right.**

1. `application_data.csv` contains all the information of the client at the time of application.
   The data is about wheather a client has payment difficulties.
2. `previous_application.csv` contains information about the client's previous loan data. It contains the data whether the previous application had been Approved, Cancelled, Refused or Unused offer.
3. `columns_descrption.csv` is data dictionary which describes the meaning of the variables.

**You are required to provide a detailed report for the below data record mentioning the answer to the questions that follows:**

- Present the overall approach of the **analysis**. Mention the problem statement and the analysis approach briefly
- **Identify** the missing data and use appropriate method to deal with it. (Remove columns/or replace it with an appropriate value)
  *Hint: Note that in EDA, since it is not necessary to replace the missing value, but if you have to replace the missing value, what should be the approach. Clearly mention the approach.*
- Identify if there are **outliers** in the dataset. Also, mention why do you think it is an outlier. Again, remember that for this exercise, it is not necessary to remove any data points.
- Identify if there is data imbalance in the data. Find the ratio of data imbalance.
  *Hint: Since there are a lot of columns, you can run your analysis in loops for the appropriate columns and find the insights.*
- Explain the **results of univariate, segmented univariate, bivariate analysis, etc.** in business terms.
- Find the top 10 **correlation** for the Client with payment difficulties and all other cases (Target variable). Note that you have to find the top correlation by segmenting the data frame w.r.t to the target variable and then find the top correlation for each of the segmented data and find if any insight is there. Say, there are 5+1(target) variables in a dataset: Var1, Var2, Var3, Var4, Var5, Target. And if you have to find top 3 correlation, it can be: Var1 & Var2, Var2 & Var3, Var1 & Var3. Target variable will not feature in this correlation

as it is a categorical variable and not a continuous variable which is increasing or decreasing.

- **Include visualizations** and **summarize** the most important results in the presentation. You are free to choose the graphs which explain the numerical/categorical variables. Insights should explain why the variable is important for differentiating the clients with payment difficulties with all other cases.

**Project Description:**

The project is about carrying out Exploratory Data Analysis (EDA) for the bank to understand the defaulters and people who will be able to repay the loans.

**Exploratory Data Analysis (EDA):**

It is an approach to analyse the data using visual techniques. It is used to discover trends, patterns or to check assumptions with the help of statistical summary and graphical representations.

**Approach:**

- Analyse the data
- Clean the data
- Filling the missing values
- Finding the outliers
- Pivot tables
- Charts to visualise

**REMOVING BLANK CELLS**

Select the cells the particular column you see blank cells as V1:V307512

Click **F5** and go to **special** there select **blanks** and click **ok**



All cells in that particular column are selected



Now Press **F2** and write the value you wish to enter which we have taken **0** and then click **ctrl+enter** all blank cells get filled by 0 repeat the step in all data sets.

**Approach to replacing the missing values:**

First we select the column we would want to fill the missing values with the mean of the data.

Then we go to the Data toolbar and select Data Analysis from the Analyse tool. The following popup will show up.

Next step is to select Descriptive statistics to get the complete data of the column

Click Ok

Then screen opens which asks you to select the column and other attributes



Select summary Statistics and click ok

**Descriptive Statistics** is used to **summarize** a given **data** found from any study. It can provide **basic information** and the **internal relationship** between the **variables** in a dataset.

We have to convert the general formatted data to number format for the descriptive statistics too to work as it does not work on non numeric data.

For that we have to go to the home tab and change the column property from general to number format, then proceed with the analysis

| APARTMENTS_AVG | |
| --- | --- |
| Mean | 0.05784 |
| Standard Error | 0.000173 |
| Median | 0 |
| Mode | 0 |
| Standard Deviation | 0.096007 |
| Sample Variance | 0.009217 |
| Kurtosis | 14.67023 |
| Skewness | 3.03722 |
| Range | 1 |
| Minimum | 0 |
| Maximum | 1 |
| Sum | 17786.36 |
| Count | 307511 |

Now select the colums with null values and fill it by mean, median or mode as per the requirement. Continue this process with all the datasets.

**OUTLIERS in EXCEL**

1. Review your data
2. Sort your data values
3. Analyse your values
4. Identify your data quartiles
5. Define the interquartile range
6. Calculate the upper and lower bound
7. Remove your outliers

**What is an Outlier?**

An outlier is a data point within a data set that lies outside of the range of most of the other data points. For example in the date set of ages 12, 13, 15, 16, 52, 14, and 11, you can see that 52 is the outlier age. This is because the other ages fall within the range of 11 to 16 and 52 is outside of that range.

It's important for statisticians to be able to identify outliers like this since they can dramatically alter their calculations. In this example, if you include the outlier the mean average age is 19. If you exclude the outlier, the mean is 13.5.

You need to investigate any outliers carefully before removing them. It's possible the outliers are simple mistakes that should be excluded. Alternatively, those outliers may

contain important statistical information such as a new trend or a significant insight that needs to be considered.

**How to find outliers?**
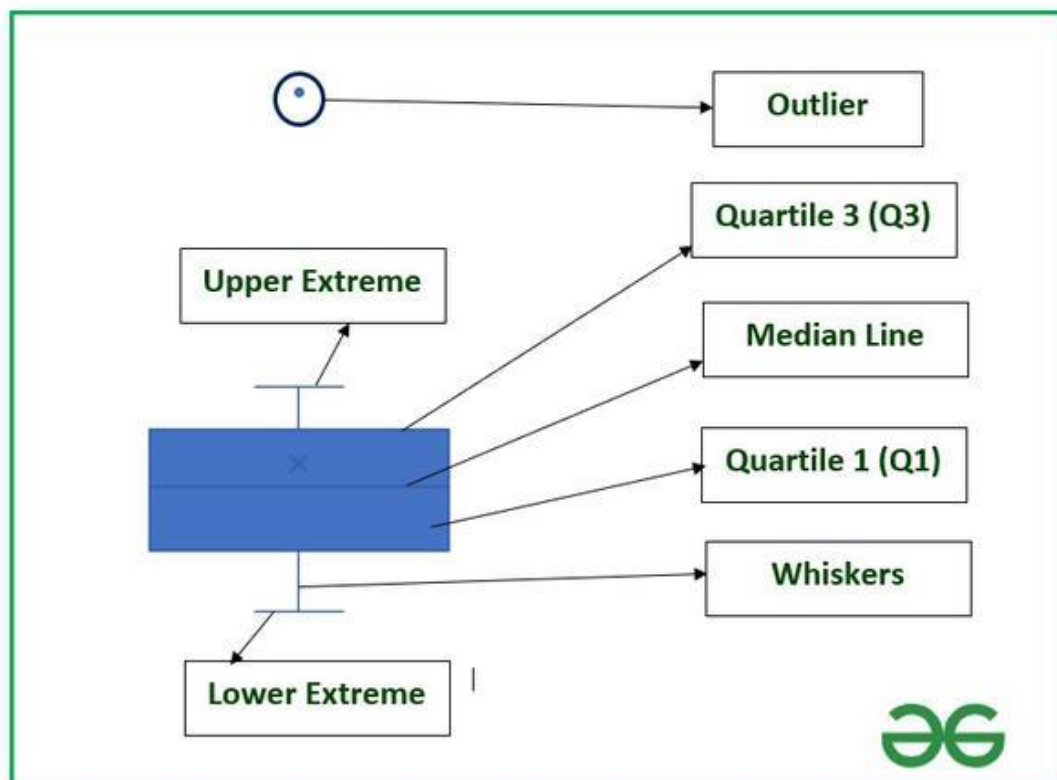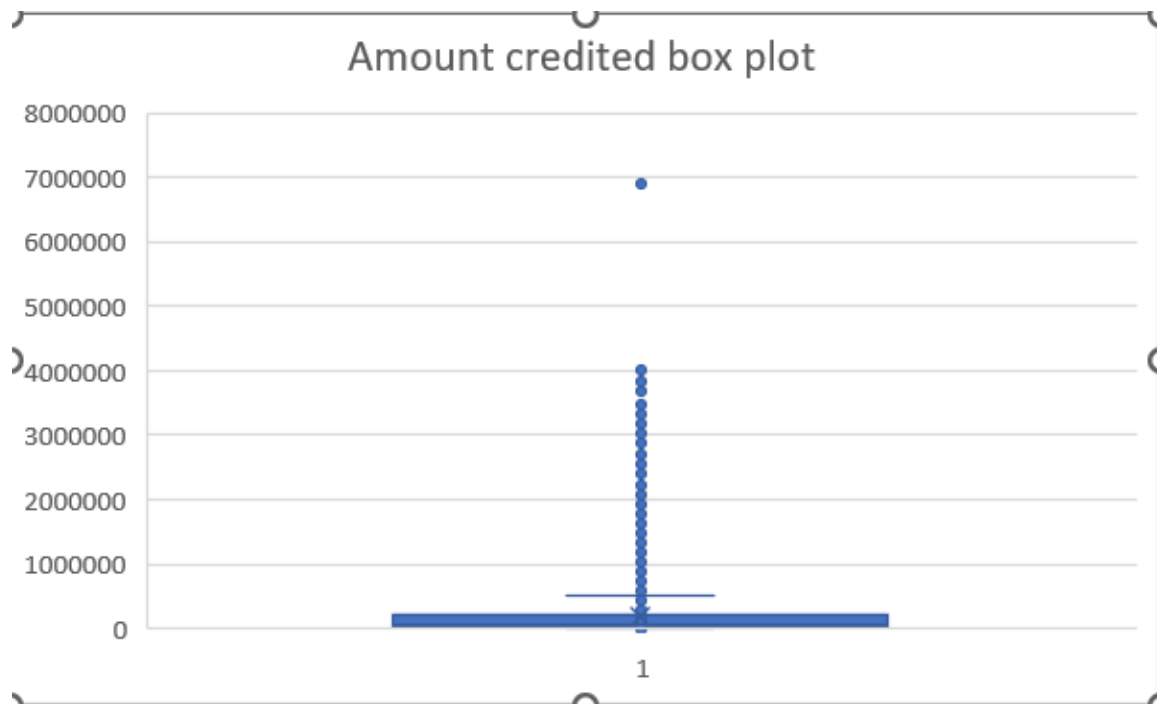
**Statistical concepts to know outliers**

1**. Box Plots** – in the image below you can see that several points exist outside of the box. The box is the central tendency of the data.  It is clustered around a middle value. The upper bound line is the limit of the centralization of that data.
2**. Quartiles**– represent how the data is broken up into quarters.

Plotting the box plot for the amount credited and identifying the outliers

A box plot shows 5 values:

- Minimum
- First quartile
- Median
- Third quartile
- Maximum

The above box plot is showing the outliers in the amount credit column similarly we perform the same for all the datasets

**What is data imbalance**

Imbalanced data refers to those types of datasets where the target class has an uneven distribution of observations, i.e one class label has a very high number of observations and the other has a very low number of observations. We can better understand imbalanced dataset handling with an example.

Let's assume that XYZ is a bank that issues a credit card to its customers. Now the bank is concerned that some fraudulent transactions are going on and when the bank checks their data they found that for each 2000 transaction there are only 30 Nos of fraud recorded. So, the number of fraud per 100 transactions is less than 2%, or we can say more than 98% transaction is "No Fraud" in nature. Here, the class "No Fraud" is called the majority class, and the much smaller in size "Fraud" class is called the minority class.

**Finding data Imbalance:**

**Dataset:** previous_application.csv

The column is distributed in binomial values which are 0 and 1

Considering the count of 1s and 0s we use countif() function

Here is the following result:

| value | count |
|-------|--------|
| 1 | 208157 |
| 0 | 840418 |

Now we plot a bar chart to see imbalance in the dataset:

**Dataset:** application.csv

Here we will consider AMT_REQ_CREDIT_BUREAU_YEAR column

The column is distributed in binomial values which are 0 and 1

Considering the count of 1s and 0s we use countif() function

Here is the following result:

| value | Count |
|-------|-------|
| 1     | 63405 |
| 0     | 71801 |

Now we plot a bar chart to see imbalance in the dataset:



**Results of univariate, segmented univariate, bivariate analysis**

**Univariate Analysis:**

Uni means one and variate means variable, so in univariate analysis, there is only one dependable variable. The objective of univariate analysis is to derive the data, define and summarize it, and analyze the pattern present in it. In a dataset, it explores each variable separately. It is possible for two kinds of variables- Categorical and Numerical.

Some patterns that can be easily identified with univariate analysis are Central Tendency (mean, mode and median), Dispersion (range, variance), Quartiles (interquartile range), and Standard deviation.

For finding out univariate analysis of application.csv we will use **descriptive statistics** tool from **Data Analysis** in **Data tab**

| AMT_INCOME_TOTAL | |
|---|---|
| Mean | 168797.9193 |
| Standard Error | 427.6058332 |
| Median | 147150 |
| Mode | 135000 |
| Standard Deviation | 237123.1463 |
| Sample Variance | 56227386501 |
| Kurtosis | 191786.5544 |
| Skewness | 391.5596541 |
| Range | 116974350 |
| Minimum | 25650 |
| Maximum | 117000000 |
| Sum | 51907216961 |
| Count | 307511 |

**Bivariate Analysis:**

**Bivariate analysis** is one of the statistical analysis where two variables are observed. One variable here is dependent while the other is independent. These variables are usually denoted by X and Y. So, here we analyse the changes occured between the two variables and to what extent. Apart from bivariate, there are other two statistical analyses, which are Univariate (for one variable) and Multivariate (for multiple variables).

In statistics, we usually interpret the given set of data and make statements and predictions about it. During the research, an analysis attempts to determine the impact and cause in order to conclude the given variables.

**Scatter Plot between AMT_INCOME_TOTAL and SK_ID_CURR**

**Find the top 10 correlation for the Client with payment difficulties and all other cases (Target variable).**

To find 10 top correlations we will include 10 columns from both previous_application.csv dataset and application.csv dataset



NAME_EDUCATION_TYPE

| Qualification | count |
|---|---|
| Secondary / secondary special | 218391 |
| Higher education | 74863 |
| Incomplete higher | 10277 |
| Lower secondary | 3816 |
| Academic Degree | 164 |

NAME_FAMILY_STATUS



NAME_HOUSING_TYPE

NAME_TYPE_SUITE



NAME_INCOME_TYPE

| income type | Count |
|---|---|
| Businessman | 10 |
| Commercial associate | 71617 |
| Maternity Leave | 5 |
| Pensioner | 55362 |
| State Servant | 21703 |
| Unemployeed | 22 |
| Working | 158774 |

Similarly we continue the same with all correlations and establish relations of variables.

Visualization using different graphs is also shown here.