

Bank Loan Case Study

Link For Excel Sheet:

<https://docs.google.com/spreadsheets/d/1mlateqxBvMcygMzOsjjxkTRCFzyoXbFu/edit?usp=sharing&ouid=107365393175079460343&rtpof=true&sd=true>

Excel file contains different worksheets which have results of different tasks in them.

Project Description:

This project aims to analyze a dataset containing information about various bank loan applications. The goal is to gain insights about approval of bank loans, such as the relation between income and credit. The data provided has various missing or null Data, our task is to handle those missing values appropriately, by either deleting or imputing these data. There are various outliers in data, we have to find these outliers. We also have to check for data imbalance and perform various analyses on data, such as univariate and bivariate analysis. Finding correlation between various parameters would help us understand what factors affect most in bank loan application approval. Thus, by employing statistics and Excel formulas, we will extract meaningful conclusions to help understand the factors that contribute to a bank loan getting approved.

Approach:

As an individual working on this project, I followed a structured approach to analyze data about bank loan applications. I began by carefully examining the provided database and familiarizing myself with its structure and columns. I tried to find columns which had the most significance in the dataset. I handled missing values by eliminating columns which had most empty cells, and were not significant. And imputed data into cells that were necessary for analysis. Then, I utilized Excel fundamentals to retrieve the necessary information for each task, employing appropriate functions and statistical methods. I focused on data accuracy and quality throughout the project, ensuring reliable results. By leveraging my Excel skills and maintaining a systematic workflow, I successfully executed the project and created a comprehensive report that fulfilled the objectives of providing marketing insights and investor metrics.

Tech-Stack Used:

For this project, I utilized Microsoft Excel as the primary software tool.

Insights:

Task 1:

Identify Missing Data and Deal with it Appropriately (Data Cleaning):

To find data having missing values we utilized COUNTA function in Excel, which returns no. of cells which are not blank.

Formula:

=COUNTA(A4:A50002)

This gave us the number of rows in the TARGET column, which is the total number of rows which we have to consider for analysis.

C17					
	A	B	C	D	E
1	49999	49999	49999	49999	49999
2	0	0	0	0	0
3	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR
4	100002	1	Cash loans	M	N
5	100003	0	Cash loans	F	N
6	100004	0	Revolving loans	M	Y
7	100006	0	Cash loans	F	N
8	100007	0	Cash loans	M	N
9	100008	0	Cash loans	M	N
10	100009	0	Cash loans	F	Y

The columns which had missing data in them were found out by using the formula:

=(100-(V1/\$A1)*100)

This formula gives us the percentage of missing values in the column.

Alignment	Number	Styles	Cells
AO	AP	AQ	AR
49999	21827	49873	40055
0	56.3451269	0.25200504	19.88839777
50.77101542	58.39916798	48.78897578	66.47932959
ORGANIZATION_TYPE	EXT_SOURCE_1	EXT_SOURCE_2	EXT_SOURCE_3
business Entity Type 3	0.083036967	0.262948593	0.13937578
chool	0.311267311	0.622245775	0.0959
overnment	0.555912083	0.729566691	
usiness Entity Type 3	0.65044169		
eligion	0.322738287		
ther	0.354224732	0.621226338	
usiness Entity Type 3	0.774761413	0.723999852	0.492060094
ther	0.714279286	0.54065445	
NA	0.587334047	0.205747288	0.751723715
lectricity	0.746643629		
edicine	0.319760172	0.651862333	0.363945239
NA	0.72204445	0.555183162	0.652896552
usiness Entity Type 2	0.464831117	0.715041819	0.176652579
elf-employed	0.566906613	0.77008707	0.1474
ransport: type 2	0.721939769	0.642656205	0.3495
usiness Entity Type 2	0.115634337	0.346633981	0.678567689
overnment	0.23637784	0.062103038	
onstruction	0.683513346		
ousing	0.706428403	0.556727426	0.0278
indergarten	0.58661714	0.477649155	0.0617
elf-employed	0.565654882	0.113374513	0.0722
rade: type 7	0.43770902	0.233766958	0.542445144
elf-employed	0.457142972	0.358951229	0.0907
NA	0.624304737	0.669056695	0.1443

We highlighted columns with missing values by using conditional formatting.

We found no. of columns with missing values greater than 10% by this formula:

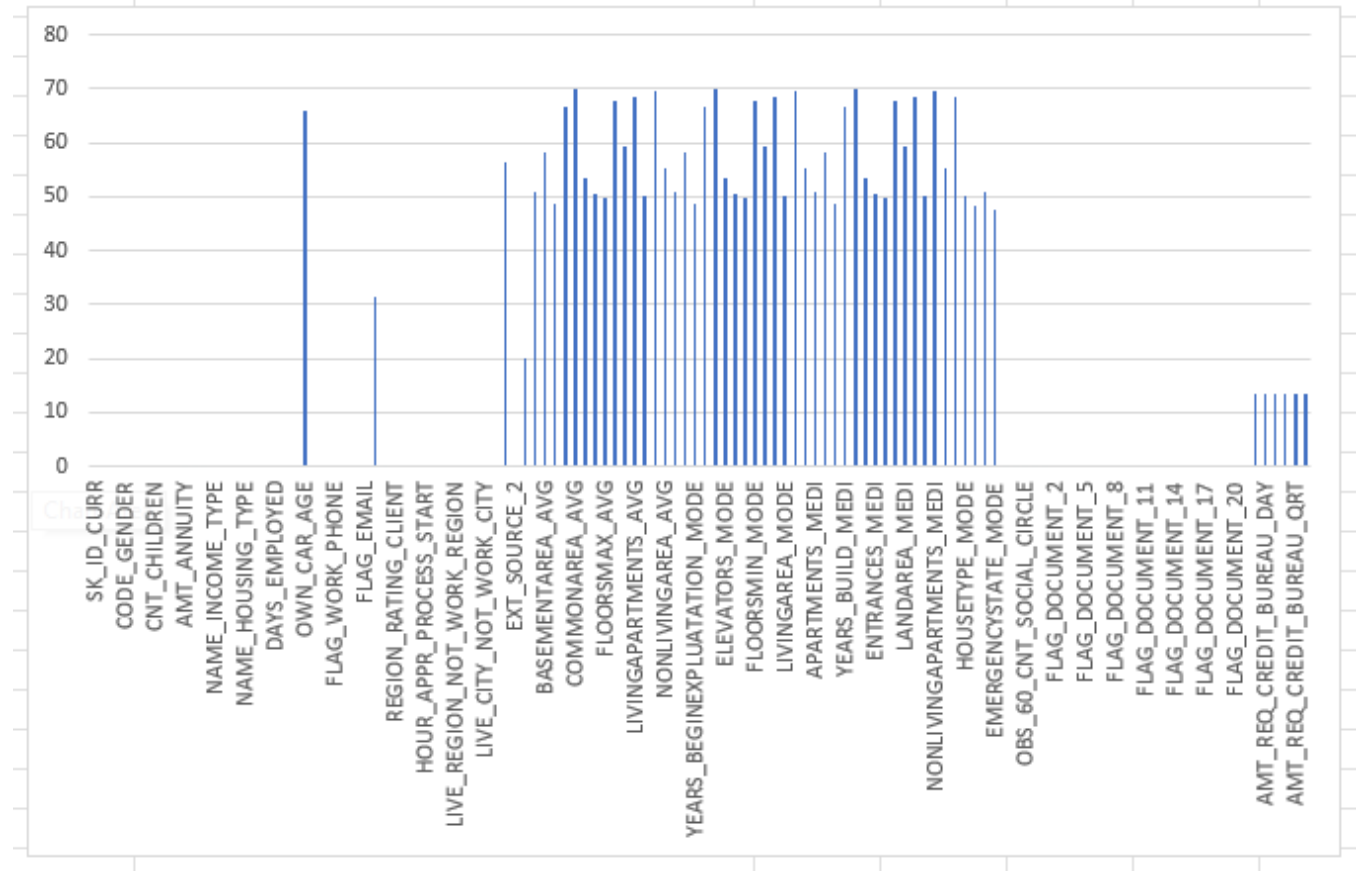
=COUNTIF(A2:DR2,">10")
 =COUNTIF(A2:DR2,"<10")

No. of columns with missing data more than 10%	57
No. of columns with missing data less than 10%	65

We plotted a bar graph to better understand the number of columns containing missing values



To Visualize columns and their respective missing values.



We saved all 65 columns with least missing values in a new worksheet called Cleaned Data.

Task 2:

Identify Outliers in the Dataset:

To find Outliers in the Dataset we utilized functions like QUARTILE, IQR, and conditional formatting to identify potential outliers.

We first copied the columns of interest into a new worksheet for finding Outliers.

Columns Copied are:

SK_ID_CURR	TARGET	AMT_INCOME_TOTAL	CNT_CHILDREN	DAYS_EMPLOYED	DAYS_EMPLOYED(ABS)
100002	1	202500	0	-637	637
100003	0	270000	0	-1188	1188

We used the QUARTILE function to find quartile 1 and quartile 3, along with the IQR and upper limit and lower limit ranges.

Formulae:

=QUARTILE.EXC(Table5[AMT_INCOME_TOTAL],1)

=QUARTILE.EXC(Table5[[#All],[AMT_INCOME_TOTAL]],3)

=I4-I2 (IQR)

=I4+1.5*I6 (Upper limit)

=I2-1.5*I6 (Lower limit)

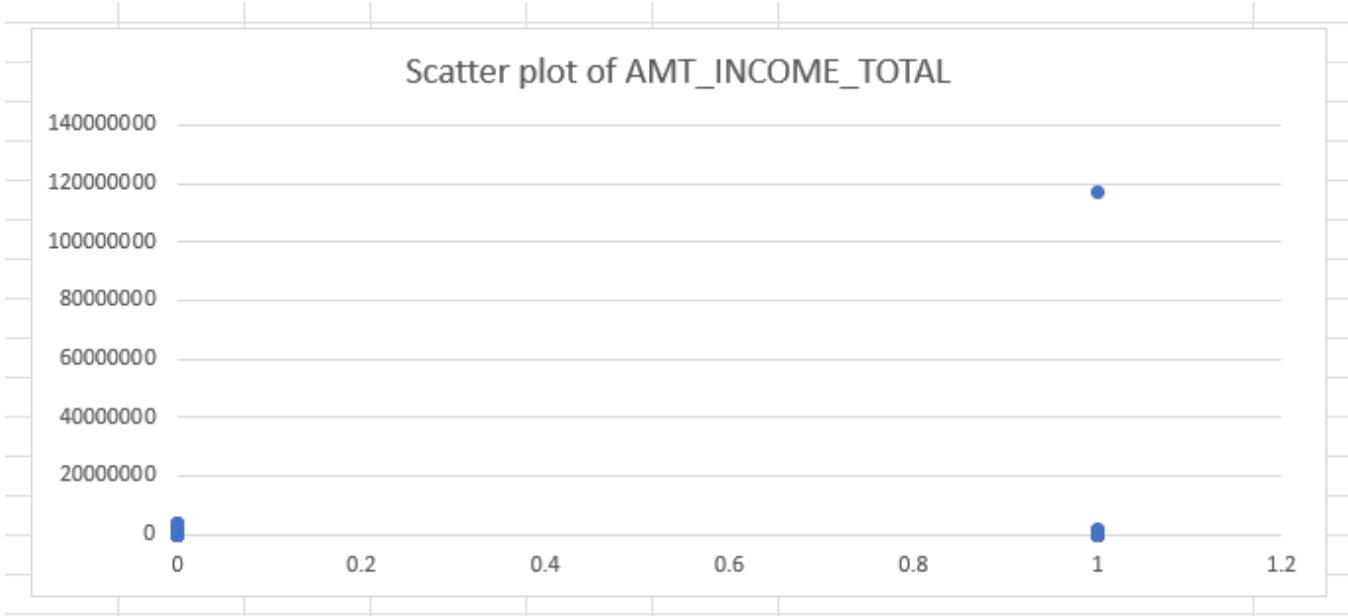
=COUNTIF(C2:C50000,">337500") (Count of elements outside limits)

Outliers in AMT_INCOME_TOTAL					
Quartile 1	112500		Upper Limit	337500	
Quartile 3	202500				Count of elements above upper limit
			Lower Limit	-22500	2295
IQR	90000				

We used Conditional Formatting to highlight the cells which contain values outside the limits.

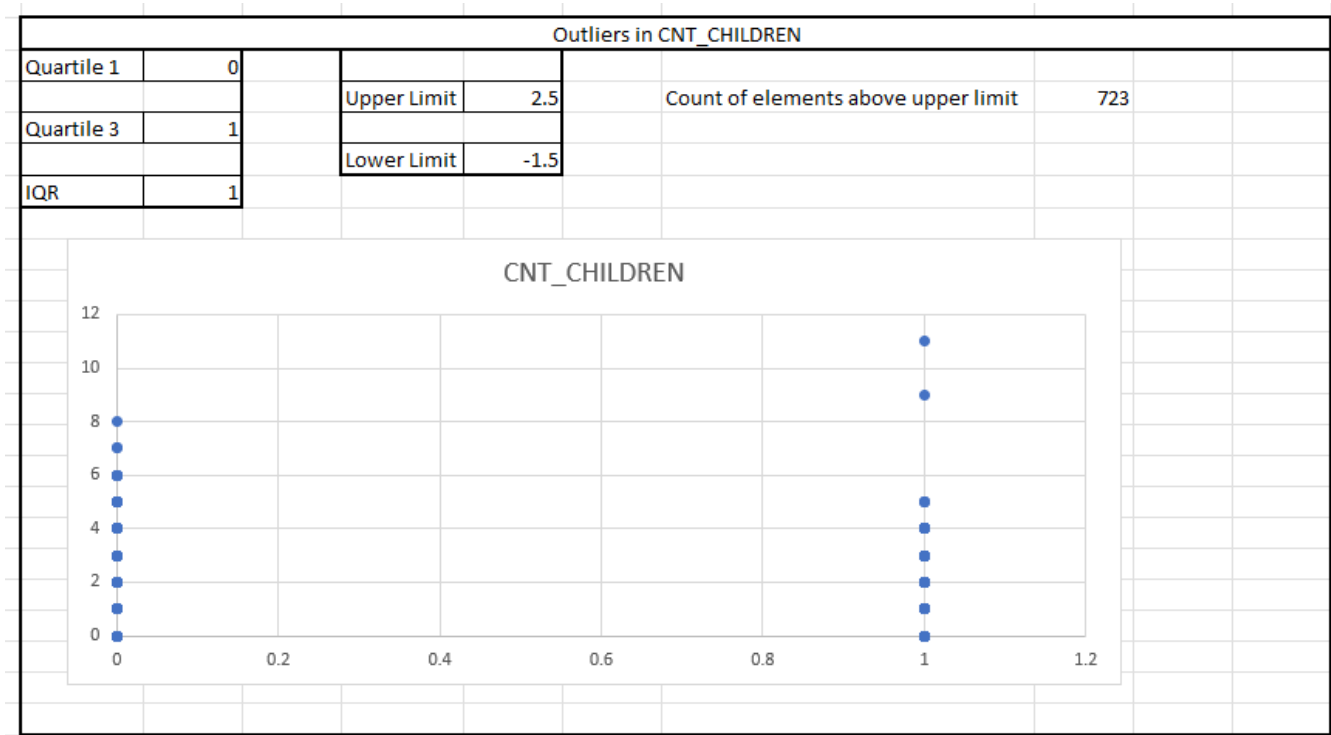
100007	0	121500
100008	0	99000
100009	0	171000
100010	0	360000
100011	0	112500
100012	0	135000
100014	0	112500
100015	0	33410.155

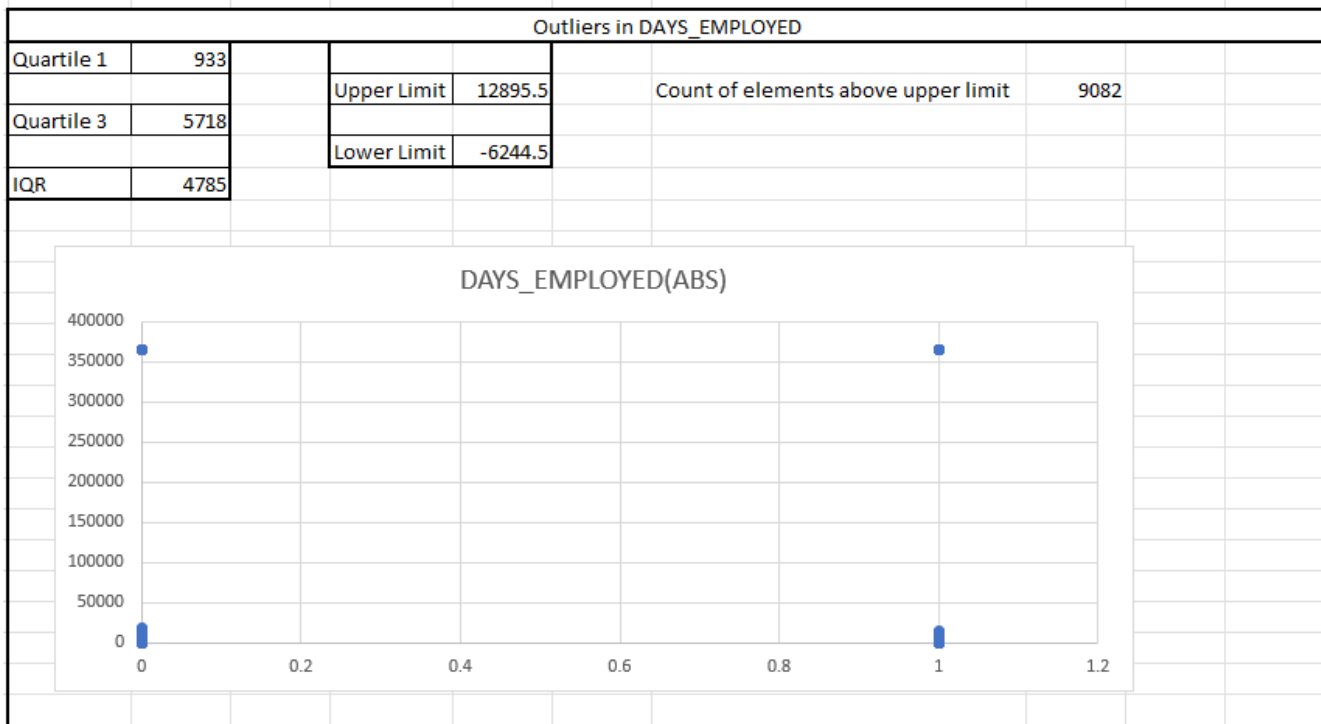
We also plotted A scatter plot to visualize the outliers



In the above plot the point which lies outside the general trend, and is very much out of the scope can be called an outlier.

Similar Steps were done for other columns and following results were obtained.





	A	B	C	D	E	F
1	SK_ID_CURR	TARGET	AMT_INCOME_TOTAL	CNT_CHILDREN	DAYS_EMPLOYED	DAYS_EMPLOYED(ABS)
2	100002	1	202500	0	-637	637
3	100003	0	270000	0	-1188	1188
4	100004	0	67500	0	-225	225
5	100006	0	135000	0	-3039	3039
6	100007	0	121500	0	-3038	3038
7	100008	0	99000	0	-1588	1588
8	100009	0	171000	1	-3130	3130
9	100010	0	360000	0	-449	449
10	100011	0	112500	0	365243	365243
11	100012	0	135000	0	-2019	2019
12	100014	0	112500	1	-679	679
13	100015	0	38419.155	0	365243	365243
14	100016	0	67500	0	-2717	2717
15	100017	0	225000	1	-3028	3028
16	100018	0	189000	0	-203	203
17	100019	0	157500	0	-1157	1157
18	100020	0	108000	0	-1317	1317
19	100021	0	81000	1	-191	191
20	100022	0	112500	0	-7804	7804
21	100023	0	90000	1	-2038	2038
22	100024	0	135000	0	-4286	4286
23	100025	0	202500	1	-1652	1652
24	100026	0	450000	1	-4306	4306
25	100027	0	83250	0	365243	365243
26	100029	0	135000	2	-746	746
27	100030	0	90000	0	-3494	3494
28	100031	1	112500	0	-2628	2628
29	100032	0	112500	1	-1234	1234
30	100033	0	270000	0	-1796	1796
31	100034	0	90000	0	-1010	1010
32	100035	0	292500	0	-2668	2668
33	100036	0	112500	0	-1104	1104
34	100037	0	90000	0	-4404	4404
35	100039	0	360000	1	-2060	2060
36	100040	0	135000	0	-4585	4585
37	100041	0	112500	0	-1275	1275
38	100043	0	198000	2	-768	768
39	100044	0	121500	0	-1288	1288

In the above image all highlighted cells are outliers.

Task 3:

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.

To find Data Imbalance we find the number of each element in the TARGET column. For doing this we use the COUNTIF formula.

Formulae:

=COUNTIF(B:B,0)

=COUNTIF(B:B,1)

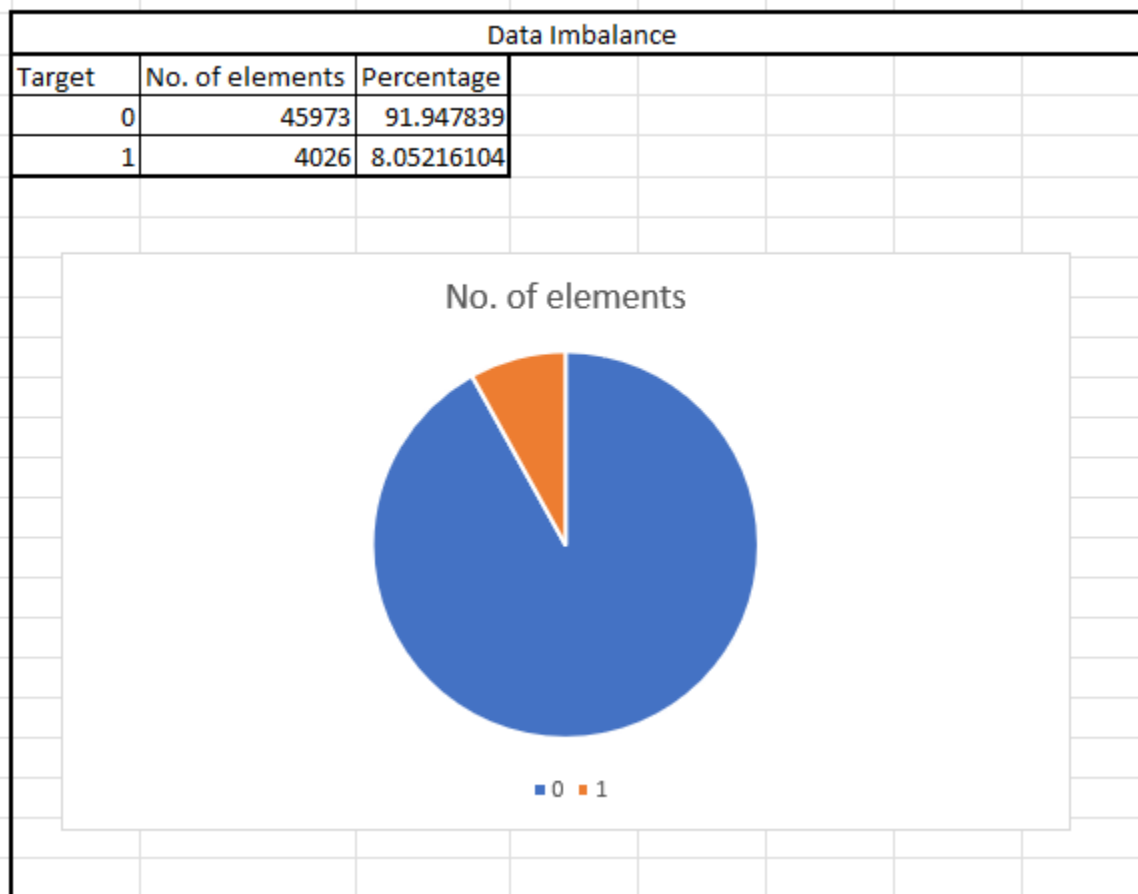
=F4/49999*100

(Percentage)

=F5/49999*100

(Percentage)

We also plot this Data in Pie Chart to visualize the Data Imbalance.



As we can see the no. of 0 in TARGET is very large compared to no. of 1. This will result in a very large data imbalance. Which might skew the results and give less accurate results.

Task 4:

Perform Univariate, Segmented Univariate, and Bivariate Analysis:

To perform Univariate/ Segmented Univariate analysis, we have to utilize functions such as COUNT, AVERAGE, or MEDIAN to find out the total number of applicants over a particular range or how much credit one shall receive according to their income, and other such relations.

We start by selecting two columns, Credit and Income, we have selected this columns as they have higher correlation. We find maximum and minimum values of these columns excluding outliers.

Maximum Income	117000000	Maximum Credit	4050000
Excluding Outlier	3825000		
Minimum Income	25650	Minimum Credit	45000

This Data helps us to define ranges to find how many applicants fall in each range.

We Define Ranges on particular Intervals

Income Ranges	Credit Ranges
25000-50000	0-200000
50000-75000	200000-400000
75000-100000	400000-600000
100000-125000	600000-800000
125000-150000	800000-1000000
150000-175000	1000000-1200000
175000-200000	1200000-1400000
200000-225000	1400000-1600000
225000-250000	1600000-1800000
250000-275000	1800000-2000000
275000-300000	2000000-2200000
300000-325000	2200000-2400000
325000-350000	2400000-2600000
350000-375000	2600000-2800000
375000-400000	2800000-3000000
400000-425000	3000000-3200000
425000-450000	3200000-3400000
450000-475000	3400000-3600000
475000-500000	3600000-3800000
500000+	3800000+

We find number of applicants over these ranges by utilizing functions such as:

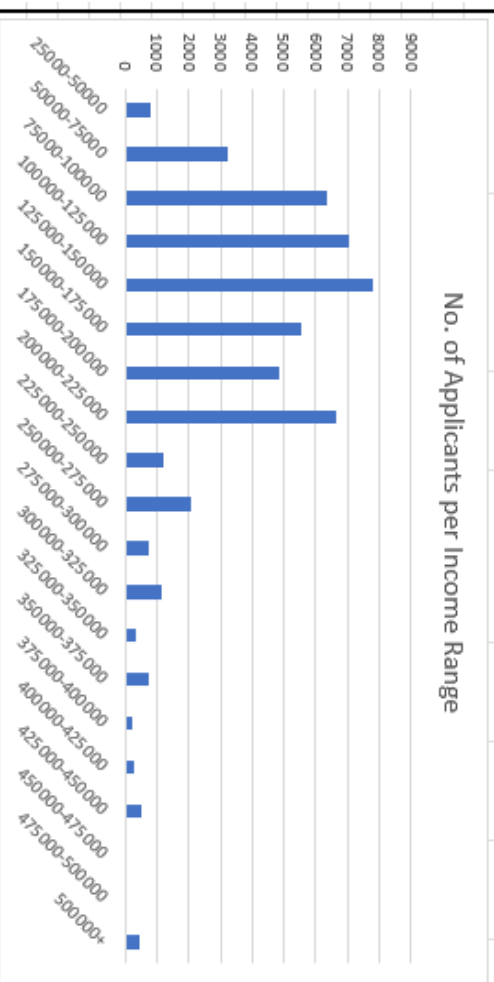
=FREQUENCY(C:C,X5:X23)

=FREQUENCY(D:D,Y5:Y23)

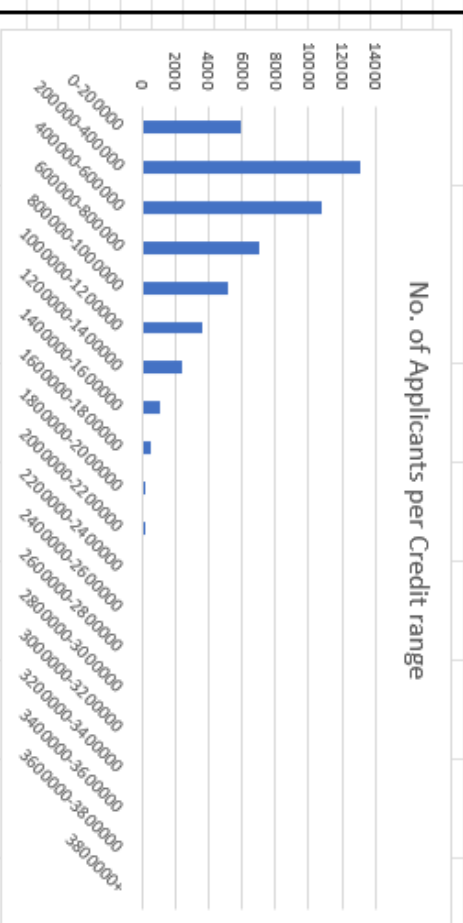
We also plot bar charts to visualize the frequency of applicants in each range.

Univariate Analysis

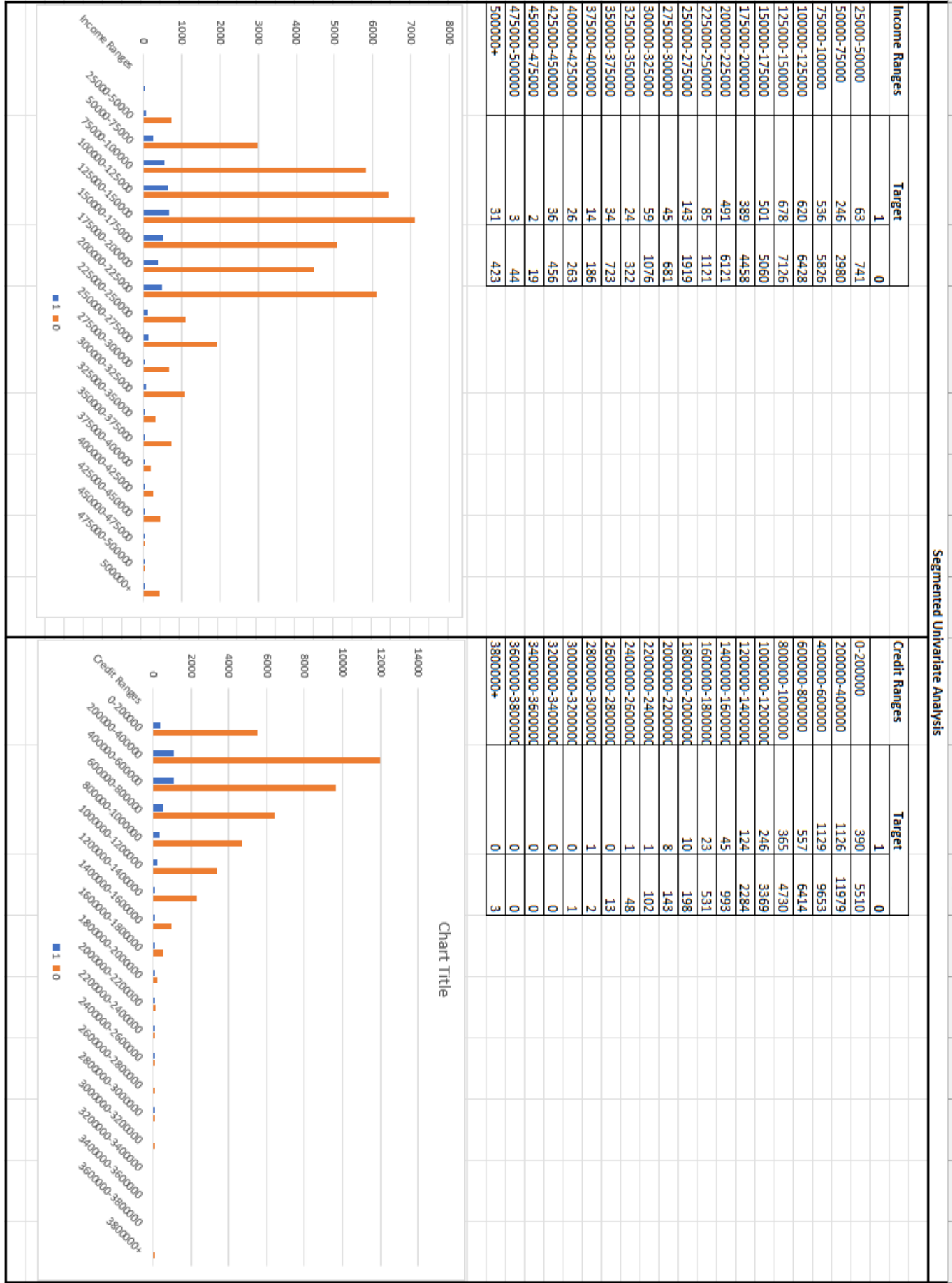
Income Ranges	No. of Applicants
25000-50000	804
50000-75000	3226
75000-100000	6362
100000-125000	7048
125000-150000	7804
150000-175000	5561
175000-200000	4847
200000-225000	6612
225000-250000	1206
250000-275000	2062
275000-300000	726
300000-325000	1135
325000-350000	346
350000-375000	757
375000-400000	200
400000-425000	289
425000-450000	492
450000-475000	21
475000-500000	47
500000+	454



Credit Ranges	No. of Applicants
0-200000	5900
200000-400000	13105
400000-600000	10782
600000-800000	6971
800000-1000000	5095
1000000-1200000	3615
1200000-1400000	2408
1400000-1600000	1038
1600000-1800000	554
1800000-2000000	208
2000000-2200000	151
2200000-2400000	103
2400000-2600000	49
2600000-2800000	13
2800000-3000000	3
3000000-3200000	1
3200000-3400000	0
3400000-3600000	0
3600000-3800000	0
3800000+	3



Similarly for segmented univariate analysis we split the Data into two classes according to TARGET.

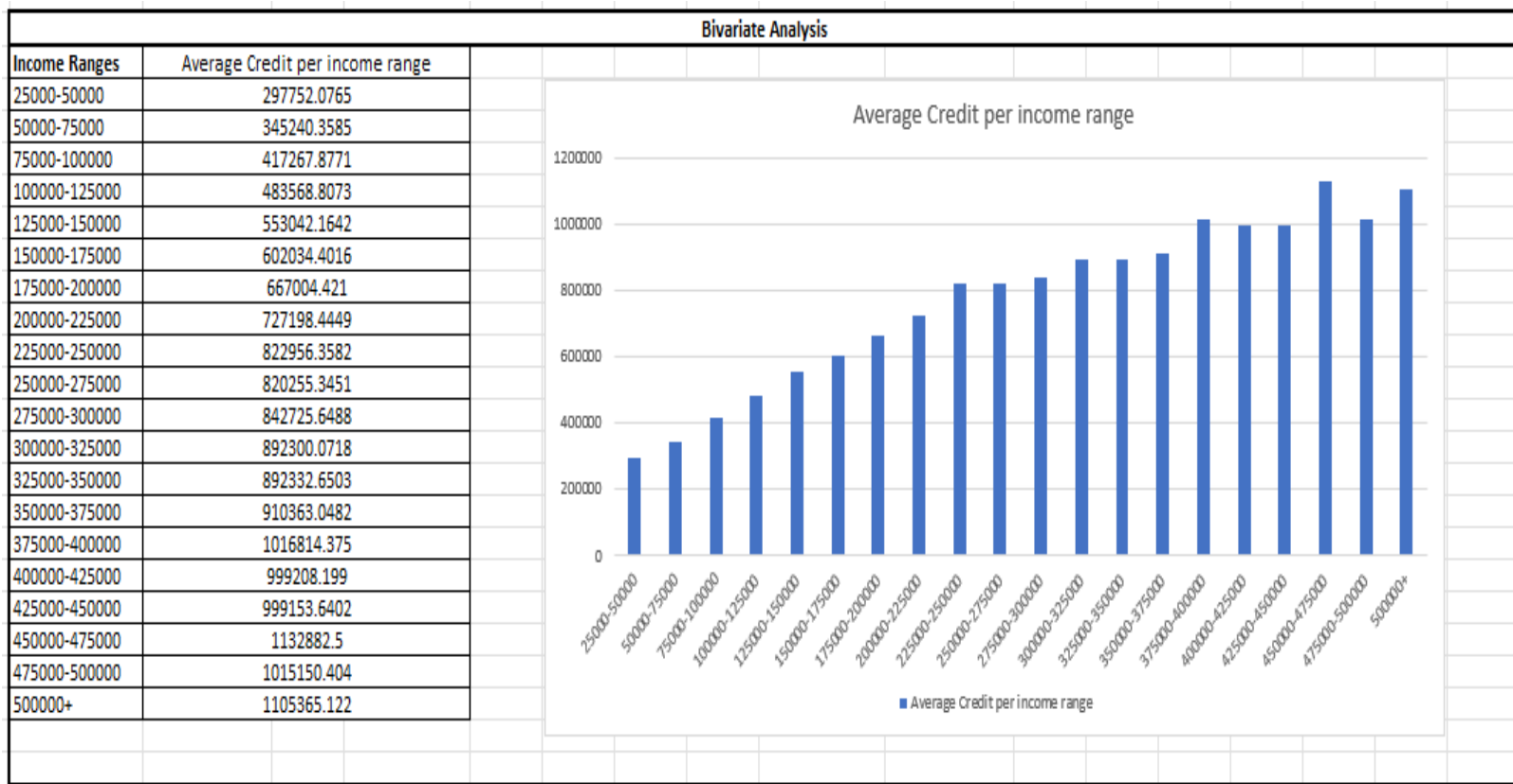


To perform Bivariate Analysis, we need to find the average of credit per income range, for that we use the AVERAGEIF function.

=AVERAGEIFS(\$D\$2:\$D\$50000,\$C\$2:\$C\$50000,">"&X4,\$C\$2:\$C\$50000,"<="&X5)

Above formula checks two conditions, if element is greater than lower limit and smaller than upper limit, and only then is considered for average.

We plotted a Bar Graph Similar to above analysis



Task 5:

Identify Top Correlations for Different Scenarios:

To find Correlation of different columns we utilized the CORREL function of Excel. We first separated the data into three tables, one having only 0 Target, one having 1 Target and both combined Target. We found correlation tables for all these by using the CORREL function, and made it better for visualization using conditional formatting.

Formulae:

=CORREL(\$C:\$C,B:B)

=CORREL('Target 0 Data'!C:C,'Target 0 Data'!\$F:\$F)

=CORREL('Target 1 Data'!B:B,'Target 1 Data'!\$F:\$F)

For better visualization heatmaps of correlation matrix were created.

Correlation for Target 1												
TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	REGION_POPULATION_RELATIVE	DAYS_BIRTH	DAYS_EMPLOYED	DAYS_REGISTRATION	DAYS_ID_PUBLISH	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL
1	0.026363931	0.010893745	-0.03242835	-0.01239904	-0.040799172	0.07678765	-0.040294905	0.04242679	0.046926745	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL
CNT_CHILDREN	1	0.010110177	0.007601905	0.029172977	-0.020395154	0.2496732	-0.189324194	0.15213117	-0.042360717	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT
AMT_INCOME_TOTAL	0.01089	1	0.015271444	0.018004594	-0.006180303	0.099333662	-0.01555963	-0.09360152	-0.00912206	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY
AMT_CREDIT	-0.0324	0.007601905	1	0.749665201	0.067775624	-0.14250603	0.01603971	-0.0424404	-0.04377901	AMT_CREDIT	AMT_ANNUITY	REGION_POPULATION_RELATIVE
AMT_ANNUITY	-0.0124	0.029172977	0.018004594	1	0.073123998	-0.00875171	-0.07356008	0.021581654	-0.0132109	AMT_ANNUITY	REGION_POPULATION_RELATIVE	DAYS_BIRTH
REGION_POPULATION_RELATIVE	-0.0408	-0.020395154	-0.006180303	0.067775624	1	-0.01546873	0.007742909	-0.046130288	-0.00518563	REGION_POPULATION_RELATIVE	DAYS_BIRTH	DAYS_EMPLOYED
DAYS_BIRTH	0.07679	0.2496732	0.009033662	-0.008751713	-0.01546873	1	-0.58147941	0.288437837	0.247896571	DAYS_BIRTH	DAYS_EMPLOYED	DAYS_REGISTRATION
DAYS_EMPLOYED	-0.0403	0.15213117	-0.009561152	-0.07356008	0.021581654	-0.58147941	1	-0.188719457	-0.188719457	DAYS_EMPLOYED	DAYS_REGISTRATION	DAYS_ID_PUBLISH
DAYS_REGISTRATION	0.04234	0.009033662	0.01603971	-0.07356008	0.021581654	0.007742909	-0.58147941	1	0.09029149	DAYS_REGISTRATION	DAYS_ID_PUBLISH	DAYS_ID_PUBLISH
DAYS_ID_PUBLISH	0.04693	-0.042360717	-0.00912206	-0.02132109	-0.00518563	0.247896571	-0.23063668	0.09029149	1	DAYS_ID_PUBLISH	DAYS_ID_PUBLISH	DAYS_ID_PUBLISH

Correlation for Target 0												
	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	REGION_POPULATION_RELATIVE	DAYS_BIRTH	DAYS_EMPLOYED	DAYS_REGISTRATION	DAYS_ID_PUBLISH	TARGET	
TARGET	1	0.026363931	0.010893745	-0.03242835	-0.01239904	-0.040799172	0.07678765	-0.040294905	0.04242679	0.046926745	TARGET	
CNT_CHILDREN	0.02636	1	0.008319722	0.005705458	0.02638217	-0.02912809	0.33876269	-0.243591518	0.18372478	-0.032537221	CNT_CHILDREN	AMT_INCOME_TOTAL
AMT_INCOME_TOTAL	0.01089	0.008319722	1	0.377965752	0.45113596	0.183941261	0.073709425	-0.162702675	0.06939375	0.03236556	AMT_INCOME_TOTAL	AMT_CREDIT
AMT_CREDIT	-0.0324	0.005705458	0.377965752	1	0.770772965	0.09359444	-0.05108418	-0.07367219	0.00805158	-0.008290189	AMT_CREDIT	AMT_ANNUITY
AMT_ANNUITY	-0.0124	0.02638217	0.45113596	0.770772965	1	0.117280752	0.009915685	-0.113007146	0.034609089	0.009426196	AMT_ANNUITY	REGION_POPULATION_RELATIVE
REGION_POPULATION_RELATIVE	-0.0408	-0.02912809	0.183941261	0.09359444	0.117280752	1	-0.03045442	-0.006510653	-0.05601361	-0.002236288	REGION_POPULATION_RELATIVE	DAYS_BIRTH
DAYS_BIRTH	0.07679	0.33876269	0.073709425	-0.05108418	0.009915685	-0.03045442	1	-0.615289978	0.35028046	0.270073313	DAYS_BIRTH	DAYS_EMPLOYED
DAYS_EMPLOYED	-0.0403	-0.243591518	-0.162702675	-0.07736712	-0.113007146	-0.006510653	-0.615289978	1	-0.204370881	-0.27224239	DAYS_EMPLOYED	DAYS_REGISTRATION
DAYS_REGISTRATION	0.04234	0.18397478	0.06939375	0.008051758	0.034609089	-0.05601361	0.35028046	-0.204370881	1	0.103548902	DAYS_REGISTRATION	DAYS_ID_PUBLISH
DAYS_ID_PUBLISH	0.04693	-0.032537221	0.03238656	-0.00829019	0.009426196	-0.002236288	0.270073313	-0.27224239	0.103548902	1	DAYS_ID_PUBLISH	

Correlation for All Targets												
	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	REGION_POPULATION_RELATIVE	DAYS_BIRTH	DAYS_EMPLOYED	DAYS_REGISTRATION	DAYS_ID_PUBLISH		
TARGET	1	0.026363931	0.010893745	-0.03242835	-0.01239904	-0.040799172	0.07678765	-0.040294905	0.04242679	0.046926745	TARGET	
CNT_CHILDREN	0.02636	1	0.009588558	0.00497156	0.026178823	-0.02555565	0.329268754	-0.238993041	0.181217383	-0.032115773	CNT_CHILDREN	
AMT_INCOME_TOTAL	0.01089	0.009588558	1	0.069315897	0.083080508	0.02984469	0.016002774	-0.031615555	0.009952379	0.003506656	AMT_INCOME_TOTAL	
AMT_CREDIT	-0.0324	0.00497156	0.069315897	1	0.758498914	0.09511221	-0.05934266	-0.070471393	0.00348569	-0.01222876	AMT_CREDIT	
AMT_ANNUITY	-0.0124	0.026178823	0.083080508	0.758498914	1	0.11511507	0.00717245	-0.11448038	0.03318936	0.06716544	AMT_ANNUITY	
REGION_POPULATION_RELATIVE	-0.0408	-0.02555565	0.02984469	0.09511221	0.11511507	1	-0.03251375	-0.00410866	-0.05932244	-0.004945136	REGION_POPULATION_RELATIVE	
DAYS_BIRTH	0.07679	0.329268754	0.016002774	-0.05934266	0.00717245	-0.03251375	1	-0.61535972	0.33832509	0.270825141	DAYS_BIRTH	
DAYS_EMPLOYED	-0.0403	-0.238993041	-0.031615555	-0.07047139	-0.11448038	-0.00410866	-0.61535972	1	-0.20480611	-0.27082022	DAYS_EMPLOYED	
DAYS_REGISTRATION	0.04234	0.181217383	0.009952379	0.00348569	0.03318936	-0.05932244	0.33832509	-0.20480611	1	0.104298561	DAYS_REGISTRATION	
DAYS_ID_PUBLISH	0.04693	-0.032115773	0.003506656	-0.01222876	0.06716544	-0.004945136	0.270825141	-0.27082022	0.104298561	1	DAYS_ID_PUBLISH	
	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	REGION_POPULATION_RELATIVE	DAYS_BIRTH	DAYS_EMPLOYED	DAYS_REGISTRATION	DAYS_ID_PUBLISH		

Results:

While working on this project, I have gained a better understanding of Bank Loan Application Process and Analytics and Advanced Excel methodologies. By analyzing Application Data, I was able to provide insights on various aspects such as Cleaning the Data, Outliers in the Data, Data Imbalance, Univariate and Bivariate Analysis, and correlation between various parameters in bank loan application.

This project has helped me enhance my Excel skills, particularly in functions and data visualization to derive meaningful insights. It has also improved my ability to interpret data and provide actionable recommendations based on the analysis.