

SUMIT SHAMLAL CHAURE

# Bank Loan Case Study

## Trainity Project 6 (Final Project -2)

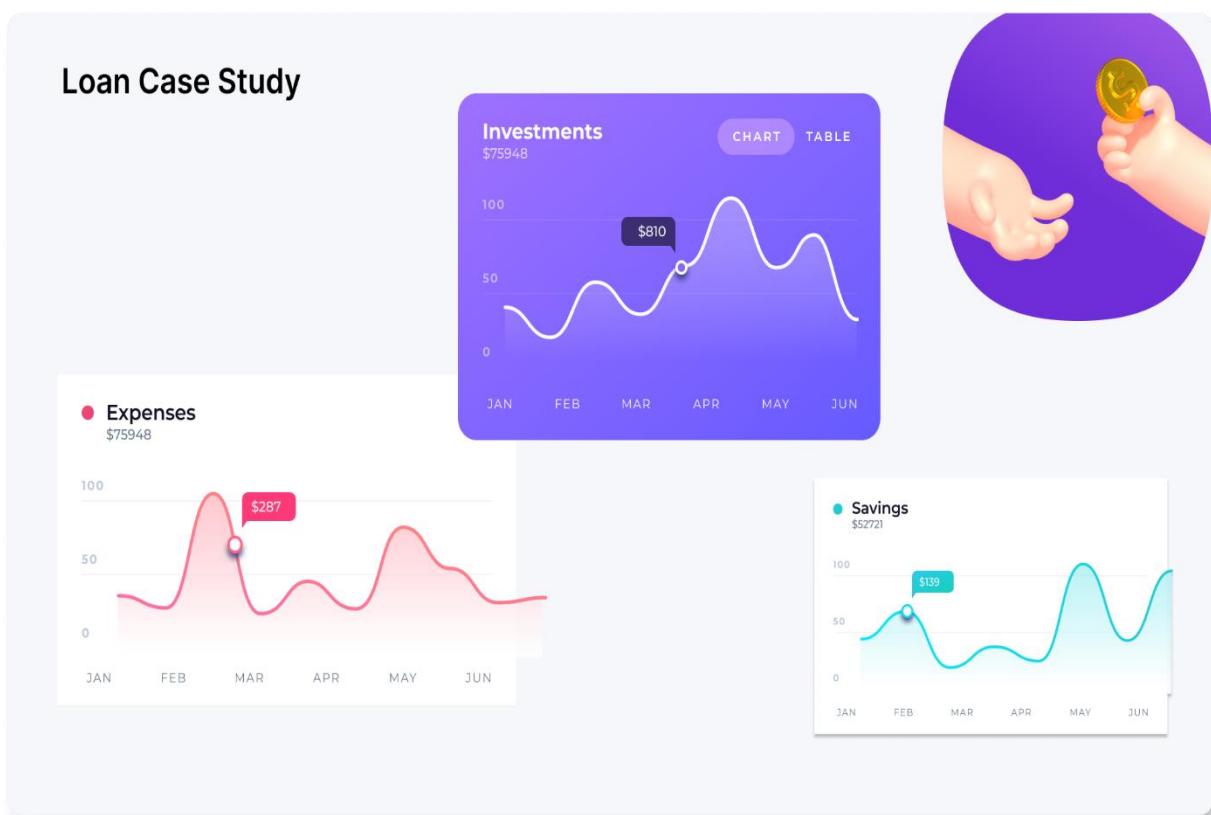


Figure 1 - <https://trainity.link/data/project06>

**Note** – Docx files uploaded on drive gets converted to sheets so much of the connections and pivot relations are lost so if possible download the zip file or main file so that all the things are intact and graphs or pivots show data accordingly while I have attached individual questions the sheets might not show up everything so download the excel files not Gdocs file to see the analysis.

Important Links ([Tap Here](#)) ([Drive Folder](#))

---

## INDEX

### Contents

<a href="#"><u>Introduction :</u></a>	1
<a href="#"><u>DESCRIPTION</u></a>	3
<a href="#"><u>Problem Statement:</u></a>	3
<a href="#"><u>Business Objectives:</u></a>	3
<a href="#"><u>Requirements –</u></a>	4
<a href="#"><u>Data Analytics Tasks :</u></a>	8
A) Identify Missing Data and Deal with it Appropriately: <i>As a data analyst, you come across missing data in the loan application dataset. It is essential to handle missing data effectively to ensure the accuracy of the analysis.</i>	8
B) Identify Outliers in the Dataset: <i>Outliers can significantly impact the analysis and distort the results. You need to identify outliers in the loan application dataset.</i>	11
C) : Analyze Data Imbalance: <i>Data imbalance can affect the accuracy of the analysis, especially for binary classification problems. Understanding the data distribution is crucial for building reliable models.</i>	15
D) Perform Univariate, Segmented Univariate, and Bivariate Analysis: <i>To gain insights into the driving factors of loan default, it is important to conduct various analyses on consumer and loan attributes.</i>	18
E): Identify Top Correlations for Different Scenarios: <i>Understanding the correlation between variables and the target variable can provide insights into strong indicators of loan default.....</i>	26
<a href="#"><u>Important Links &amp; Greetings :</u></a>	29

---

## **DESCRIPTION**

### **Problem Statement:**

Imagine you're a data analyst at a finance company that specializes in lending various types of loans to urban customers. Your company faces a challenge: some customers who don't have a sufficient credit history take advantage of this and default on their loans. Your task is to use Exploratory Data Analysis (EDA) to analyze patterns in the data and ensure that capable applicants are not rejected.

### **When a customer applies for a loan, your company faces two risks:**

1. If the applicant can repay the loan but is not approved, the company loses business.
2. If the applicant cannot repay the loan and is approved, the company faces a financial loss.

The dataset I will be working with contains information about loan applications ie. Current application or **application\_data.csv file**. It includes two types of scenarios:

1. Customers with payment difficulties: These are customers who had a late payment of more than X days on at least one of the first Y installments of the loan.
2. All other cases: These are cases where the payment was made on time.

### **Business Objectives:**

The main aim of this project is to identify patterns that indicate if a customer will have difficulty paying their installments. This information can be used to make decisions such as denying the loan, reducing the amount of loan, or lending at a higher interest rate to risky applicants. The company wants to understand the key factors behind loan default so it can make better decisions about loan approval.

---

## **Requirements –**

### **1) Project Description:**

The main aim of the project is to find the use the knowledge of Descriptive statistics & make use of visualization tools to predict the loan payments and the factors that may otherwise hinder with the timely payment of debts. In order to do that we need find out necessary insights from excel sheet to grab the details and trends of the applicants like his income, lifestyle and factors that may affect his budgets and spending's like family members and children's count for which the client need to make spends in the long run also factors like the place where the applicant is planning to buy the house or utilizing the credit amount so as to make assumptions of his finances.

After looking at the data we plan a format to operate on the data, tools needed and charts that can be useful for trends etc. while keeping in the questions in mind. The descriptive analysis at various stages of the questions helps us to get the major trends and their effects on loan repayment analysis as well gain deep insights for deciding the factors for future loan approvals by analyzing trends according to income, profession, education , family members and region counts and taking into account the defaulting factors and other correlated factors that we gain from the overall dataset analysis.

### **2) Approach:**

I first analyzed the data and looked for null values, blanks, duplicates and treated them using basic functions like delete cells, find & replace, remove blank rows , imputing with appropriate statistics like mean, mode, median values etc. For certain cell values I changed them to suit better with other values like – (age in years instead of days) etc. After confirming that the data has little to less outliers (single high income or credit limit etc.) and saved the raw data to work on with the operations. Also, I removed the most irrelevant columns from the dataset as the presence of such redundant data will not help us in the further analysis like few categorical columns had large amounts of blanks or numerical data with 50% and above missing values.

The analysis based on excel functions, Descriptive analysis, pivot tables has insights at the bottom of the screengrabs to let the others understand the aim of each analysis.

### 3) Tech-Stack Used:

**Excel** – The basic data manipulation, handling and overall pivot charts and the statistics has been handled using MS Excel also the descriptive statistics and correlation matrix have been done on same.

**Google Sheets** – Used to do the basic data manipulation and to get column stats (gif added) and also to verify some statistics before treating them.

**Word** – The report is written in word/docx format using MS Word and then exported to pdf.

**Drive** – To upload all the essential files attached in the report for reference & report upload.

**Loom** – To make the video presentation of the analysis of each individual questions (links attached for the same below)

### 4) Insights:

I have done the analysis on **application\_data.csv** file which pertains to the current application details of the clients.

The summary for each query is given with the screenshot but to summarize the overall dataset I came to the conclusion that among the total data set we had **49999 data rows, 122 columns & a total of 4735935 cells** in the given dataset.

A1:D1	SK_ID_CURR	TARGET	NAME_CONTR_CODE	GENDER	FLAG_OWN_CASH	FLAG_OWN_REVOLVING	RE_CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_FNAME_TYPE_S
1	100002	1	Cash loans	M	N	Y	0	202500	406597.5	24700.5	351000 Unaccompanied
2	100003	0	Cash loans	F	N	N	0	270000	1293502.5	35698.5	1129500 Family
3	100004	0	Revolving loans	M	Y	Y	0	67500	135000	6750	135000 Unaccompanied
4	100006	0	Cash loans	F	N	Y	0	135000	312682.5	20645.5	291000 Unaccompanied
5	100007	0	Cash loans	M	N	Y	0	121000	613000	30885.5	515000 Unaccompanied
6	100008	0	Cash loans	M	N	Y	0	99000	498195.5	27511.5	464500 Spouse/partner
7	100009	0	Cash loans	F	Y	Y	1	171000	1560726	41301	1395000 Unaccompanied
8	100010	0	Cash loans	M	Y	Y	0	360000	1530000	42075	1530000 Unaccompanied
9	100011	0	Cash loans	F	N	Y	0	112500	1019610	33826.5	913500 Children
10	100012	0	Revolving loans	M	N	Y	0	135000	405000	20250	405000 Unaccompanied
11	100014	0	Cash loans	F	N	Y	1	112500	652500	21177	652500 Unaccompanied
12	100015	0	Cash loans	F	N	Y	0	38419.155	148365	10679.5	135000 Children
13	100016	0	Cash loans	F	N	Y	0	67500	80865	5881.5	67500 Unaccompanied
14	100017	0	Cash loans	M	Y	N	1	225000	918468	28066.5	697500 Unaccompanied
15	100018	0	Cash loans	F	N	Y	0	189000	773680.5	32778	679500 Unaccompanied
16	100019	0	Cash loans	M	Y	Y	0	157500	299772	20160	247500 Family
17	100020	0	Cash loans	M	N	N	0	108000	509602.5	26149.5	387000 Unaccompanied
18	100021	0	Cash loans	F	N	Y	0	108000	509602.5	26149.5	387000 Unaccompanied

After removing the blanks, duplicates(no duplicates were found in dataset) and adjusting the non-relevant columns we made the dataset to **49999 rows, 45 columns & 2250000 cells** for our final calculation. Few of the columns had no direct use in the analysis but I have kept them as the values were used in correlation calculation for factors impacting defaults.

A	B	C	
Column No.	Table	Row	Description
1	application_data	SK_ID_CURR	ID of loan in our sample
2	application_data	TARGET	Target variable (1 - client with payment difficulties: he/she had late payment more than X days on at least one of the first Y installments of the loan)
3	application_data	NAME_CONTRACT_TYPE	Identification if loan is cash or revolving
4	application_data	CODE_GENDER	Gender of the client
5	application_data	FLAG_OWN_CAR	Flag if the client owns a car
6	application_data	FLAG_OWN_REALTY	Flag if client owns a house or flat
7	application_data	CNT_CHILDREN	Number of children the client has
8	application_data	AMT_INCOME_TOTAL	Income of the client
9	application_data	AMT_CREDIT	Credit amount of the loan
10	application_data	AMT_ANNUITY	Loan annuity
11	application_data	AMT_GOODS_PRICE	For consumer loans it is the price of the goods for which the loan is given
12	application_data	NAME_INCOME_TYPE	Clients income type (businessman, working, maternity leave,...)
13	application_data	NAME_EDUCATION_TYPE	Level of highest education the client achieved
14	application_data	NAME_FAMILY_STATUS	Family status of the client
15	application_data	NAME_HOUSING_TYPE	What is the housing situation of the client (renting, living with parents, ...)
16	application_data	REGION_POPULATION_RELATIVE	Normalized population of region where client lives (higher number means the client lives in more populated region)
17	application_data	DAYS_BIRTH	Client's age in days at the time of application
18	application_data	DAYS_EMPLOYED	How many days before the application the person started current employment
19	application_data	DAYS_REGISTRATION	How many days before the application did client change his registration
20	application_data	OWN_CAR_AGE	Age of client's car
21	application_data	OCCUPATION_TYPE	What kind of occupation does the client have
22	application_data	CNT_FAM_MEMBERS	How many family members does client have
23	application_data	REGION_RATING_CLIENT	Our rating of the region where client lives (1,2,3)
24	application_data	REGION_RATING_CLIENT_W_CITY	Our rating of the region where client lives with taking city into account (1,2,3)
25	application_data	WEEKDAY_APPR_PROCESS_START	On which day of the week did the client apply for the loan
26	application_data	HOUR_APPR_PROCESS_START	Approximately at what hour did the client apply for the loan
27	application_data	ORGANIZATION_TYPE	Type of organization where client works
28	application_data	NAME_TYPE_SUITE	Who was accompanying client when he was applying for the loan
29	application_data		

Figure 1 - Columns From application dataset to work (after cleaning)

For question 5 as we needed to calculate the correlation factor we only needed numerical columns in analysis so I made a separate cleaned data for same by removing the categorical columns and the columns which did not show any direct impact on loan payment as a factor so I trimmed it down to 15 or so columns to do correlation of the top factors that may hinder in the repayment process.

For more data insight on the questions look at the respective questions for screenshot, pivot charts and or pivot tables and descriptive statistics given at necessary places.

## 5) Result:

The Project has given me a good idea about the importance and vast variety of excel usage which helps us to look deep into plain numbers and generate a visually insightful data which can help business to gain knowledge and prepare for future as well as give out trends to focus on from the numerical data. The statistics section has helped me learn about the various concept

---

which are useful for majority of the operation for handling and displaying basic charts and generate a meaningful dashboard as well as the use of Descriptive stats to get more deeper insights in the data. ***The overall result of the project has helped me gain knowledge about the real-life data operations that BFSI*** (Banking Financial services & Insurance) go through to analyze and decide whether to approve loans in future for similar types of clients or to change the amounts by checking for defaults, income and living standards of an individual.

## **Data Cleaning Task:**

As it is asked in the first question to do analysis and find out about the missing values and ways to handle them I will try to put all the necessary steps involved in data cleaning there itself instead of placing a separate block in the explanation.

Will not explain much just have added the links of gif/process. (These were the basic steps to clean and adjust the data)

**Note:** I have linked many files in the report like the separate excel sheet file for questions or High Quality image for graphs or GIF link to showcase cleaning process, also loom ppt links for reference and at the end the drive folder links. So the places highlighted by blue color and underline are hyperlink for sources which either gives the files, video or good quality of the screenshots and can be opened by tapping. Same in the main excel file I have linked the pages in index page for easy referencing so do download the xlsx files as it has pivot connections and other things while the files on google drive automatically converts to googledocs which loses graphs and pivot relations(individual questions might show error as data might be missing thing like that so download the [Main Excel sheet](#) or zip file)

## **Data Analytics Tasks :**

**A) Identify Missing Data and Deal with it Appropriately:** As a data analyst, you come across missing data in the loan application dataset. It is essential to handle missing data effectively to ensure the accuracy of the analysis.

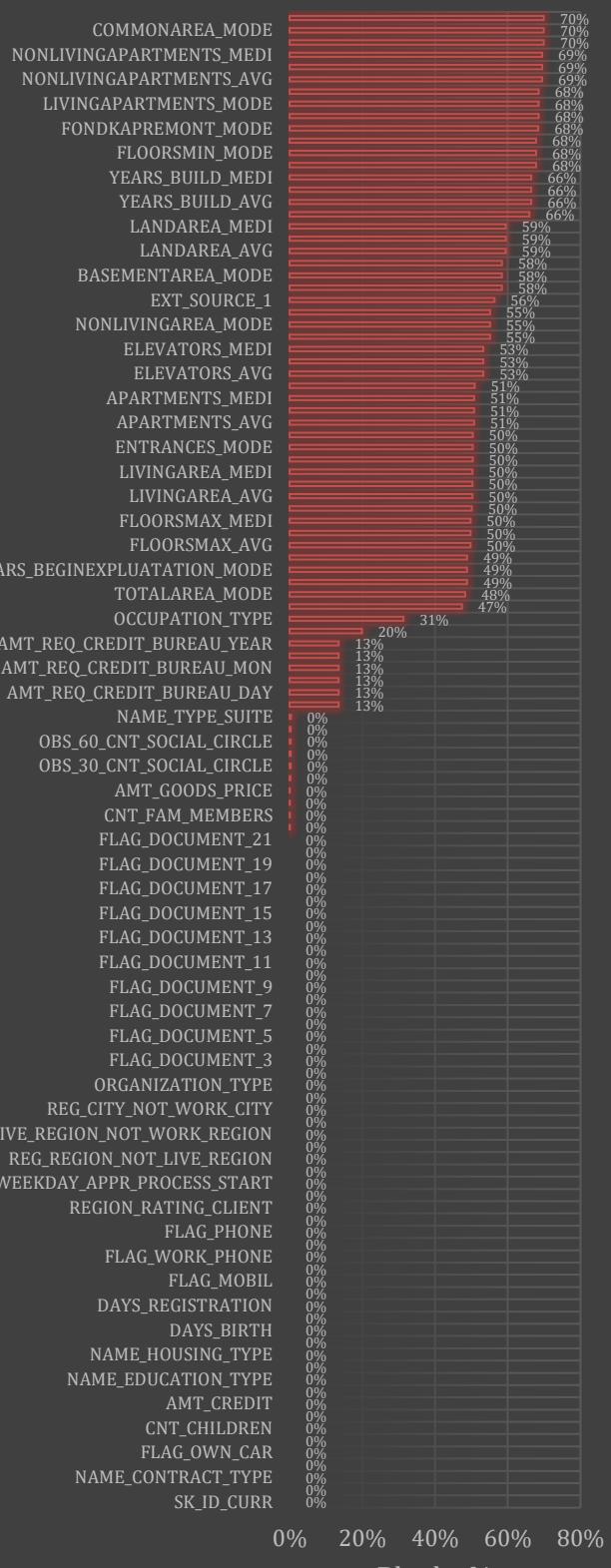
- **Task:** Identify the missing data in the dataset and decide on an appropriate method to deal with it using Excel built-in functions and features.

C93	:	X	✓	f/x	HOUSETYPE_MODE
1	Column No	Table	Row	▼	Description
2	1	application_data	SK_ID_CURR		ID of loan in our sample
3	2	application_data	TARGET		Target variable (1 - client with payment difficulties: he/she had late payment more than X days on at least one of the first Y installments of the
4	5	application_data	NAME_CONTRACT_TYPE		Identification if loan is cash or revolving
5	6	application_data	CODE_GENDER		Gender of the client
6	7	application_data	FLAG_OWN_CAR		Flag if the client owns a car
7	8	application_data	FLAG_OWN_REALTY		Flag if client owns a house or flat
8	9	application_data	CNT_CHILDREN		Number of children the client has
9	10	application_data	AMT_INCOME_TOTAL		Income of the client
10	11	application_data	AMT_CREDIT		Credit amount of the loan
11	12	application_data	AMT_ANNUITY		Loan annuity
12	13	application_data	AMT_GOODS_PRICE		For consumer loans it is the price of the goods for which the loan is given
13	15	application_data	NAME_INCOME_TYPE		Clients income type (businessman, working, maternity leave,...)
14	16	application_data	NAME_EDUCATION_TYPE		Level of highest education the client achieved
15	17	application_data	NAME_FAMILY_STATUS		Family status of the client
16	18	application_data	NAME_HOUSING_TYPE		What is the housing situation of the client (renting, living with parents, ...)
17	19	application_data	REGION_POPULATION_RELATIVE		Normalized population of region where client lives (higher number means the client lives in more populated region)
18	20	application_data	DAYS_BIRTH		Client's age in days at the time of application
19	21	application_data	DAYS_EMPLOYED		How many days before the application the person started current employment
20	22	application_data	DAYS_REGISTRATION		How many days before the application did client change his registration
21	24	application_data	OWN_CAR_AGE		Age of client's car
22	31	application_data	OCCUPATION_TYPE		What kind of occupation does the client have
23	32	application_data	CNT_FAM_MEMBERS		How many family members does client have
24	33	application_data	REGION_RATING_CLIENT		Our rating of the region where client lives (1,2,3)
25	34	application_data	REGION_RATING_CLIENT_W_CITY		Our rating of the region where client lives with taking city into account (1,2,3)
26	35	application_data	WEEKDAY_APPR_PROCESS_START		On which day of the week did the client apply for the loan
27	36	application_data	HOUR_APPR_PROCESS_START		Approximately at what hour did the client apply for the loan
28	43	application_data	ORGANIZATION_TYPE		Type of organization where client works
29	94	application_data	OBS_30_CNT_SOCIAL_CIRCLE		How many observation of client's social surroundings with observable 30 DPD (days past due) default
30	95	application_data	DEF_30_CNT_SOCIAL_CIRCLE		How many observation of client's social surroundings defaulted on 30 DPD (days past due)
31	96	application_data	OBS_60_CNT_SOCIAL_CIRCLE		How many observation of client's social surroundings with observable 60 DPD (days past due) default
32	97	application_data	DEF_60_CNT_SOCIAL_CIRCLE		How many observation of client's social surroundings defaulted on 60 (days past due) DPD

These are the remaining columns in our table after cleaning the missing data like blanks with deletion or appropriate imputation methods. Blanks and missing value shown below with a bar graph for comparison. (there are more columns which can not be fit in the screenshot)

	A	B	C	E	F
1	Table Stats - No of blanks and their percentage to determine whether to impute them or drop the Columns (Blanks above 30 % will be dropped - Red & Orange Color) Also Categorical Column distinction to decide whether to keep or drop them as most are not useful in analysis - 0 or 1 in values or just names.				
2	Column Names	No. Of Blanks	% of blanks	Is Categorical	To Drop or Not
3	SK_ID_CURR	0	0%	Categorical Column	No
4	TARGET	0	0%	Categorical Column	No
5	NAME_CONTRACT_TYPE	0	0%	Categorical Column	No
6	CODE_GENDER	0	0%	Categorical Column	No
7	FLAG_OWN_CAR	0	0%	Categorical Column	No
8	FLAG_OWN_REALTY	0	0%	Categorical Column	No
9	CNT_CHILDREN	0	0%	Categorical Column	No
10	AMT_INCOME_TOTAL	0	0%	Categorical Column	No
11	AMT_CREDIT	0	0%	Categorical Column	No
12	AMT_ANNUITY	1	0%	Not Categorical	No
13	AMT_GOODS_PRICE	38	0%	Not Categorical	No
14	NAME_TYPE_SUITE	192	0%	Not Categorical	No
15	NAME_INCOME_TYPE	0	0%	Categorical Column	No
16	NAME_EDUCATION_TYPE	0	0%	Categorical Column	No
17	NAME_FAMILY_STATUS	0	0%	Categorical Column	No
18	NAME_HOUSING_TYPE	0	0%	Categorical Column	No
19	REGION_POPULATION_RELATIVE	0	0%	Categorical Column	No
20	DAYS_BIRTH	0	0%	Categorical Column	No
21	DAYS_EMPLOYED	0	0%	Categorical Column	No
22	DAYS_REGISTRATION	0	0%	Categorical Column	No
23	DAYS_ID_PUBLISH	0	0%	Categorical Column	No
24	FLAG_MOBIL	0	0%	Categorical Column	No
25	FLAG_EMP_PHONE	0	0%	Categorical Column	No
26	FLAG_WORK_PHONE	0	0%	Categorical Column	No
27	FLAG_CONT_MOBILE	0	0%	Categorical Column	No
28	FLAG_PHONE	0	0%	Categorical Column	No
29	FLAG_EMAIL	0	0%	Categorical Column	No
30	OCCUPATION_TYPE	15654	31%	Not Categorical	Yes
31	CNT_FAM_MEMBERS	1	0%	Not Categorical	No
32	REGION_RATING_CLIENT	0	0%	Categorical Column	No
33	REGION_RATING_CLIENT_W_CITY	0	0%	Categorical Column	No
34	WEEKDAY_APPR_PROCESS_START	0	0%	Categorical Column	No
35	HOUR_APPR_PROCESS_START	0	0%	Categorical Column	No
36	REG_REGION_NOT_LIVE_REGION	0	0%	Categorical Column	No
37	REG_REGION_NOT_WORK_REGION	0	0%	Categorical Column	No
38	LIVE_REGION_NOT_WORK_REGION	0	0%	Categorical Column	No
39	REG_CITY_NOT_LIVE_CITY	0	0%	Categorical Column	No
40	REG_CITY_NOT_WORK_CITY	0	0%	Categorical Column	No
41	ORGANIZATION_TYPE	0	0%	Categorical Column	No
42	EXT_SOURCE_1	28172	56%	Not Categorical	Yes
43	EXT_SOURCE_2	126	0%	Not Categorical	No
44	EXT_SOURCE_3	9944	20%	Not Categorical	No
45	APARTMENTS_AVG	25385	51%	Not Categorical	Yes
46	BASEMENTAREA_AVG	29199	58%	Not Categorical	Yes
47	BALCONY_AVG	25344	49%	Not Categorical	Yes
48	KITCHENLIVINGPULATION_AVG	33239	66%	Not Categorical	Yes
49	YEARS_BUILD_AVG	33239	66%	Not Categorical	Yes
50	COMMONAREA_AVG	34960	70%	Not Categorical	Yes
51	ENTRANCES_AVG	26651	53%	Not Categorical	Yes
52	FLOORSMAX_AVG	25195	50%	Not Categorical	Yes
53	FLOORSMIN_AVG	24875	50%	Not Categorical	Yes
54	LANDAREAVG	33894	68%	Not Categorical	Yes
55	LANDAREA_AVG	29721	59%	Not Categorical	Yes
56	LIVINGAREA_AVG	34226	68%	Not Categorical	Yes
57	LIVINGAPARTMENTS_AVG	25137	50%	Not Categorical	Yes
58	LIVINGAPARTMENTS_MODE	34714	69%	Not Categorical	Yes
59	NONLIVINGAPARTMENTS_AVG	34714	69%	Not Categorical	Yes
60	NONLIVINGAPARTMENTS_MODE	27572	55%	Not Categorical	Yes
61	NONLIVINGAREA_AVG	25385	51%	Not Categorical	Yes
62	NONLIVINGAREA_MODE	29199	58%	Not Categorical	Yes
63	APARTMENTS_MODE	25385	51%	Not Categorical	Yes
64	BASEMENTAREA_MODE	29199	58%	Not Categorical	Yes
65	YEARS_BEGINEXPLUATATION_MODE	24394	49%	Not Categorical	Yes
66	YEARS_BUILD_MODE	33239	66%	Not Categorical	Yes
67	COMMONAREA_MODE	34960	70%	Not Categorical	Yes
68	ENTRANCES_MODE	26651	53%	Not Categorical	Yes
69	FLOORSMAX_MODE	25195	50%	Not Categorical	Yes
70	FLOORSMIN_MODE	24875	50%	Not Categorical	Yes
71	LANDAREAMODE	33894	68%	Not Categorical	Yes
72	LANDAREA_MODE	29721	59%	Not Categorical	Yes
73	LIVINGAPARTMENTS_MODE	34226	68%	Not Categorical	Yes
74	LIVINGAREA_MODE	25137	50%	Not Categorical	Yes
75	NONLIVINGAPARTMENTS_MODE	34714	69%	Not Categorical	Yes
76	NONLIVINGAREA_MODE	27572	55%	Not Categorical	Yes
77	APARTMENTS_MEDI	25385	51%	Not Categorical	Yes
78	BASEMENTAREA_MEDI	29199	58%	Not Categorical	Yes
79	YEARS_BEGINEXPLUATATION_MEDI	24394	49%	Not Categorical	Yes
80	YEARS_BUILD_MEDI	33239	66%	Not Categorical	Yes
81	COMMONAREA_MEDI	34960	70%	Not Categorical	Yes
82	ELEVATORS_MEDI	26651	53%	Not Categorical	Yes
83	ENTRANCES_MEDI	25195	50%	Not Categorical	Yes
84	FLOORSMAX_MEDI	24875	50%	Not Categorical	Yes
85	FLOORSMIN_MEDI	33894	68%	Not Categorical	Yes
86	LANDAREAMEDI	29721	59%	Not Categorical	Yes
87	LIVINGAPARTMENTS_MEDI	34226	68%	Not Categorical	Yes
88	LIVINGAREA_MEDI	25137	50%	Not Categorical	Yes
89	NONLIVINGAPARTMENTS_MEDI	34714	69%	Not Categorical	Yes
90	NONLIVINGAREA_MEDI	27572	55%	Not Categorical	Yes
91	FONDKAPREMONT_MODE	34191	68%	Not Categorical	Yes
92	HOUSETYPE_MODE	25075	50%	Not Categorical	Yes
93	TOTALAREA_MODE	24148	48%	Not Categorical	Yes
94	WALLSMATERIAL_MODE	25459	51%	Not Categorical	Yes
95	EMERGENCYSTATE_MODE	23698	47%	Not Categorical	Yes
96	OBS_30_CNT_SOCIAL_CIRCLE	168	0%	Not Categorical	No
97	DEF_30_CNT_SOCIAL_CIRCLE	168	0%	Not Categorical	No
98	OBS_60_CNT_SOCIAL_CIRCLE	168	0%	Not Categorical	No
99	DEF_60_CNT_SOCIAL_CIRCLE	168	0%	Not Categorical	No
100	DAYS_LAST_PHONE_CHANGE	1	0%	Not Categorical	No
101	FLAG_DOCUMENT_2	0	0%	Categorical Column	No
102	FLAG_DOCUMENT_3	0	0%	Categorical Column	No
103	FLAG_DOCUMENT_4	0	0%	Categorical Column	No
104	FLAG_DOCUMENT_5	0	0%	Categorical Column	No
105	FLAG_DOCUMENT_6	0	0%	Categorical Column	No
106	FLAG_DOCUMENT_7	0	0%	Categorical Column	No
107	FLAG_DOCUMENT_8	0	0%	Categorical Column	No
108	FLAG_DOCUMENT_9	0	0%	Categorical Column	No
109	FLAG_DOCUMENT_10	0	0%	Categorical Column	No
110	FLAG_DOCUMENT_11	0	0%	Categorical Column	No
111	FLAG_DOCUMENT_12	0	0%	Categorical Column	No
112	FLAG_DOCUMENT_13	0	0%	Categorical Column	No
113	FLAG_DOCUMENT_14	0	0%	Categorical Column	No
114	FLAG_DOCUMENT_15	0	0%	Categorical Column	No
115	FLAG_DOCUMENT_16	0	0%	Categorical Column	No
116	FLAG_DOCUMENT_17	0	0%	Categorical Column	No
117	FLAG_DOCUMENT_18	0	0%	Categorical Column	No
118	FLAG_DOCUMENT_19	0	0%	Categorical Column	No
119	FLAG_DOCUMENT_20	0	0%	Categorical Column	No
120	FLAG_DOCUMENT_21	0	0%	Categorical Column	No
121	AMT_REQ_CREDIT_BUREAU_HOUR	6734	13%	Not Categorical	No
122	AMT_REQ_CREDIT_BUREAU_DAY	6734	13%	Not Categorical	No
123	AMT_REQ_CREDIT_BUREAU_WEEK	6734	13%	Not Categorical	No
124	AMT_REQ_CREDIT_BUREAU_MON	6734	13%	Not Categorical	No
125	AMT_REQ_CREDIT_BUREAU_QRT	6734	13%	Not Categorical	No
126	AMT_REQ_CREDIT_BUREAU_YEAR	6734	13%	Not Categorical	No

## % of Blank Columns - Current Application



Since the Pivot table and Bar graph is lengthy I have captured a gif to show the image here also you can find the same in attached xls file. > [Gif](#)

### **Insights:**

1. For missing value detection I first did a simple column statistics using google sheet to get info about the empty cells count of whole application dataset and understood the column relationship. For Duplicates & blanks we can use the data tabs feature to look for duplicates but there were none.

> [ColumnStats Gif](#),

> [Duplicate Values , Dups](#)

2. For missing values we used basic excel functions like **COUNT** (to get details about filled cells), **ISBLANK** or **COUNTBLANK** (to get the details about blank cells in a column). Using the counts of blank we determined the % of blanks for each column and then decided to treat the values *above 35%* and drop them for rest we did mean mode median imputation as applicable.

H	I	J	K	L	M	N	O	P	Q	R	S	T
Blanks	0	0	0	0	0	0	0	0	0	1	38	192
% Blanks	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%
Mean	129013.2106	0.0805216	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	0.419848397	170767.5905	599700.5815	27107.37736	539060.0361	#DIV/0!
Mode	#N/A	0	#N/A	#N/A	#N/A	#N/A	0	135000	450000	9000	450000	#N/A
Median	129076	0	#NUM!	#NUM!	#NUM!	#NUM!	0	145800	514777.5	24939	450000	#NUM!
SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	NAME_TYPE_SUITE	
100002	1	Cash loans	M	N	Y	0	202500	406597.5	24700.5	351000	Unaccompanied	
100003	0	Cash loans	F	N	N	0	270000	1293502.5	35698.5	1129500	Family	
												Sumif S_C

2 - For the categorical data directly doing stats will show some error like div by or n/a or num! but as most of them had no missing value that needed treatment below 35% so its fine to ignore them the formula was used to get the data for numerical missing data points calculation.

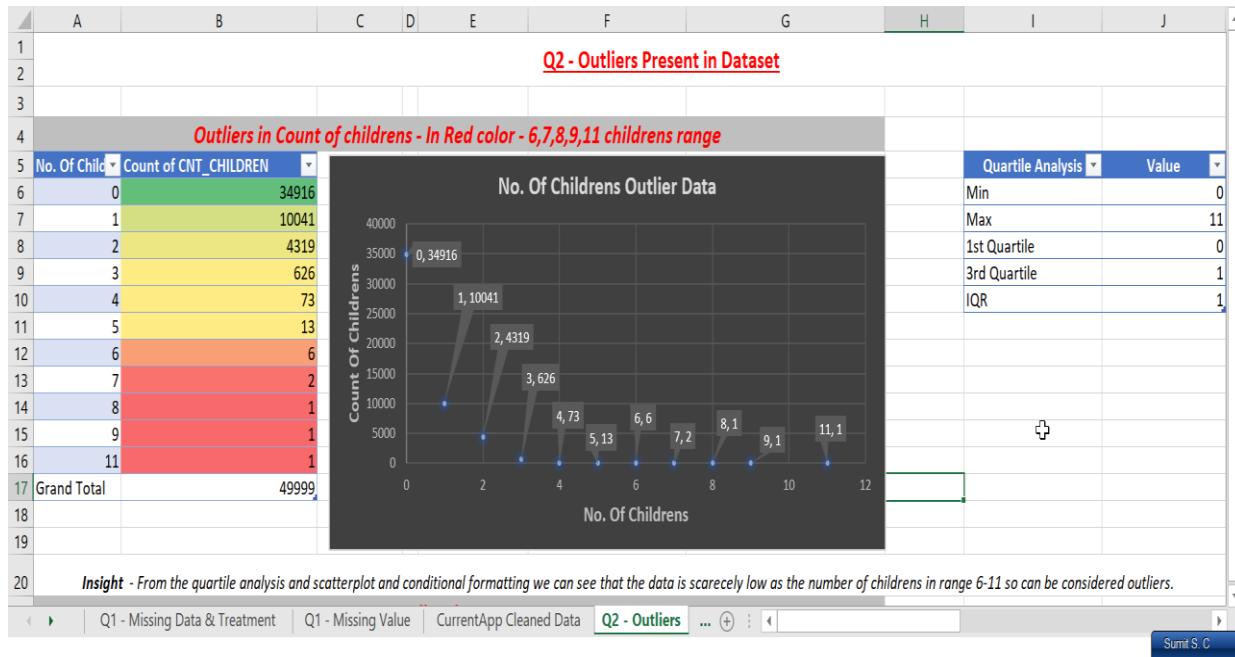
3. Looking at the Bar Graph of the missing data and doing some filtering using conditional formatting and **IF** formula we sorted the columns to be *dropped above 35%* (marked red in pivot table) and then the rest of the columns with missing values were treated and the final [application dataset file](#) (Cleaned Dataset) was made to do analysis with necessary columns.

4. Files

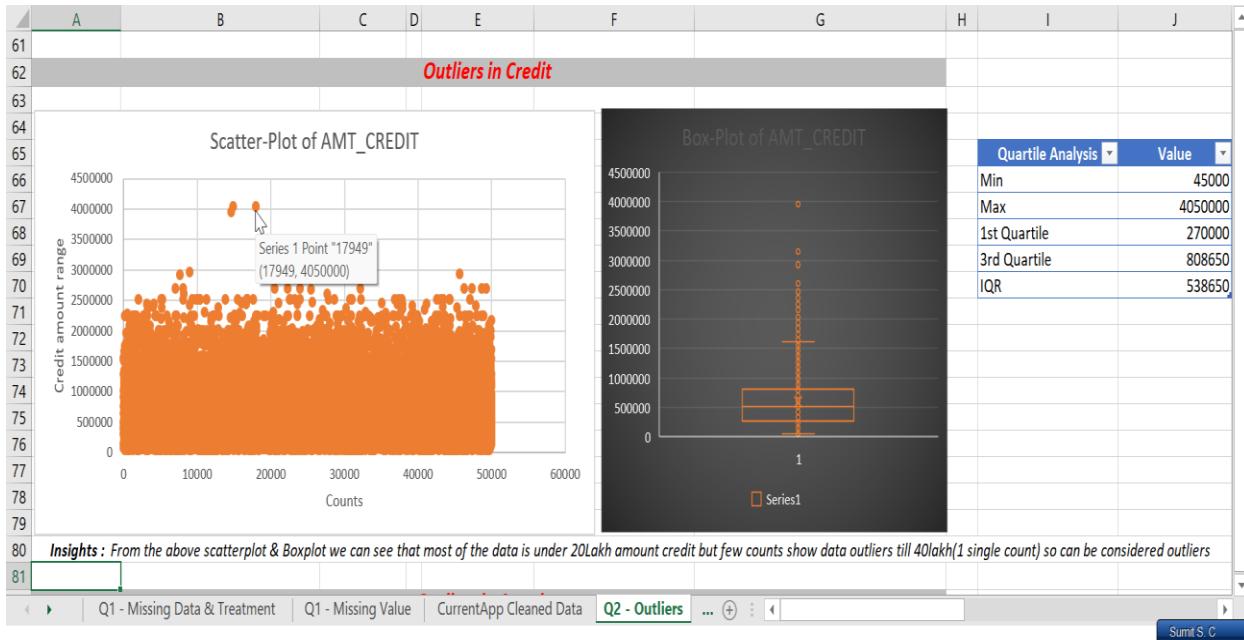
- Q 1 [Excel File](#), [Docs File](#)
- Loom Video Q1 – [Identifying Missing Values & Treating them.](#)

**B) Identify Outliers in the Dataset:** Outliers can significantly impact the analysis and distort the results. You need to identify outliers in the loan application dataset.

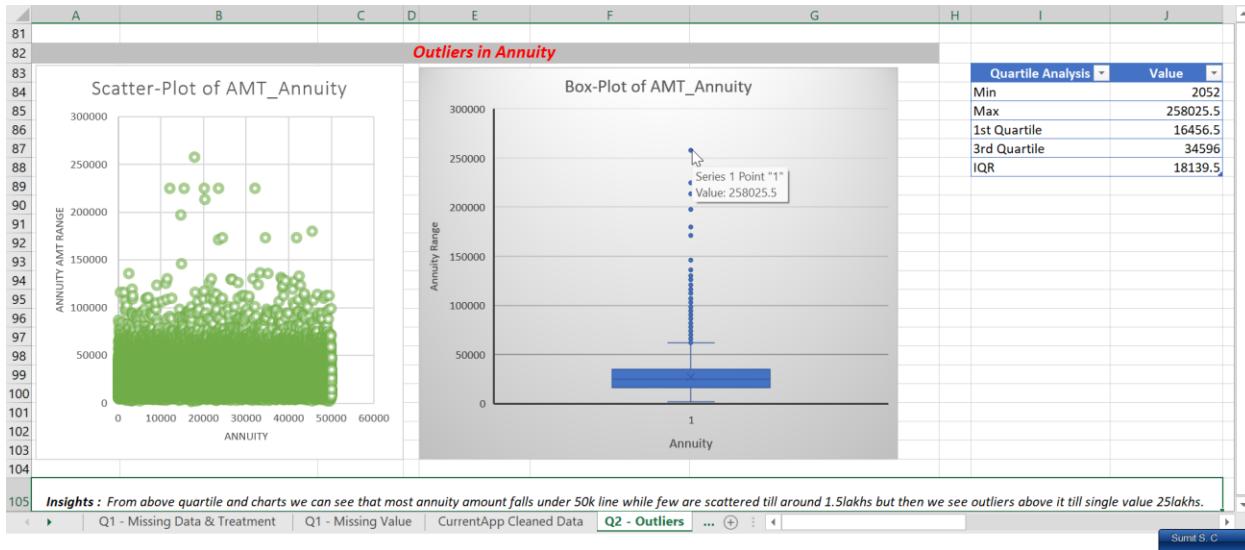
- **Task:** Detect and identify outliers in the dataset using Excel statistical functions and features, focusing on numerical variables.



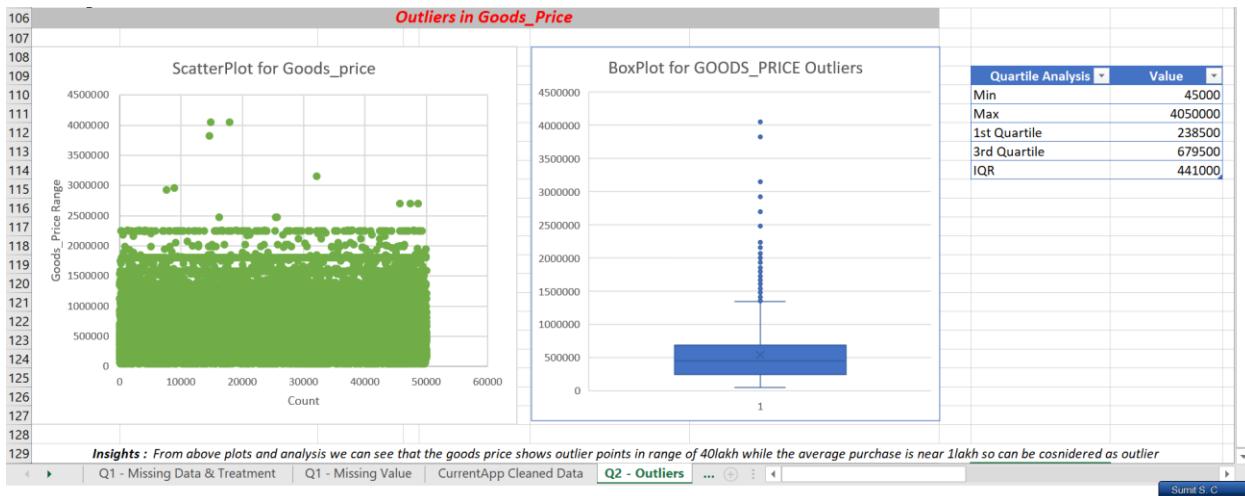
#### Outlier 1 - In No. Of childrens



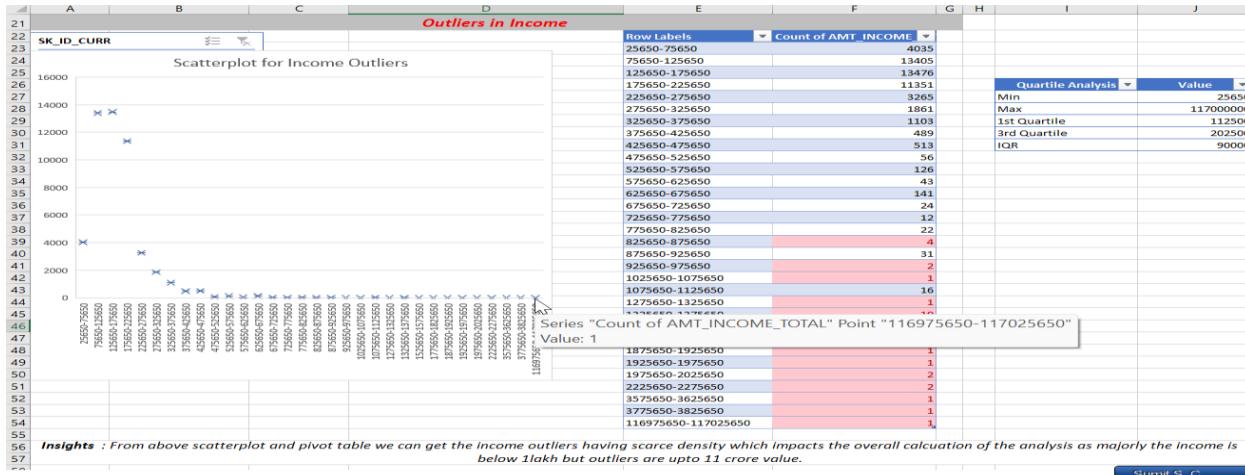
#### Outlier 2 - In Credit\_AMT



### Outlier 3 In Annuity Amount



### Outlier 4 - In Goods\_Price Column

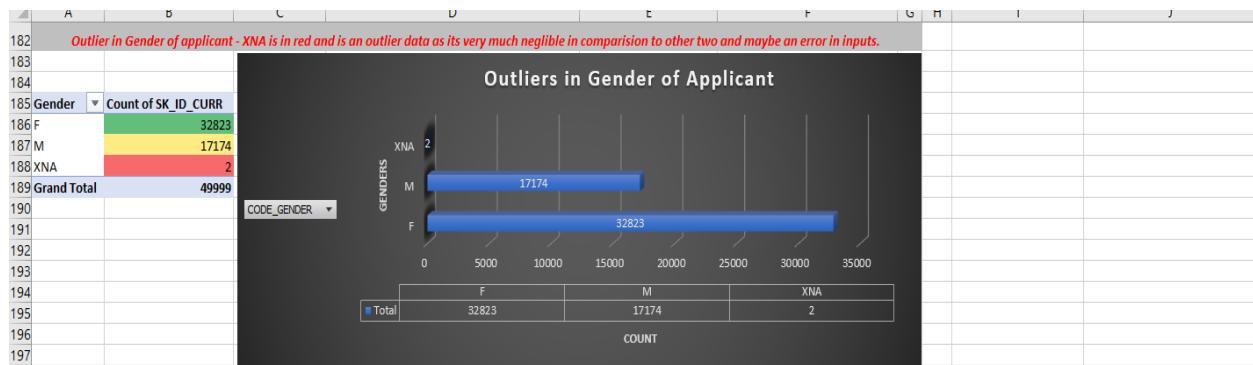


### Outlier 5 - Outliers In income

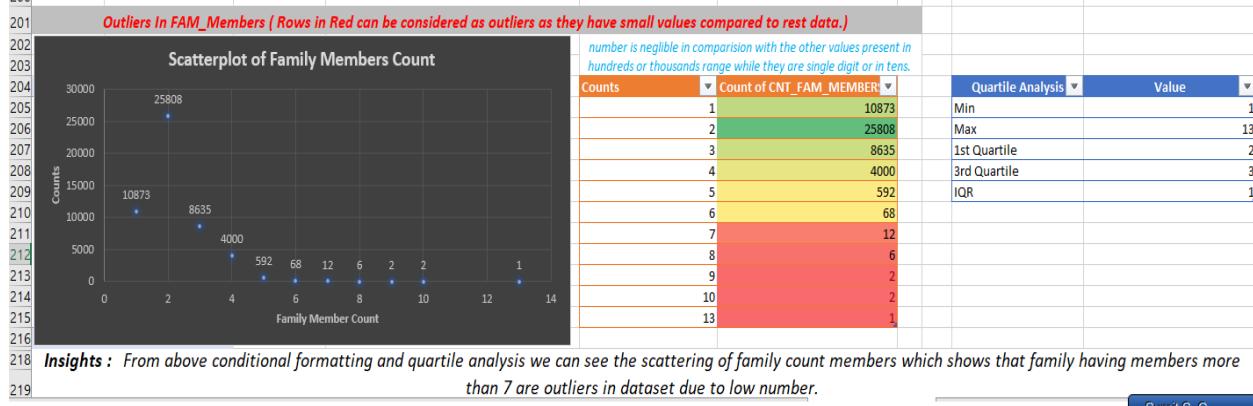


Sumit S. C

### Outlier 6 - Outliers in Tenure (In employment years and from application day)



*Insights : From the gender column we can see that the XNA group has only 2 entry and can be considered as an outlier on lower range or an error in data entry.*



Sumit S. C

### Outlier 7 - Inside Gender distribution & Family members count

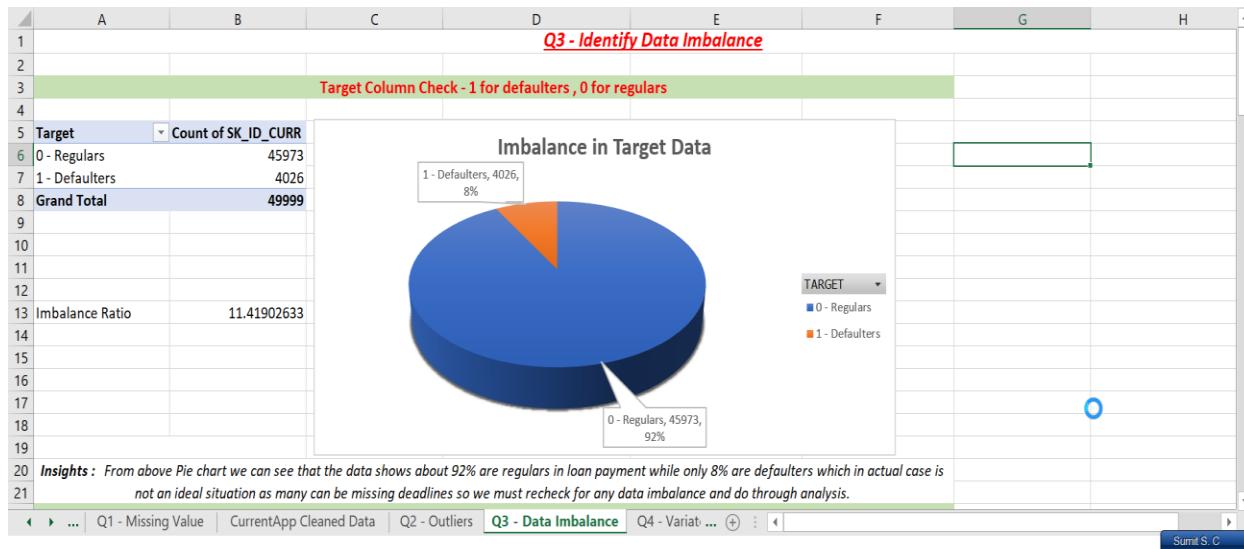
---

### **Insights:**

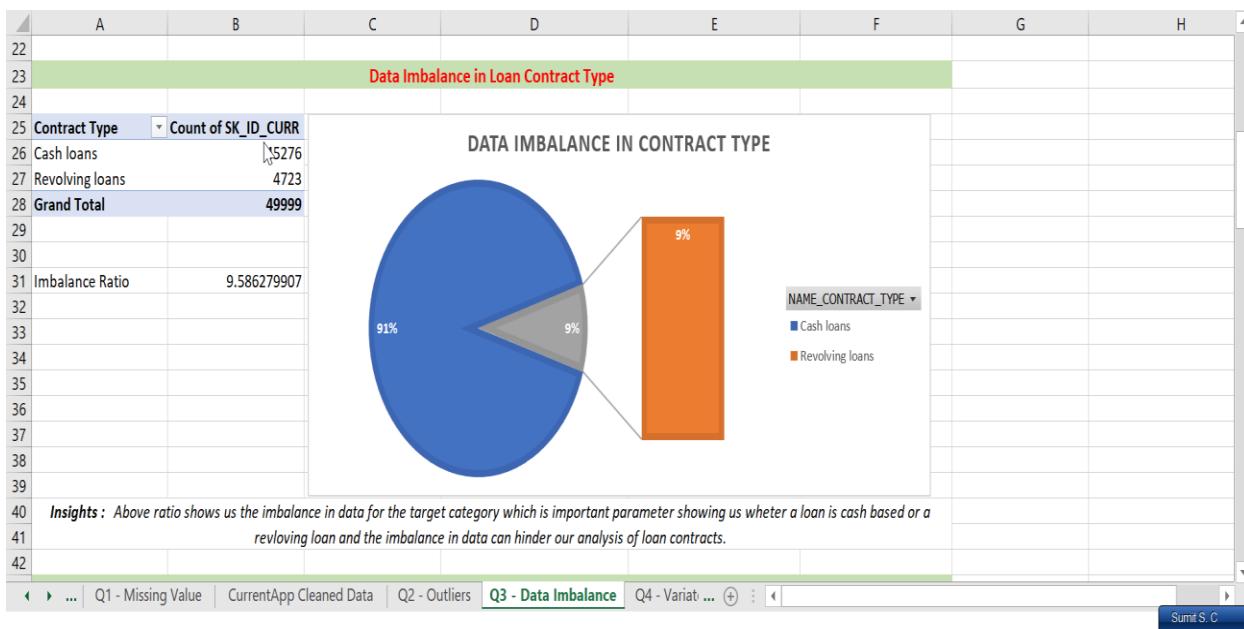
1. For every outlier I have tried giving the insight in the screenshot but to summarize the overall data we have found out that many outliers data are present in the dataset that needs to be checked or appropriately dealt with to make the analysis and overall future predictions about customer loan repayments and factors affecting their loans.
2. Certain factors like income and annuity shows us outlier points because 1 or 2 data entry shows a huge number of loan amount compared to rest of the data , same with the outliers present in tenure of loan and employment factor analysis would help us gain knowledge about trends in different category of profession and income range and corresponding habits etc so these outliers should be dealt accordingly to make analysis more righteous.
3. Few other outlier points are mentioned in screenshot and in excel report so you can find more insights from there.
4. Links
  - Q2 Excel sheet (Individual file crashes the excel program so look into [main excel file](#) and go to Q2 tab)
  - Loom Video Explanation ([Identifying Outliers](#))

**C) Analyze Data Imbalance:** Data imbalance can affect the accuracy of the analysis, especially for binary classification problems. Understanding the data distribution is crucial for building reliable models.

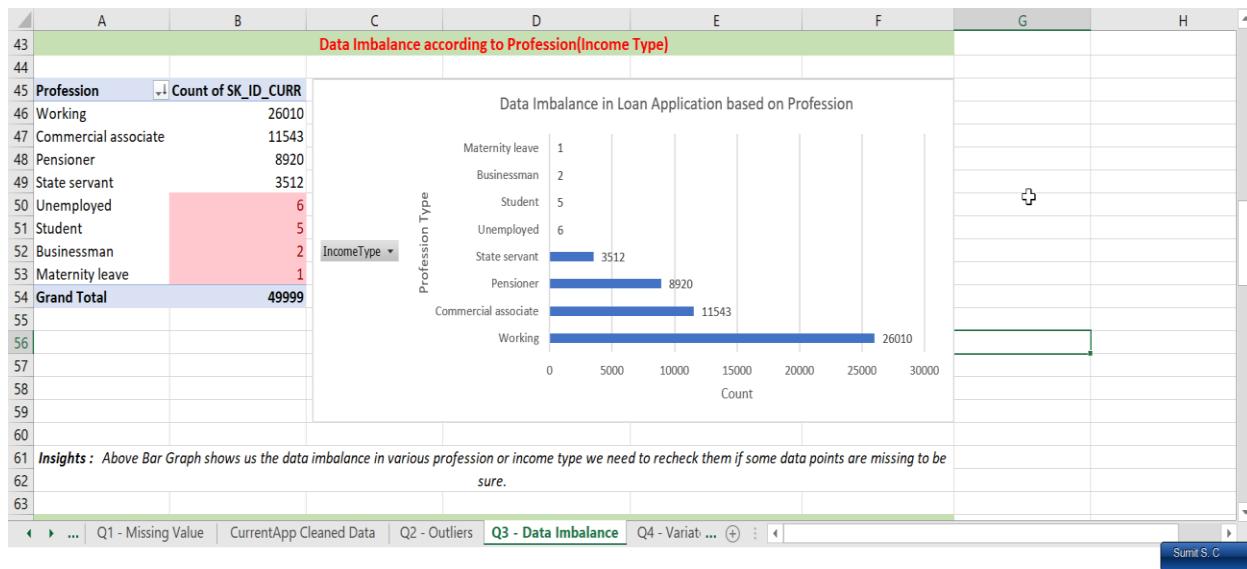
- **Task:** Determine if there is data imbalance in the loan application dataset and calculate the ratio of data imbalance using Excel functions.



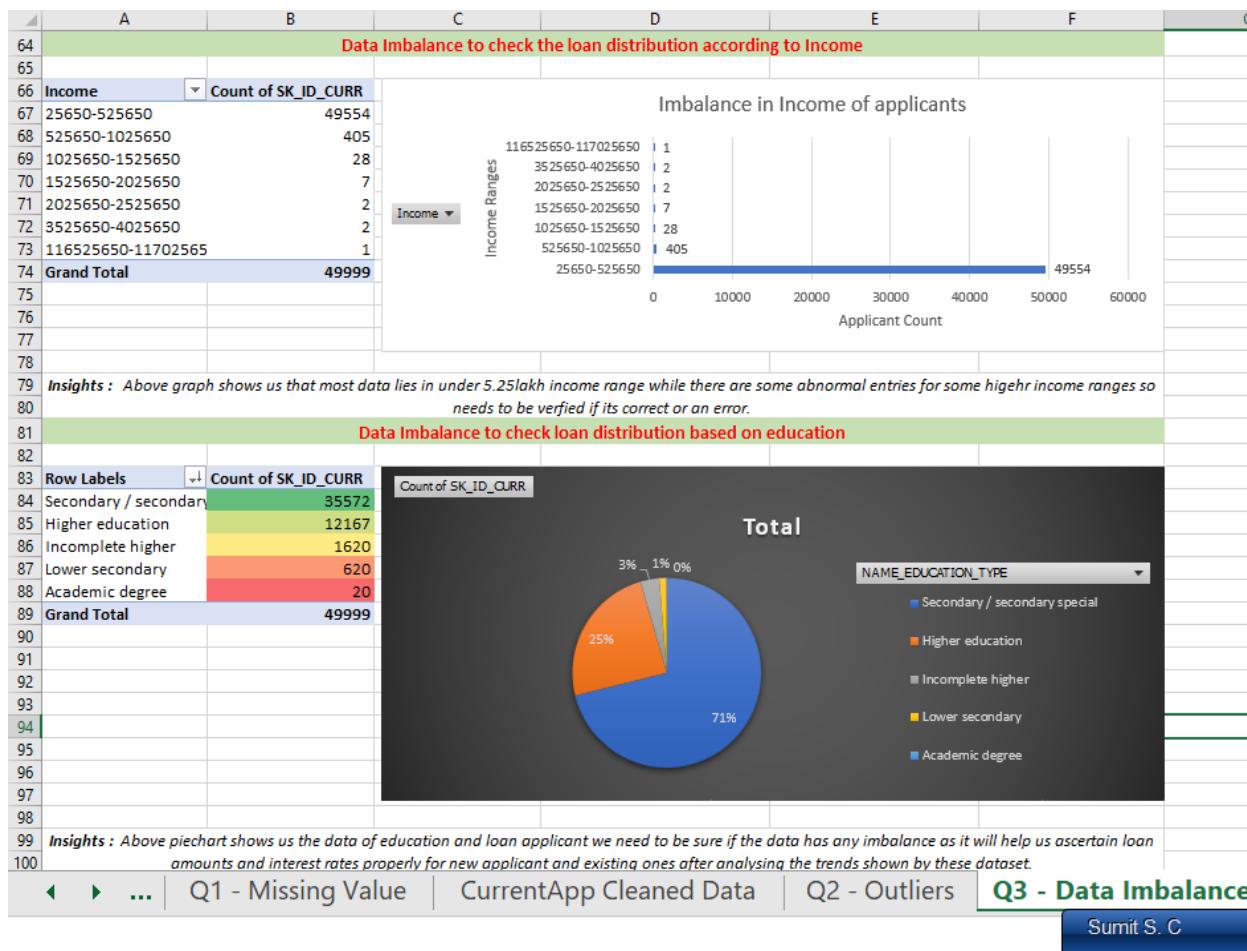
*Data Imbalance 1 - [In target Column](#)*



*Data Imbalance 2 - [In Target Category](#)*



Data Imbalance 3 - In Income Type <https://mitsus.life-is-pa.in/6gMrRvcDW.png>



#### Data Imbalance 4 - In education & Income range

---

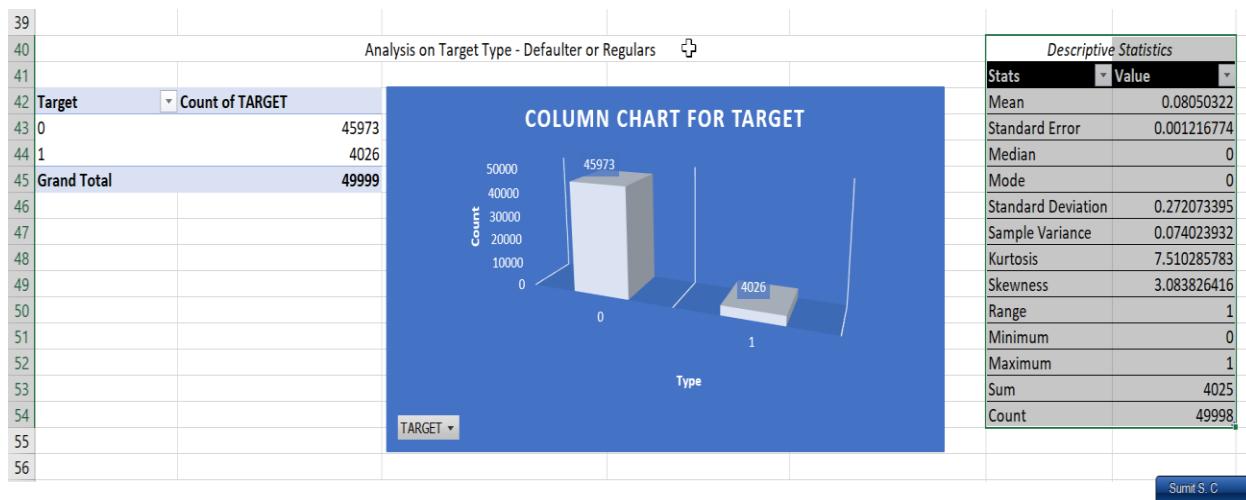
### **Insights:**

1. From the above analysis which we did using pivot table and basic functions like **COUNT, COUNTIF & SUM** we got the data imbalance.
2. For better results and analysis we must recheck the data if the given data inputs are valid and are not having error will help us to get insights which will help us determine certain factors like age, profession, income and contract types which might hamper the payment or our overall analysis if some data are incorrect.
3. Data imbalance and its ratio helps us to prevent any unwanted calculation or change the analysis due to certain factors not being entered or taken correctly. It act as an indicator for certain values like yes or no , 1 – 0 columns which while entering or copying might have some human error and which might impact overall dataset working.
4. Files
  - Q3 [Excel File](#), [Gdocs File](#)
  - Loom Video PPT [Q3 – Identify the data imbalance](#)

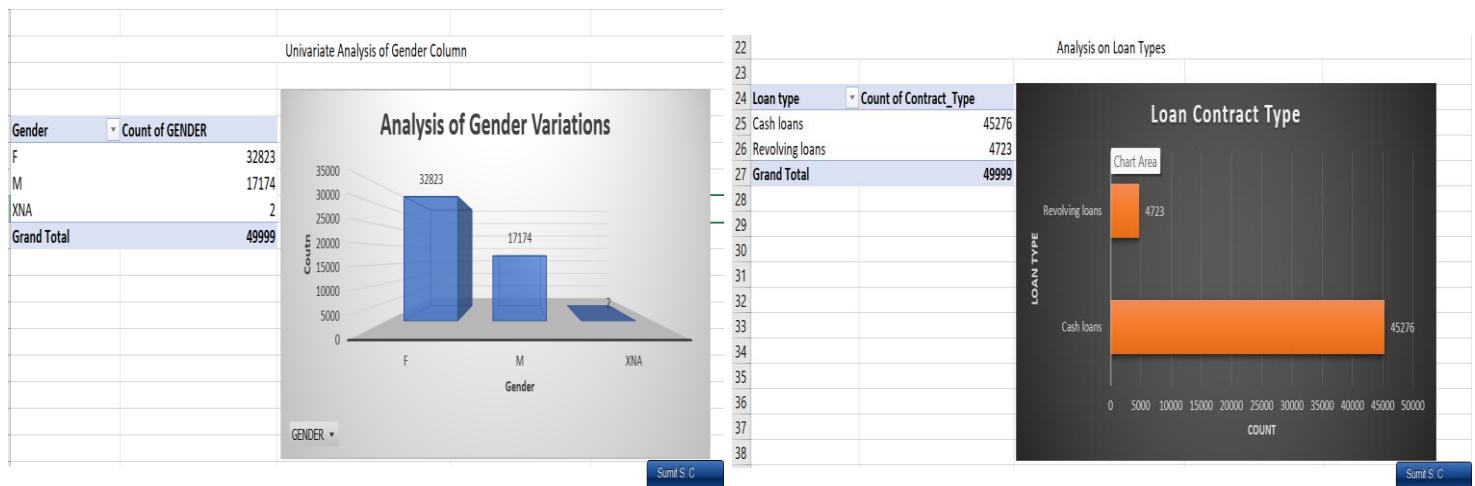
**D) Perform Univariate, Segmented Univariate, and Bivariate Analysis:** To gain insights into the driving factors of loan default, it is important to conduct various analyses on consumer and loan attributes.

- **Task:** Perform univariate analysis to understand the distribution of individual variables, segmented univariate analysis to compare variable distributions for different scenarios, and bivariate analysis to explore relationships between variables and the target variable using Excel functions and features.

**Univariate Analysis – When we examine and analyse a single variable in isolation and understand the distribution summary stats using descriptive stats and try to gain information its called as univariate analysis.**

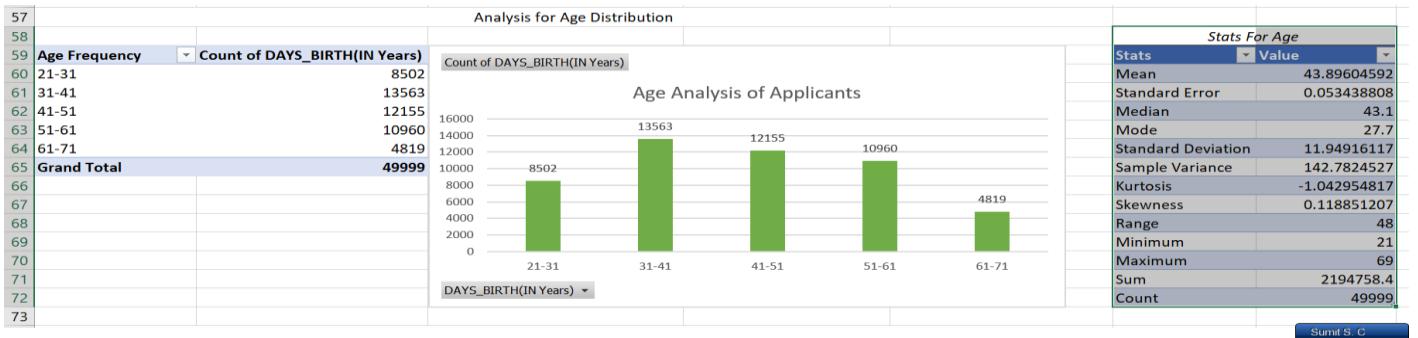


Univariate 1 - On Target



Univariate 2 On Gender

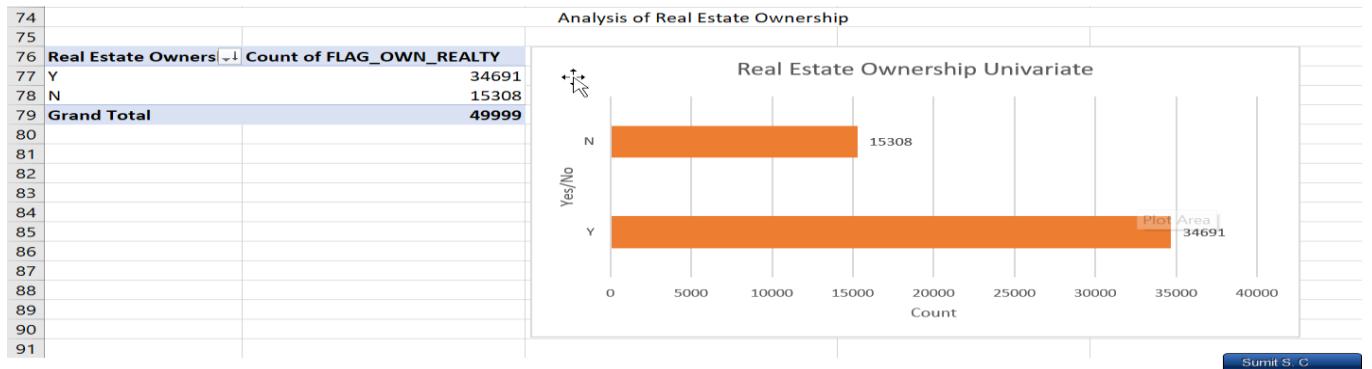
Univariate 3 On Loan type



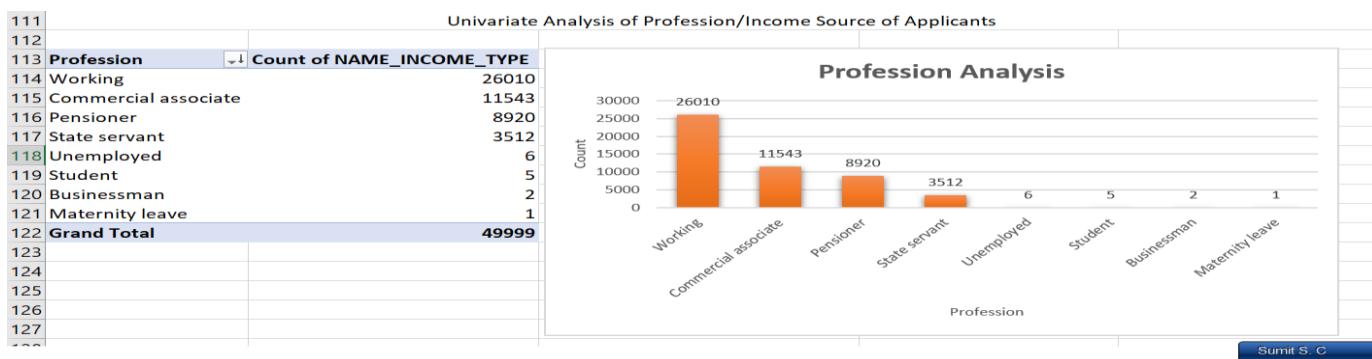
### Univariate 3 On Age



### Univariate 4 On Loan Amount

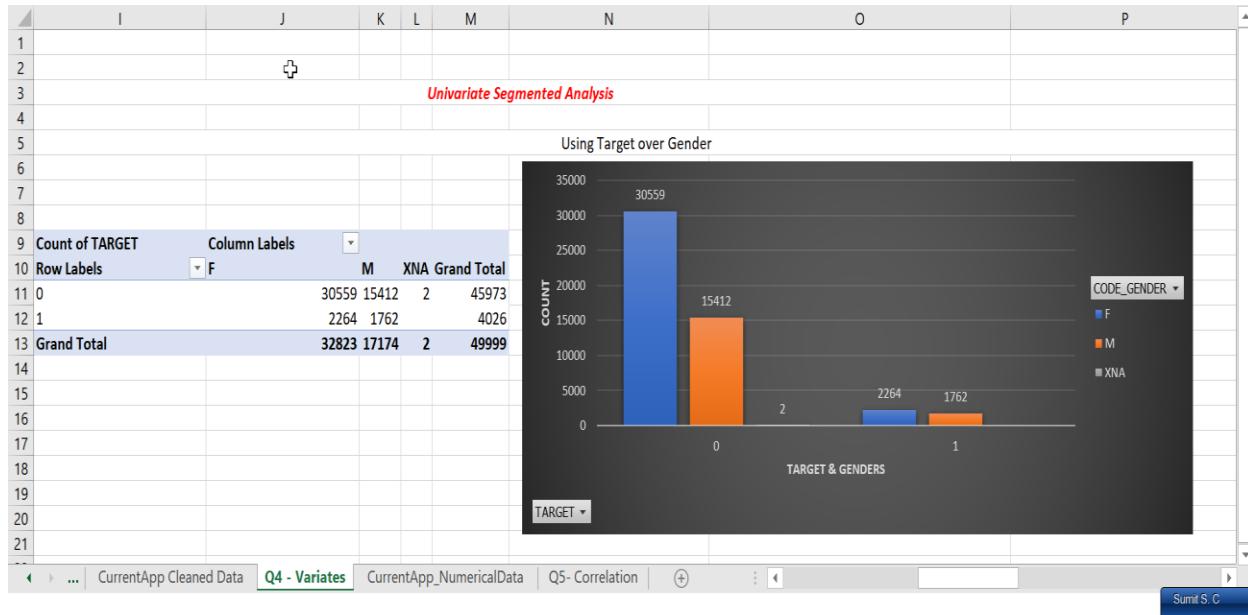


### Univariate 5 on estate ownership

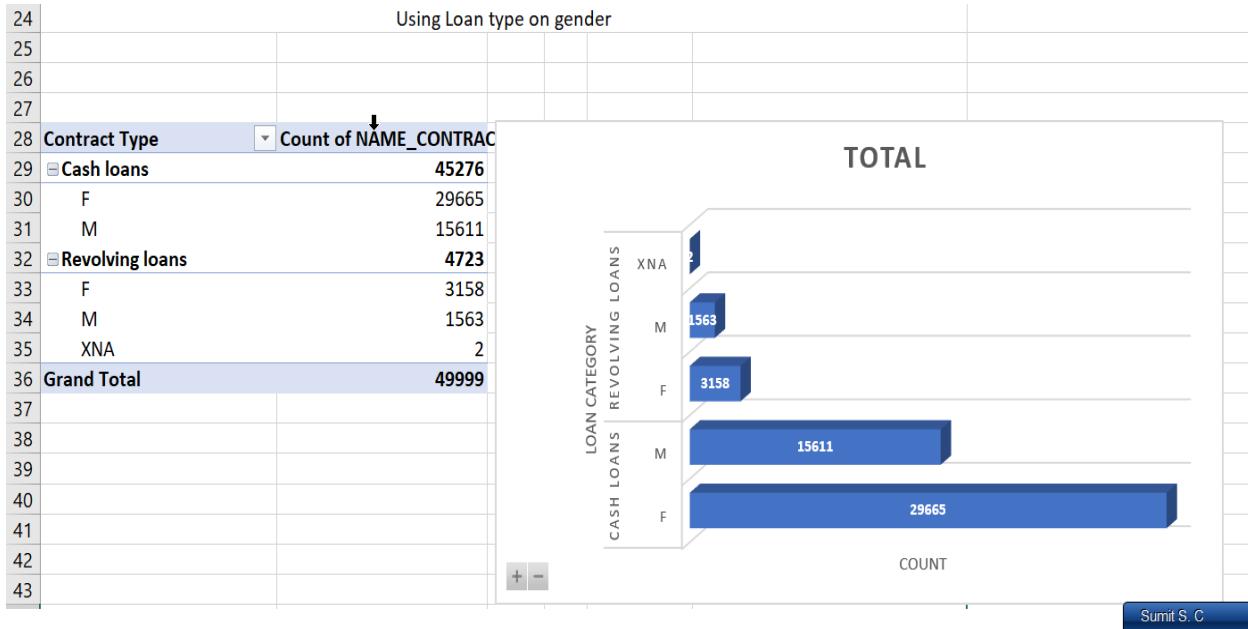


### Univariate 6 On Income Type

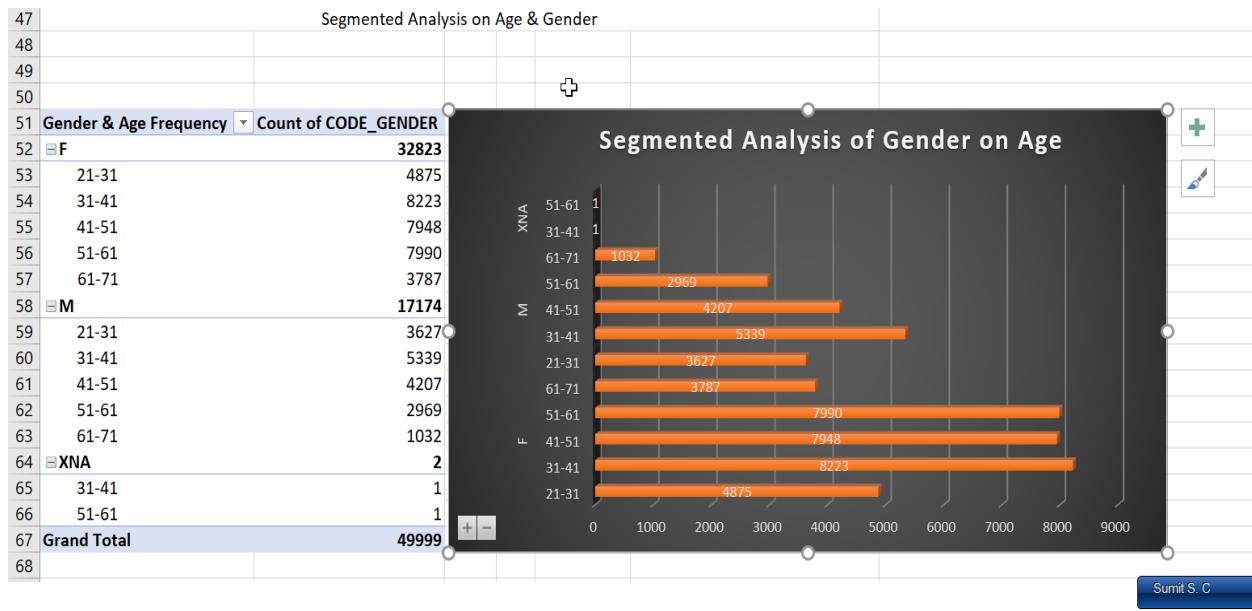
**Segmented Univariate Analysis** – When we examine and analyse a single variable with different segments or subgroups defined by another variable and understand the distribution summary stats using descriptive stats and try to gain information its called as segmented univariate analysis.



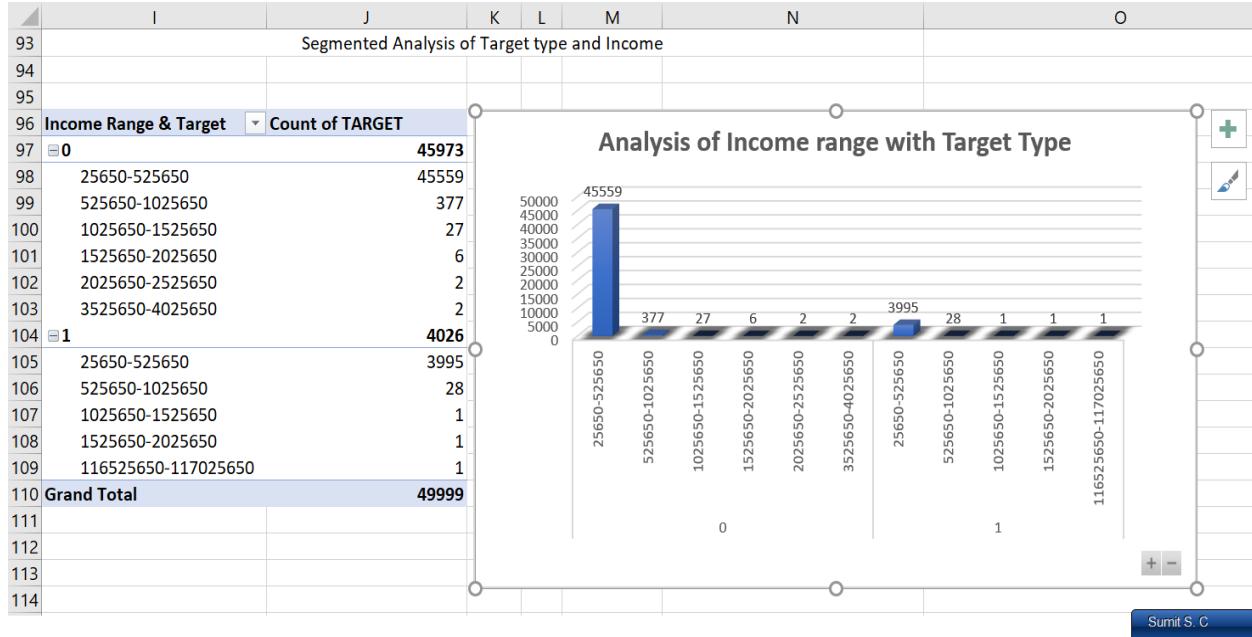
#### Univariate Segmented 1 - On Target & Gender



#### Univariate Segmented 2 - On Loan type & gender

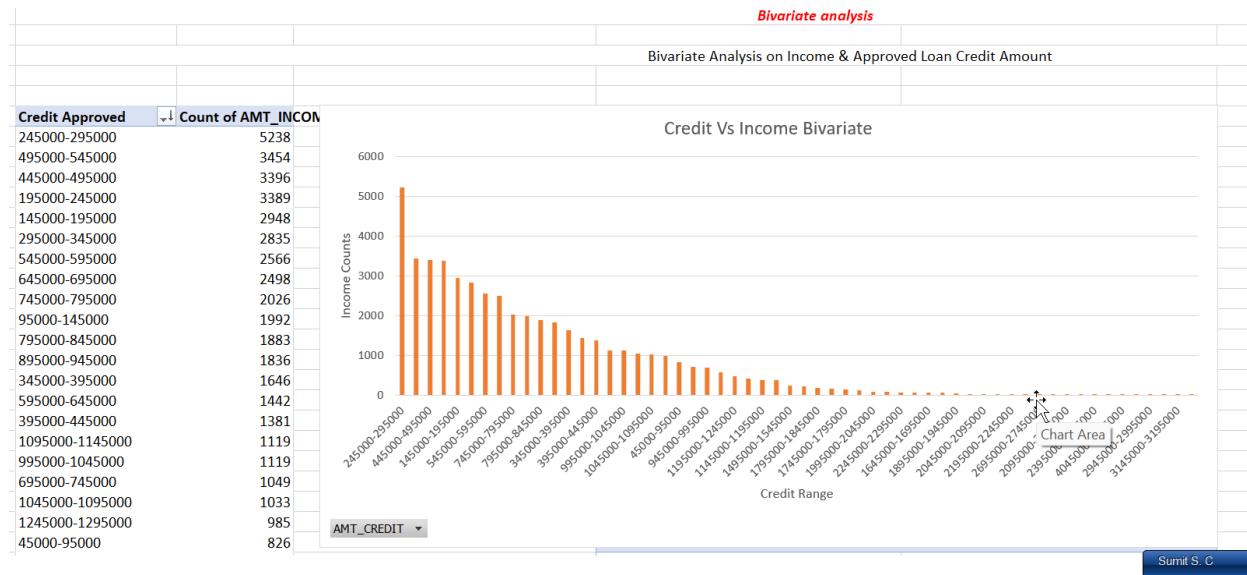


#### Univariate Segmented 3 - On Gender & Age

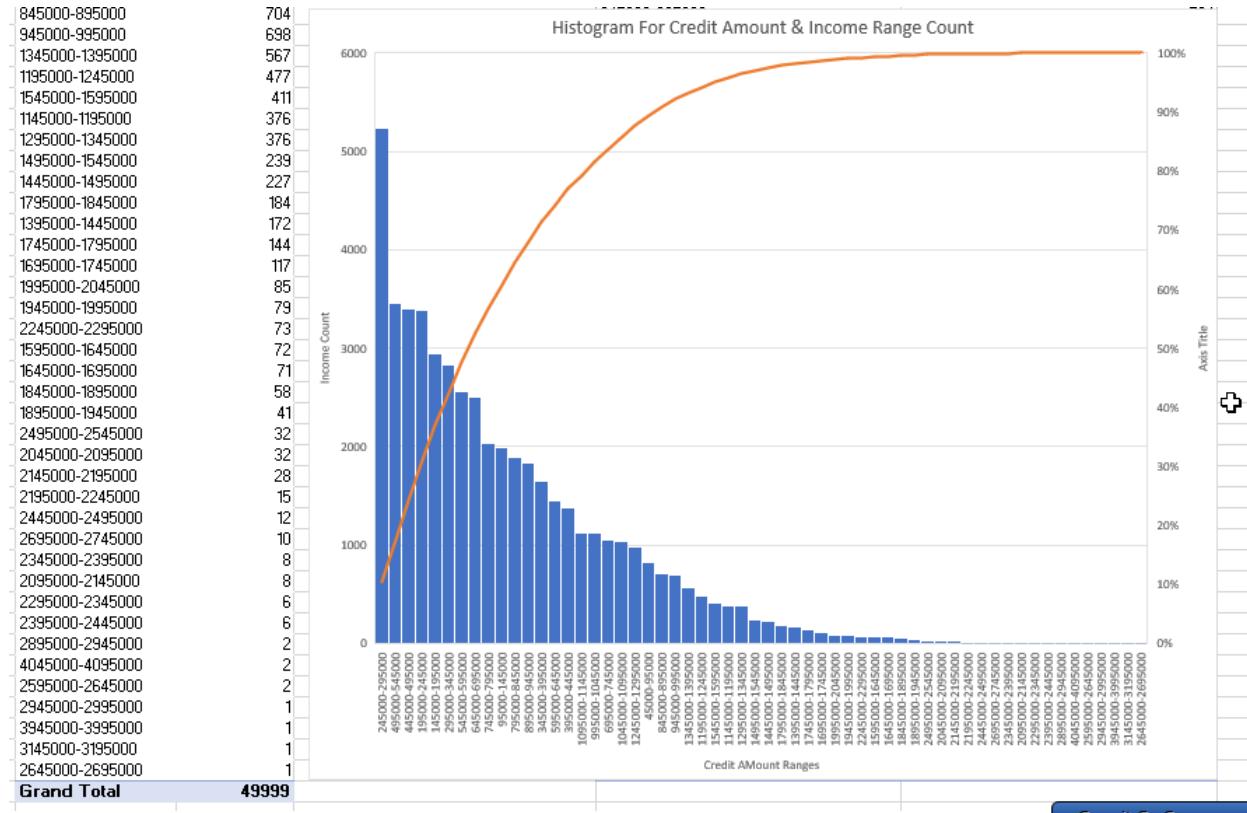


#### Univariate Segmented 4 - On target & Income Range

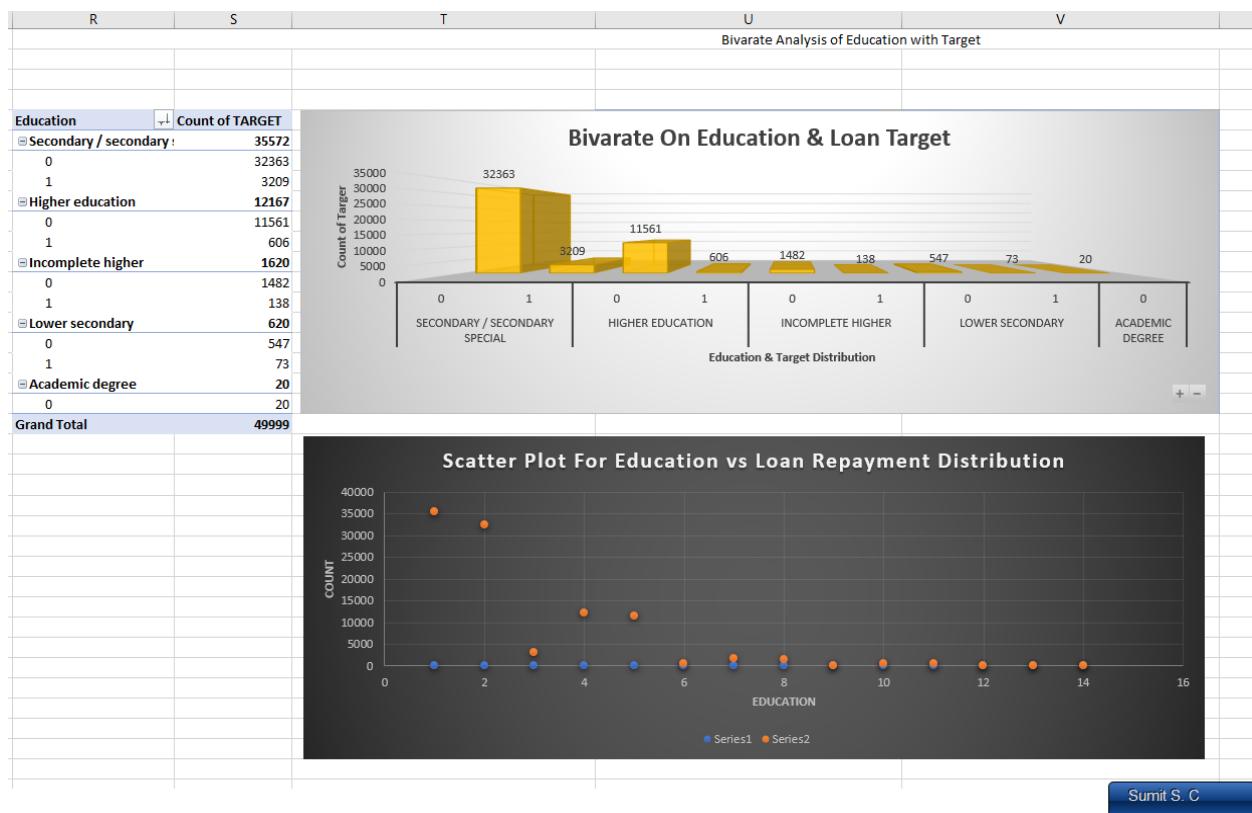
**Bivariate Analysis** – When we examine and analyse two variables simultaneously and understand the distribution summary stats using descriptive stats and try to gain information its called as bivariate analysis. The purpose is to identify the relationship between 2 variables to understand the patterns that change sthe data corresponding to other data in comparisons.



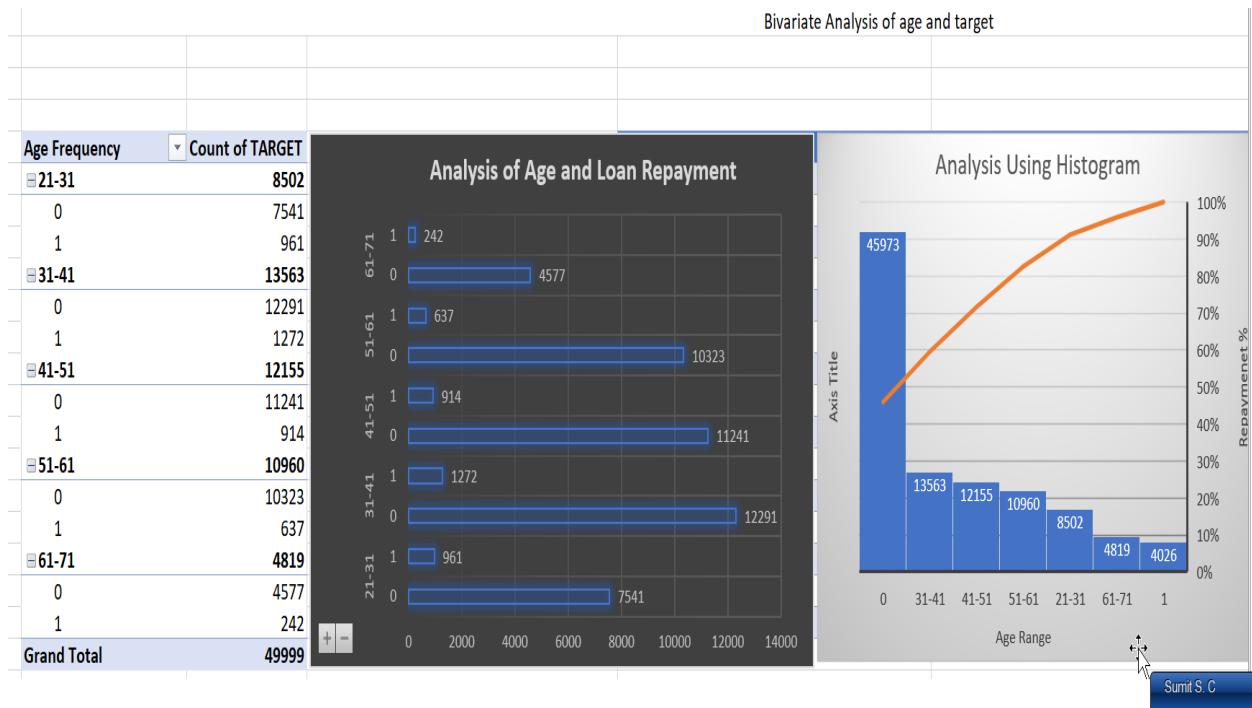
Bivariate 1 Column Bar GraphOn Credit amount & Credit approved



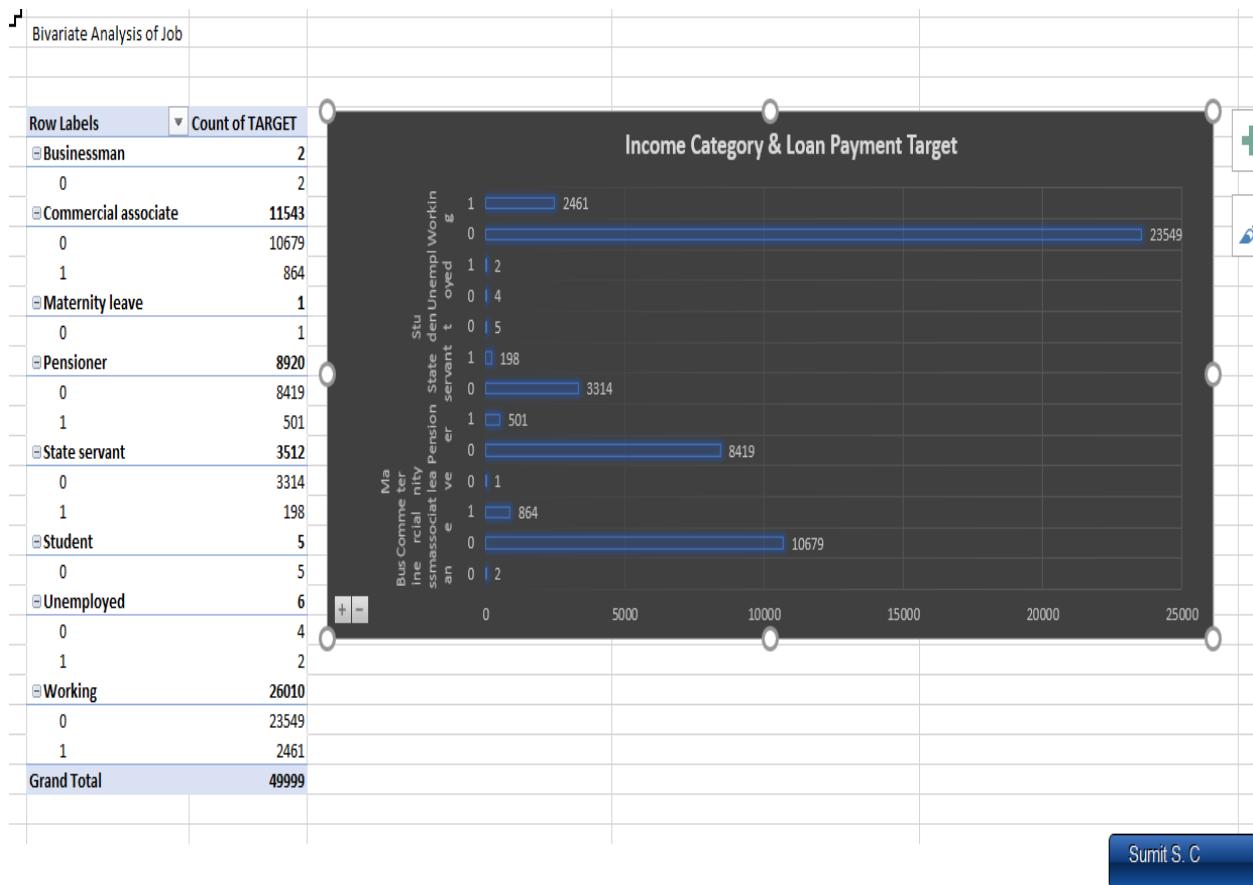
Bivariate 2 Histogram on above set



Bivariate 3 - On education and target counts (Column & Scatterplot distribution)



Bivariate 4 - On age & Target ( Histogram & Bar plot)



#### Bivariate 5 - On Profession & Target type (Bar Plotting)

**Note** – During taking the screenshots I missed the Descriptive stats of Bivariate & Segmented analysis for numerical data which are calculated in the excel sheet.

---

### **Insights:**

1. In the univariate analysis we have selected a single column and checked its descriptive stats and plotted a graph to see its distribution.
2. For univariate analysis we basically take a single variable into consideration and check the other variables that impacts each other using scatter and other plots for reference but all are individuals.
3. I have placed few univariate analyses above for reference.
4. In Segmented univariate analysis we basically divide the single variable and further deeply analysis it says on target type on gender roles or say by loan type and drill the column for more better insights on the column.
5. In segmented analysis we basically get more deeper insights of the category like just gender column gets divided into male & female for analysis , loan type gets into cash or revolving thus giving us insight about loan repayment and categories more.
6. In Bivariate analysis we compare two variables and compare them on a single feature to see the impacts of the comparison on particular analysis.
7. Say in our above analysis we prepared histogram and scatterplots/Boxplots to show variations that arises when we compare different entity and its overall impact (here on loan repayment).
8. For more insights one can go through the excel files if I have missed something in screenshot or in report because these questions is a long one and writing everything will not be feasible.
9. Links
  - Q4 [Excel Link](#), [Gdocs Links](#)
  - Loom Video PPT [Q4 – Variate analysis](#)

## E): Identify Top Correlations for Different Scenarios: Understanding the correlation between variables and the target variable can provide insights into strong indicators of loan default.

- Task:** Segment the dataset based on different scenarios (e.g., clients with payment difficulties and all other cases) and identify the top correlations for each segmented data using Excel functions.

Correlation Among the Numerical Columns Which have impact on Loan Repayment using the data analysis tool in Excel or by using =CORREL() function and passing the arrays to compare

Column	SK_ID_CURR	R	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE	DAYS_BIRTH	DAYS_EMPLOYED	DAYS_REGISTRATION	CNT_FAM_MEMBERS	OBS_30_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	OBS_60_CNT_SOCIAL_CIRCLE	DEF_60_CNT_SOCIAL_CIRCLE
SK_ID_CURR		1															
TARGET		0.0032485	0.0032485														
CNT_CHILDREN		0.0055381	0.0263639	1													
AMT_INCOME_TOTAL		-0.0030144	0.0108937	5	45	0.009588558		1									
AMT_CREDIT		0.0007523	0.0324285	87	47	0.00497156	0.068015897	1									
AMT_ANNUITY		0.0020881	-0.0123982	0.026180456	0.087008438		0.7694887	87	1								
AMT_GOODS_PRICE		0.0007407	0.0412781	72	1	0.000232954	0.068891714	0.9867045	86	0.7743134041	1						
REGION_POPULATION_RELATIVE		0.0013783	0.0407991	12	75	0.025559665	0.029841469	0.0951112	21	0.11511008	0.099196948	1					
DAY_S_BIRTH	(In Years)	0.0012087	0.0767854	64	69	0.592828045	-0.01588528	0.0951321	18	0.007711224	0.057941876	0.052514214	1				
DAY_S_EMPLOYED	(In Years)	0.0043925	0.0424724	44	89	0.241539929	-0.031510652	0.0677398	92	0.108708539	-0.06000589	-0.004138451	0.621739548	1			
DAY_S_REGISTRATION	(In Years)	0.0056884	0.0423275	98	55	0.181202379	-0.009967079	0.0054756	5	-0.03524593	-0.00610759	0.059309167	0.335623699	0.209162575	1		
CNT_FAM_MEMBERS		0.0019498	0.0129954	15	6	0.880455292	0.011225911	0.0659971	55	0.07757959	0.061572677	-0.02503741	-0.277249972	-0.250761119	-0.170107255	1	
OBS_30_CNT_SOCIAL_CIRCLE		0.0030089	0.0145479	2	79	0.016710609	-0.008758226	0.0024712	04	0.008504737	0.002115757	-0.01890262	-0.11515492	0.004888783	-0.010298782	0.025872176	1
DEF_30_CNT_SOCIAL_CIRCLE		0.0089507	0.0417682	56	52	0.028930585	-0.007701787	0.0157185	86	0.021551819	-0.01712151	0.008747605	-0.018090909	0.015891615	-0.004979494	-0.02656908	0.311824664
OBS_60_CNT_SOCIAL_CIRCLE		0.0033452	0.0143126	87	1	0.016594434	-0.00872687	0.0030271	72	0.008194734	0.002403725	-0.017789196	-0.113993506	0.004820803	-0.0105775	0.025905877	0.098332865
DEF_60_CNT_SOCIAL_CIRCLE		0.0074759	0.0455955	94	31	0.003891976	-0.007405962	0.0207732	95	-0.024862556	0.003560435	-0.002781762	0.014769193	-0.00688787	-0.004480448	0.235126904	0.856262392

Sumit S. C

Correlation of factors 1 [HD Image](#) (Considering ID)

Column	Target	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE	DAY_S_BIRTH	DAY_S_EMPLOYED	DAY_S_REGISTRATION	CNT_FAM_MEMBERS	OBS_30_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	OBS_60_CNT_SOCIAL_CIRCLE	DEF_60_CNT_SOCIAL_CIRCLE	
TARGET		0.026														
CNT_CHILDREN		0.00393	1													
AMT_INCOME_TOTAL		0.010	0.00995													
AMT_CREDIT		0.00374	0.00095	5	88168	1										
AMT_ANNUITY		0.032	0.0049	0.0090158	7	87	1									
AMT_GOODS_PRICE		42034	0.0049	0.0090158	7	87	1									
REGION_POPULATION_RELATIVE		0.00212	0.00521	0.0080884	0.0080884	0.0080884	0.0080884	0.0080884	0.0080884	0.0080884	0.0080884	0.0080884	0.0080884	0.0080884	0.0080884	
DAY_S_BIRTH	(In Years)	0.076	0.0392	0.0159885	33111	0.0077	0.0257605	0.03112	0.03112	0.03112	0.03112	0.03112	0.03112	0.03112	0.03112	
DAY_S_EMPLOYED	(In Years)	0.042	0.01415	0.0011206	73866	0.2087	0.0060005	0.03127	0.03127	0.03127	0.03127	0.03127	0.03127	0.03127	0.03127	
DAY_S_REGISTRATION	(In Years)	0.042	0.01413	0.0009670	0.0009670	0.0009670	0.0009670	0.0009670	0.0009670	0.0009670	0.0009670	0.0009670	0.0009670	0.0009670	0.0009670	
CNT_FAM_MEMBERS		0.012	0.00404	0.0112255	99715	0.0773	0.005172	0.02772	0.02772	0.02772	0.02772	0.02772	0.02772	0.02772	0.02772	
OBS_30_CNT_SOCIAL_CIRCLE		0.00531	0.00045	0.0020844	11322	0.1151	0.099196	0.002115	0.002115	0.002115	0.002115	0.002115	0.002115	0.002115	0.002115	
DEF_30_CNT_SOCIAL_CIRCLE		0.0074759	0.0455955	94	31	0.003891976	-0.007405962	0.0207732	95	-0.024862556	0.003560435	-0.002781762	0.014769193	-0.00688787	-0.004480448	0.235126904
OBS_60_CNT_SOCIAL_CIRCLE		0.00531	0.00045	0.0020844	11322	0.1151	0.099196	0.002115	0.002115	0.002115	0.002115	0.002115	0.002115	0.002115	0.002115	
DEF_60_CNT_SOCIAL_CIRCLE		0.0074759	0.0455955	94	31	0.003891976	-0.007405962	0.0207732	95	-0.024862556	0.003560435	-0.002781762	0.014769193	-0.00688787	-0.004480448	0.235126904

Sumit S. C

Correlation of factors 2 [HD Image](#) (Removing ID)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q
1	Correlation Among the Numerical Columns Which have impact on Loan Repayment using the data analysis tool in Excel or by using =CORREL() function and passing the arrays to compare																
2	Column	SK_ID_CURR	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	GENDER	POPULATION_RELATIONSHIP	DAYS_BIRTH	YEARS_EMPLOYED	YEARS_REGISTRATION	YNT_FAM_MEMBERS	DEF_30_CNT_SOCIAL_CIRCLE	DEF_60_CNT_SOCIAL_CIRCLE	DEF_90_CNT_SOCIAL_CIRCLE
3	SK_ID_CURR																
4	TARGET	0.00234977															
5	CNT_CHILDREN	0.00533179	0.02636333														
6	AMT_INCOME_TOTAL	-0.00301443	0.00083374	0.00650055													
7	AMT_CREDIT	-0.00722387	-0.0324283	0.0049756	0.03316867												
8	AMT_ANNUITY	-0.00288282	-0.0423892	0.02690456	0.00200458	0.079480767											
9	AMT_GOODS_PRICE	-0.0074072	-0.0412761	0.00023254	0.00889174	0.588704368	0.77450404										
10	REGION_POPULATION_RELATIONSHIP	0.00378512	-0.0470392	-0.025555665	0.02304468	0.0951122	0.1910008	0.03905946									
11	DAYS_BIRTH	0.0030784	-0.0767824	-0.32202495	-0.05368529	0.05933218	-0.00771124	0.05768074	0.0251424								
12	YEARS_EMPLOYED	-0.0043254	-0.042376	-0.24503932	-0.01510652	-0.00710589	-0.008010539	-0.06500689	-0.00458451	0.02175598							
13	YEARS_REGISTRATION	0.00389498	-0.042376	-0.18202379	-0.003967079	-0.00347365	-0.03324532	-0.0060739	0.05300967	0.033823699	0.20362575						
14	YNT_FAM_MEMBERS	0.0094801	0.0239346	0.00045326	0.0122511	0.03397165	0.07737193	0.08572377	-0.02303741	-0.272249972	-0.23076118	-0.1700725					
15	DEF_30_CNT_SOCIAL_CIRCLE	-0.0030082	0.0164738	0.01710083	-0.008758226	0.002872204	-0.008594737	0.00215751	-0.01880282	-0.01151542	0.004688783	-0.0028702	0.02507276				
16	DEF_60_CNT_SOCIAL_CIRCLE	-0.00830786	0.0478623	-0.002893085	-0.00770787	-0.01571858	-0.02051019	-0.071203	0.008747603	-0.00630369	0.0589168	-0.00497364	-0.002856308	0.31624684			
17	DEF_90_CNT_SOCIAL_CIRCLE	-0.002345207	0.0143261	0.01639493	-0.00872897	0.003027172	-0.008194724	0.002403725	-0.0773936	-0.01339306	0.00462003	-0.016575	0.02505877	0.9933283	0.31423286		
18	DEF_60_CNT_SOCIAL_CIRCLE	-0.007473994	0.04438533	-0.003891976	-0.007405982	-0.02072355	-0.0248626	-0.0214233	0.00360435	-0.00279782	0.04753193	-0.00688787	-0.004490448	0.23612504	0.056282382	0.23607483	

Page 3

Page 4

Page 5

Page 6

Sumit S. C

### Insights:

- We have used the Correlation analysis present in data > data analysis tab which helps us to choose the correlation factor for the selected column range or called array of numbers. Correlations are basically done by taking the 1<sup>st</sup> column into account and then identifying the changes in other columns with respect to 1<sup>st</sup> one. For our analysis we have taking target type for the correlation if we need to do more analysis we can select the specific columns such that we can make the correlation between them.
- In correlation comparison we do the comparison of each column with oneself and with other column in question and thus a matrix is formed of the same length as the column numbers and the diagonal is equal to 1 (since relation of same column comparison.)
- The same result can be achieved with making column and then applying the **CORREL** function on the columns to compare one by one which will be time consuming process. A histogram or heatmap like plot is plotted to get the relevance of data. In Excel

---

the lower triangle gets filled while upper remains blank as it's the transposed values for columns if we make it in python or BI tools they might fill the same on other sides.

4. Do note that for correlation we need numbers so we need to remove any categorical value columns first. So, I cleaned up the dataset further by removing the categorical column and the numerical column which does not directly impact on the loan behaviors.
5. From the above analysis we have found out about the top 15 factors that might impact the loan defaults as these factors play a role in individuals spending and savings habit etc.
6. Links
  - Q5 [Excel File](#) , [Gdocs Files](#)
  - Loom Video PPT [Q5 – Find the Correlation of scenario](#)

---

## **Important Links :**

[Drive Folder Link](#)

[Individual Excel Sheets](#) (Individual Questions File)

[Final Excel sheet](#) (cleaned dataset & all questions also present on drive folder page)

Word File Link & Pdf File (Will Be in the [drive folder](#) – can't add before I upload the file)

[Video Presentation\(Loom Folder 6 Videos\)](#)

[Zip Link](#) (Download these to extract the pdf and excel sheet files)

---

# **Thank You**

---