

Bank Loan Case Study

Project Description

The analysis aims to uncover patterns within a dataset containing information about loan applications, focusing on identifying factors influencing the likelihood of loan default. The dataset encompasses two distinct scenarios: when customers apply for loans and face either payment difficulties or all other cases.

In the former scenario, customers experienced late payments exceeding a certain threshold on initial instalments, while the latter encompasses cases where payments were made on time.

The project's primary objective is to employ Exploratory Data Analysis (EDA) techniques to discern how customer attributes and loan characteristics contribute to the probability of default. By understanding these patterns, the company can make informed decisions regarding loan approval, such as adjusting loan amounts, interest rates, or denying loans to high-risk applicants.

Key Business Objectives:

- Identify predictive factors indicating whether a customer is likely to encounter difficulties in repaying their loan instalments.
- Improve decision-making processes related to loan approval by leveraging insights gleaned from EDA.
- Mitigate financial risk by proactively identifying high-risk loan applicants and implementing appropriate measures, such as adjusting loan terms or denying loans altogether.

Approach

Upon reviewing the dataset description, the files were downloaded, including:

1. Application Data: The primary dataset containing comprehensive information about current loan applicants.
2. Previous Application Data: Historical data on past loan applicants.
3. Column Description: Detailed descriptions of each feature present in the dataset, aiding in understanding the data attributes.
4. Important Notes for Project 6: Providing insights and instructions pertinent to the project.

Following the download, each file was meticulously examined to grasp its contents and relevance to the project objectives. Here's a breakdown of the approach:

- Application Data: Focused on extracting insights related to current loan applicants.
- Previous Application Data: Analysis of past applicant behavior and trends.
- Column Description: Understanding the nuances of each feature to facilitate a deeper understanding of the dataset's attributes.
- Important Notes for Project 6: Utilized for crucial information and guidelines essential for executing the project effectively.

Additionally, supplementary research was conducted through relevant articles and resources on risk analysis, ensuring a comprehensive grasp of the project's objectives and methodologies. This preparatory phase equipped with the necessary insights to embark on the subsequent stages of data exploration and analysis effectively.

Tech Stack Used

Microsoft Excel 2021 version is used for this project due to its simplicity of use and extraordinary analysis and visualization capability.

A. Identify Missing Data and Deal with it Appropriately:

- In the application data file, several columns had null (NaN) values. Here are the new columns I introduced from the existing ones:

1. **Average EXT_SOURCE**: Created from EXT_SOURCE_1, EXT_SOURCE_2, and EXT_SOURCE_3. As these columns contained float scores, but with varying percentages of null data (56.35%, 0.25%, and 19.89% respectively), I computed their average to utilize all available data effectively.
2. **Total FLAG_DOCUMENT**: Consolidated from FLAG_DOCUMENT_2 to FLAG_DOCUMENT_21. While these columns had no null values, their unspecified document references prompted total aggregation, facilitating easy analysis of document submission patterns.
3. **Years Birth**: Derived from DAYS_BIRTH, converting days to years.
4. **Years Employed**: Derived from DAYS_EMPLOYED, converting days to years.
5. **Years Registration**: Derived from DAYS_REGISTRATION, converting days to years.
6. **Years ID_PUBLISH**: Derived from DAYS_ID_PUBLISH, converting days to years.

- Columns dropped due to insufficient data (with null values percentage):

OWN_CAR_AGE	(65.90%),	OCCUPATION_TYPE	(31.31%),
APARTMENTS_AVG	(50.77%),	BASEMENTAREA_AVG	(58.40%),
YEARS_BEGINEXPLUATATION_AVG	(48.79%),	YEARS_BUILD_AVG	(66.48%),
COMMONAREA_AVG	(69.92%),	ELEVATORS_AVG	(53.30%),
ENTRANCES_AVG	(50.39%),	FLOORSMAX_AVG	(49.75%),
FLOORSMIN_AVG	(67.79%),	LANDAREA_AVG	(59.44%),
LIVINGAPARTMENTS_AVG	(68.45%),	LIVINGAREA_AVG	(50.28%),
NONLIVINGAPARTMENTS_AVG	(69.43%),	NONLIVINGAREA_AVG	
(55.15%),	APARTMENTS_MODE	(50.77%),	BASEMENTAREA_MODE
(58.40%),	YEARS_BEGINEXPLUATATION_MODE	(48.79%),	
YEARS_BUILD_MODE	(66.48%),	COMMONAREA_MODE	(69.92%),

ELEVATORS_MODE (53.30%), ENTRANCES_MODE (50.39%),
FLOORSMAX_MODE (49.75%), FLOORSMIN_MODE (67.79%),
LANDAREA_MODE (59.44%), LIVINGAPARTMENTS_MODE (68.45%),
LIVINGAREA_MODE (50.28%), NONLIVINGAPARTMENTS_MODE (69.43%),
NONLIVINGAREA_MODE (55.15%), APARTMENTS_MEDI (50.77%),
BASEMENTAREA_MEDI (58.40%), YEARS_BEGINEXPLUATATION_MEDI
(48.79%), YEARS_BUILD_MEDI (66.48%), COMMONAREA_MEDI (69.92%),
ELEVATORS_MEDI (53.30%), ENTRANCES_MEDI (50.39%),
FLOORSMAX_MEDI (49.75%), FLOORSMIN_MEDI (67.79%),
LANDAREA_MEDI (59.44%), LIVINGAPARTMENTS_MEDI (68.45%),
LIVINGAREA_MEDI (50.28%), NONLIVINGAPARTMENTS_MEDI (69.43%),
NONLIVINGAREA_MEDI (55.15%), FONDKAPREMONT_MODE (68.38%),
HOUSETYPE_MODE (50.15%), TOTALAREA_MODE (48.30%),
WALLSMATERIAL_MODE (50.92%), EMERGENCYSTATE_MODE (47.40%)

- Columns dropped for least relevance:

- WEEKDAY_APPR_PROCESS_START
- HOUR_APPR_PROCESS_START

- Columns with null values imputed:

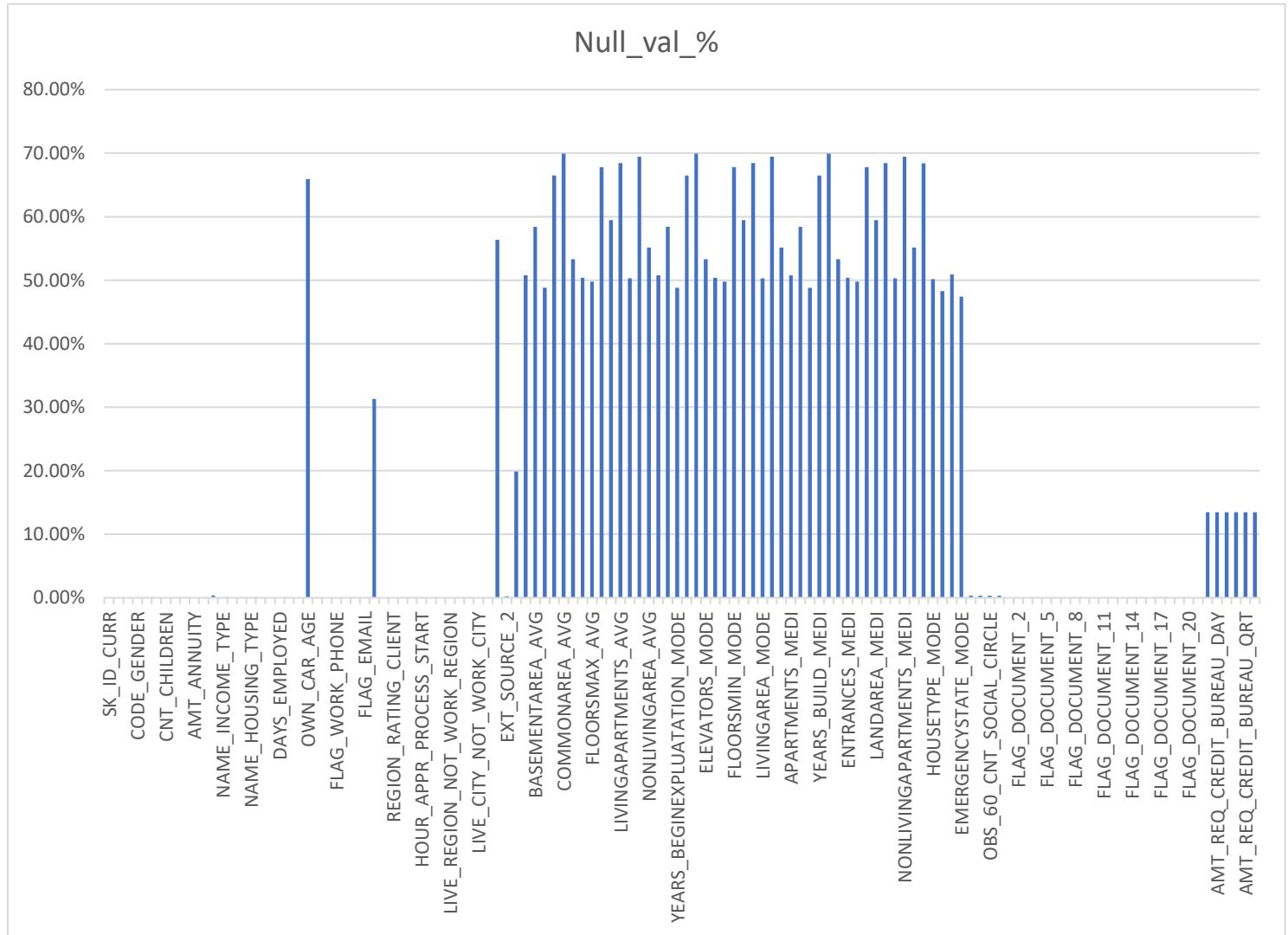
- AMT_GOODS_PRICE: Imputed with 90% of AMT_CREDIT values.
- NAME_TYPE_SUITE: Imputed with the most frequently occurring value using Index Match Max COUNTIF function.
- OBS_30_CNT_SOCIAL_CIRCLE, DEF_30_CNT_SOCIAL_CIRCLE, -
OBS_60_CNT_SOCIAL_CIRCLE, DEF_60_CNT_SOCIAL_CIRCLE: Imputed with rounded averages.

- Columns with null values removed:

- All six columns related to credit bureau inquiries - namely AMT_REQ_CREDIT_BUREAU_HOUR, AMT_REQ_CREDIT_BUREAU_DAY, AMT_REQ_CREDIT_BUREAU_WEEK, AMT_REQ_CREDIT_BUREAU_MON, AMT_REQ_CREDIT_BUREAU_QRT, and AMT_REQ_CREDIT_BUREAU_YEAR - may not appear immediately pertinent to this dataset. However, in the context of creditor analysis and borrower evaluation, they hold significant importance.

- Understanding risk analysis underscores the critical role of an applicant's credit history, akin to the concept of a credit score. These columns capture the inquiries made by applicants to the credit bureau, which are typically regarded as requests for credit-related information.
- In risk assessment, a high frequency of such inquiries over a short period often indicates a "credit-hungry" behavior, suggesting potential financial instability or overreliance on credit. Therefore, despite their seemingly indirect relevance, these columns were deemed important. Consequently, null values were dropped to ensure a comprehensive analysis, aligning with the broader objective of identifying patterns indicative of loan repayment capability and default likelihood.

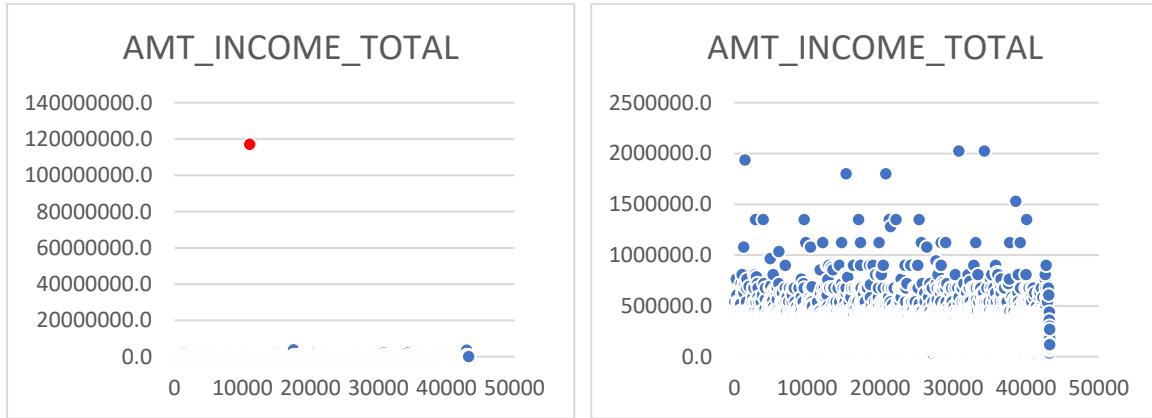
This graph indicates null values in various features:



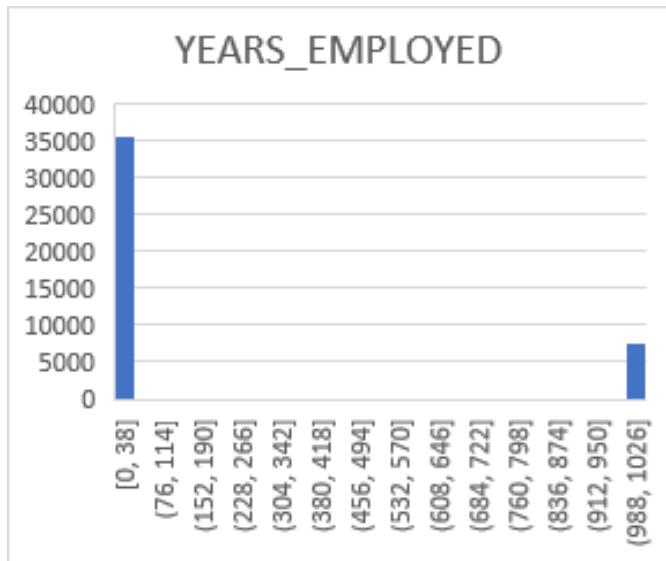
B. Identify Outliers in the Dataset:

To identify outliers within the dataset, I employed the interquartile range (IQR) method, which involves computing the first quartile (Q1) and the third quartile (Q3) at the 25th and 75th percentiles, respectively. I then utilized a multiplier value (k) of 1.5 to establish thresholds for the lower and upper limits.

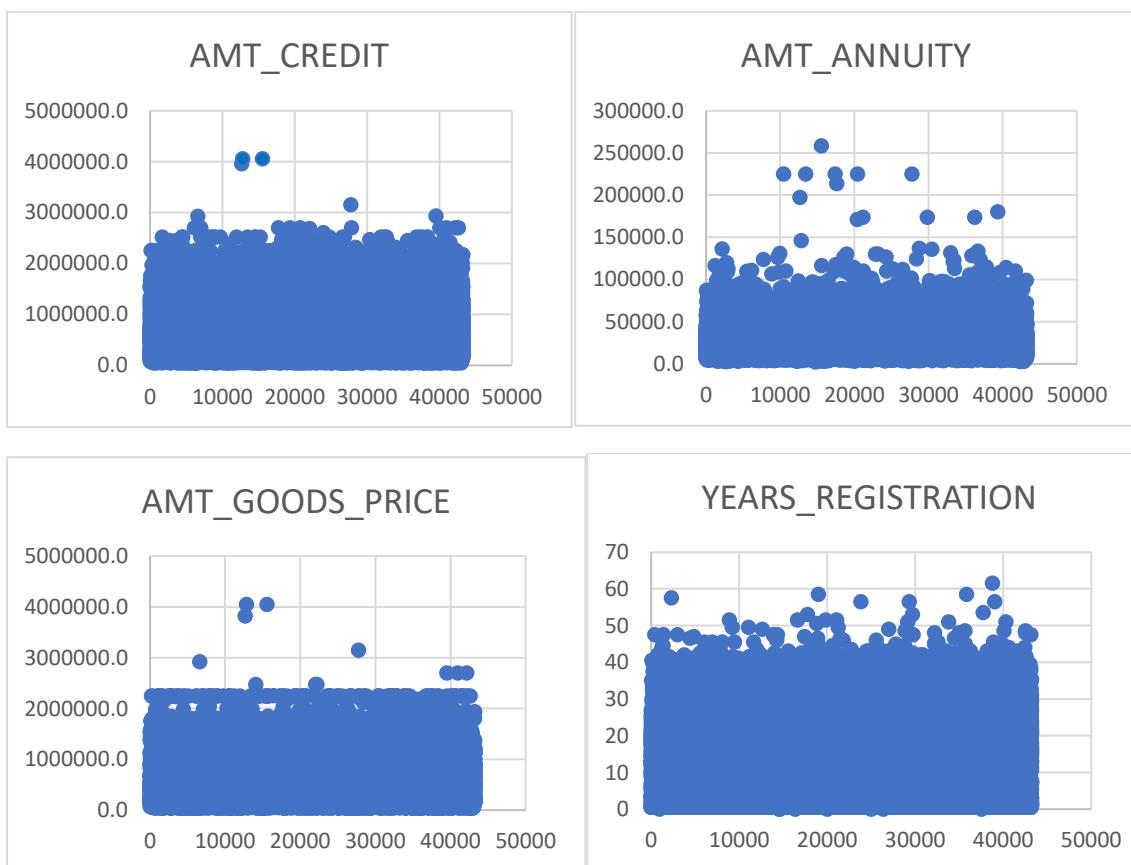
Upon examination of the column AMT_INCOME_TOTAL, an outlier was detected with a value of 117,000,000. This value stands out significantly from the rest of the dataset, as illustrated in the scatterplot graph that depicts the data both with and without the outlier [3825000, 3600000, 117000000]



Similarly, in the column YEARS_EMPLOYED, outliers were observed, suggesting potential typographical errors. Notably, one entry indicates an applicant's employment duration of 1001 years, which is clearly implausible. This anomaly is evident in the accompanying graph.



For further insight, scatterplot graphs were generated for columns such as AMT_CREDIT, AMT_ANNUITY, AMT_GOODS_PRICE, and YEARS_REGISTRATION. These visualizations aid in identifying any additional outliers and understanding their impact on the dataset's distribution and overall analysis.

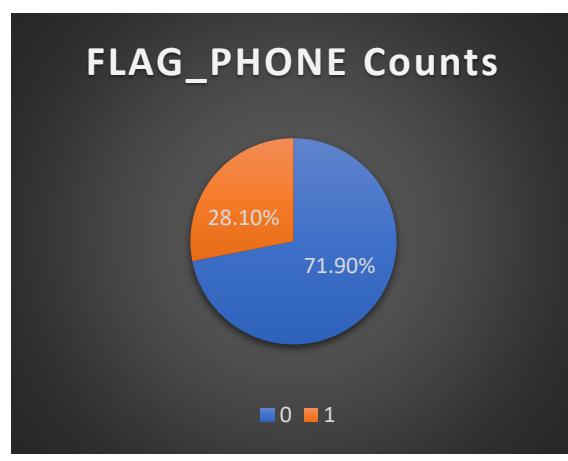
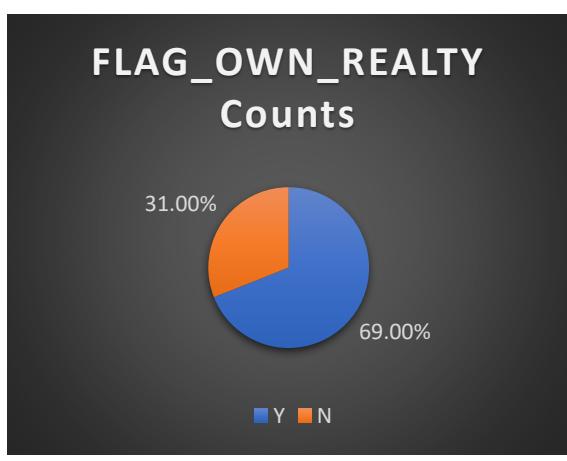
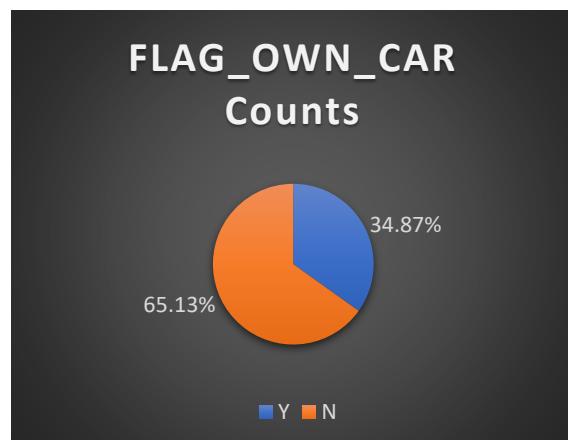
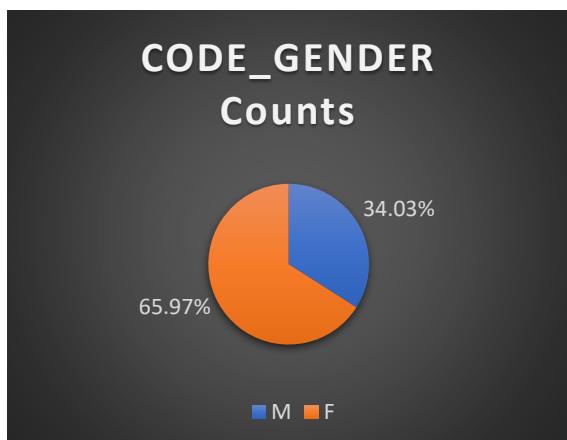


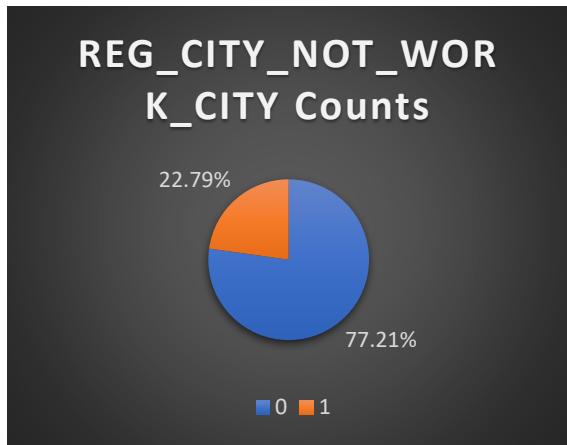
C. Analyze Data Imbalance:

The dataset reveals a notable imbalance in certain columns, with some features displaying varying degrees of disproportion between different categories. Here's an examination of the data imbalance across three categories:

1. Mild Imbalance [20-40% minority proportion]:

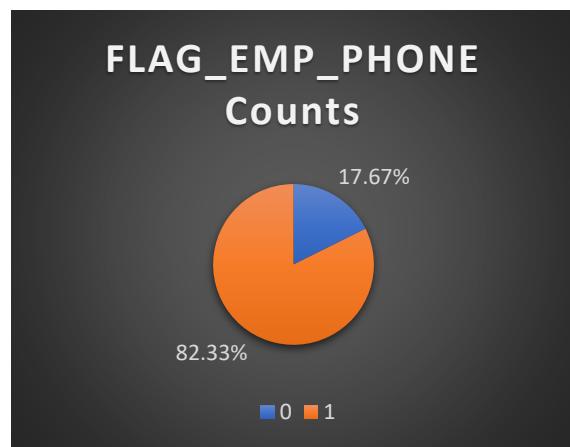
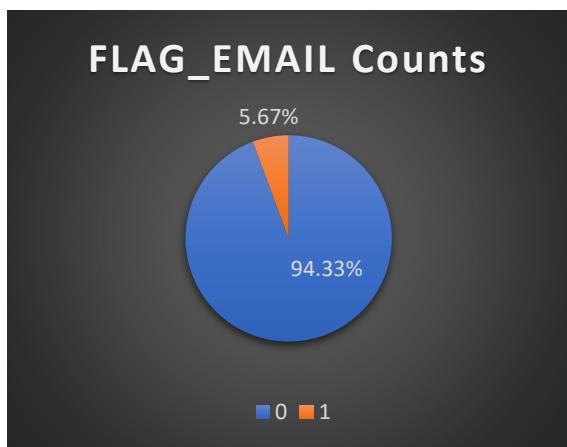
- CODE_GENDER: 34.03% Male against 65.97% Female
- FLAG_OWN_CAR: 34.87% "N" Value against 65.13% "Y" Value
- FLAG_OWN_REALTY: 31% "N" value against 69% "Y" value
- FLAG_PHONE: 28.10% "1" value against 71.90% "0" value
- REG_CITY_NOT_WORK_CITY: 22.79% "1" value against 77.21% "0" value



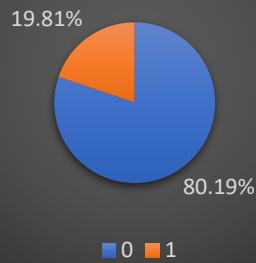


2. Moderate Imbalance [1-20% minority proportion]:

- FLAG_EMAIL: 5.67% "1" value against 94.33% "0" value
- FLAG_EMP_PHONE: 17.67% "1" value against 82.33% "0" value
- FLAG_WORK_PHONE: 19.81% "1" value against 80.19% "0" value
- REG_REGION_NOT_WORK_REGION: 4.89% "1" value against 95.11% "0" value
- LIVE_CITY_NOT_WORK_CITY: 17.78% "1" value against 82.22% "0" value



FLAG_WORK_PHONE Counts



REG_REGION_NOT_WORK_REGION Counts



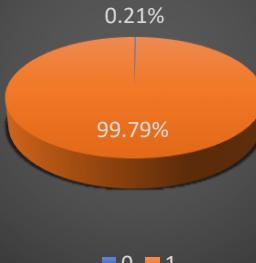
LIVE_CITY_NOT_WORK_CITY Counts



3. **Extreme Imbalance [less than 1% minority proportion]:**

- FLAG_CONT_MOBILE: 0.21% "0" value against 99.79% "1" value

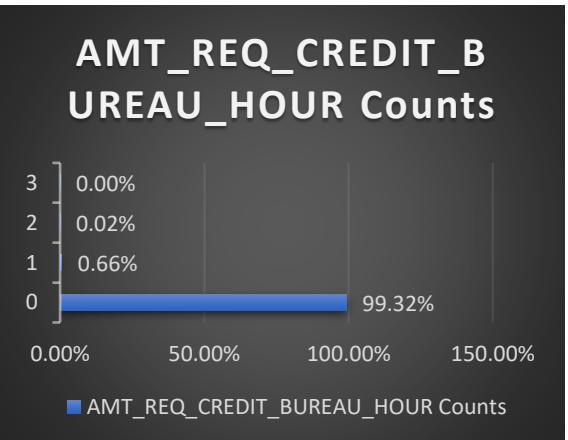
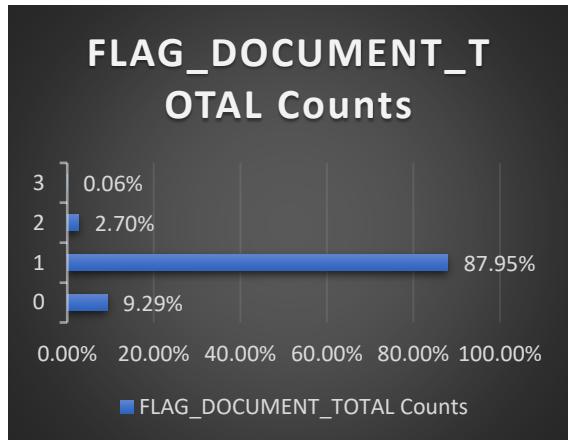
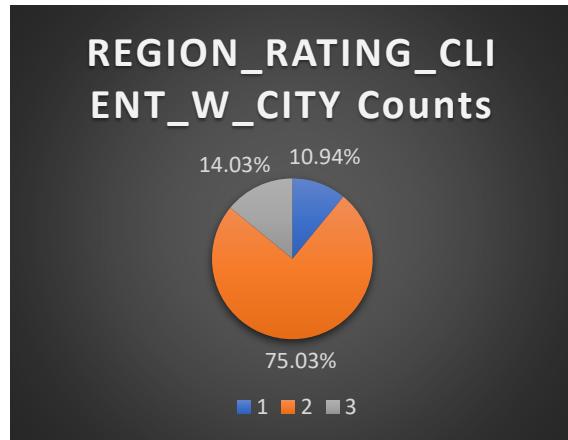
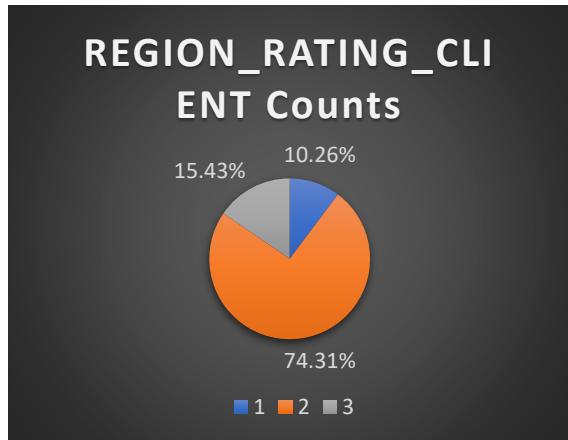
FLAG_CONT_MOBILE Counts



now here are those imbalances that has more than two variables

- REGION_RATING_CLIENT has 10.26% of "1" value, 74.31% of "2" value, and 15.43% of "3" value.
- REGION_RATING_CLIENT_W_CITY has 10.94% of "1" value, 75.03% of "2" value, and 14.03% of "3" value.

- FLAG_DOCUMENT_TOTAL has 9.29% of "0" value, 87.95% of "1" value, 2.70% of "2" value, and 0.06% of "3" value.
- AMT_REQ_CREDIT_BUREAU_HOUR has 99.32% of "0" value, 0.66% of "1" value, 0.02% of "2" value, and 0.00% of "3" value.

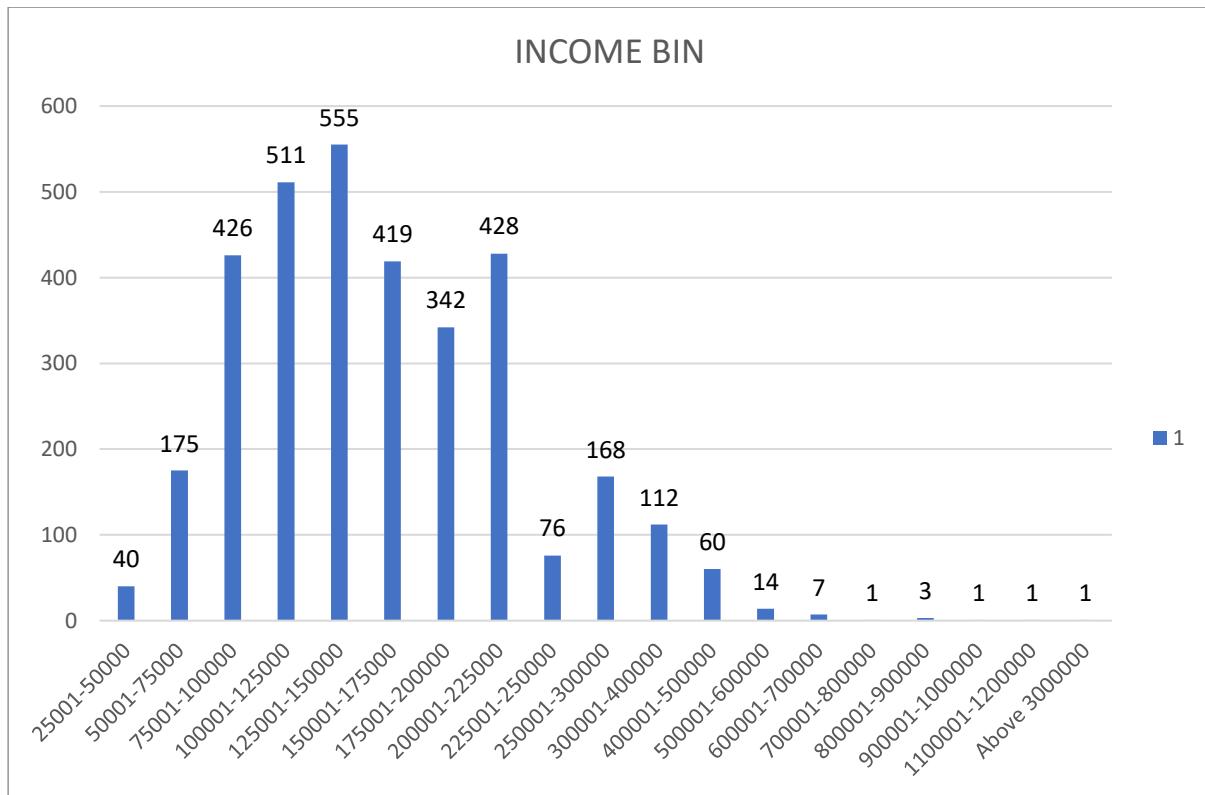


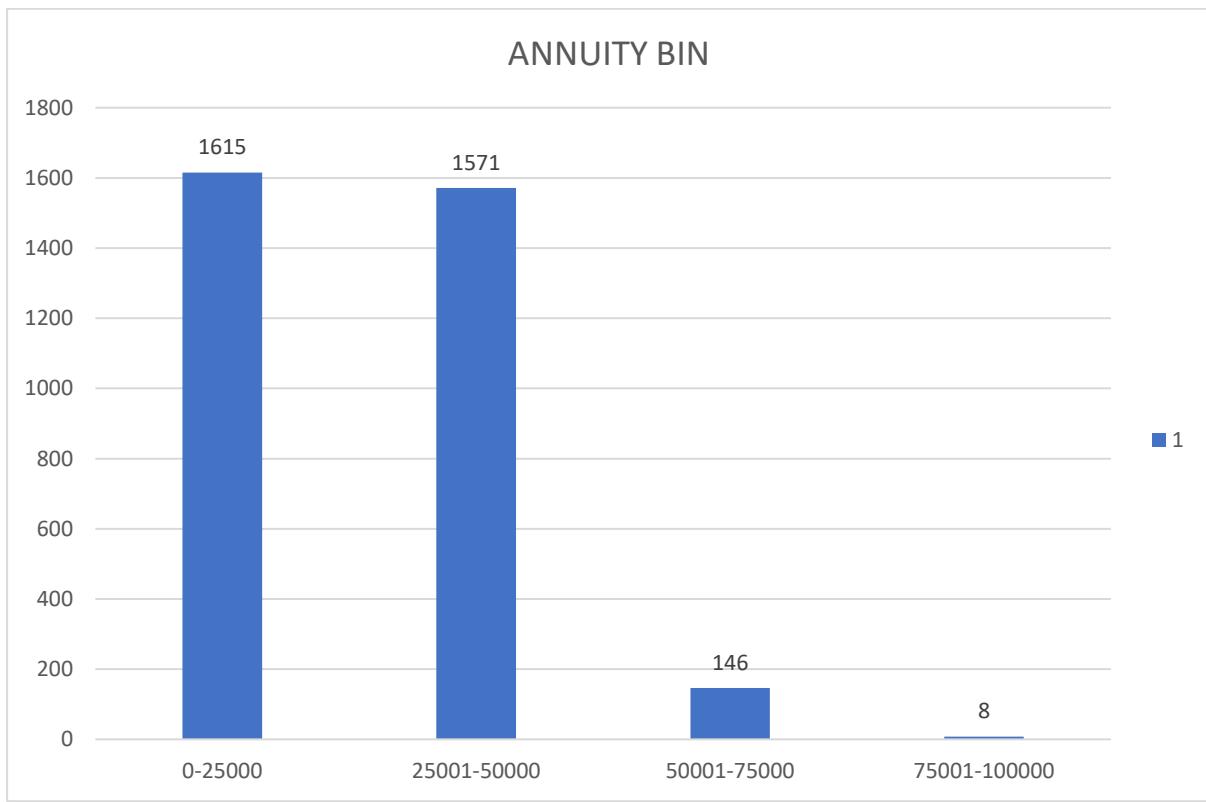
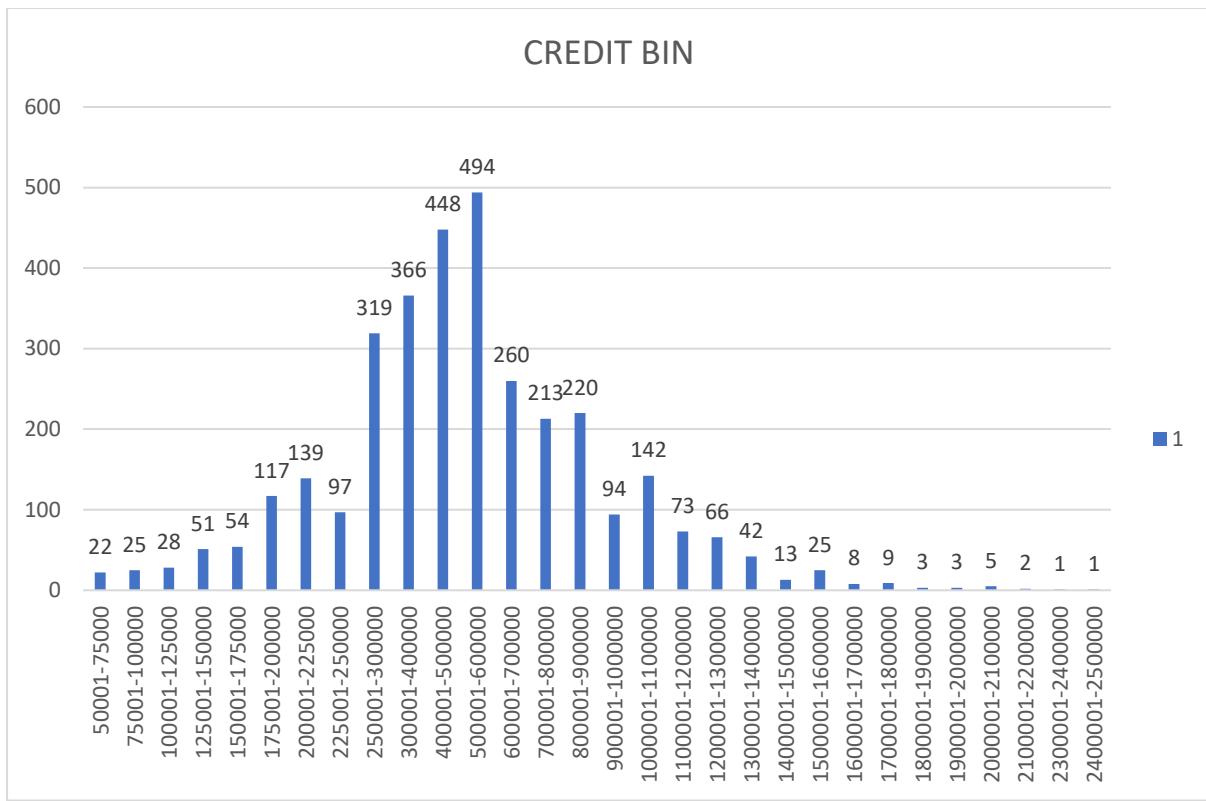
D. Perform Univariate, Segmented Univariate, and Bivariate Analysis

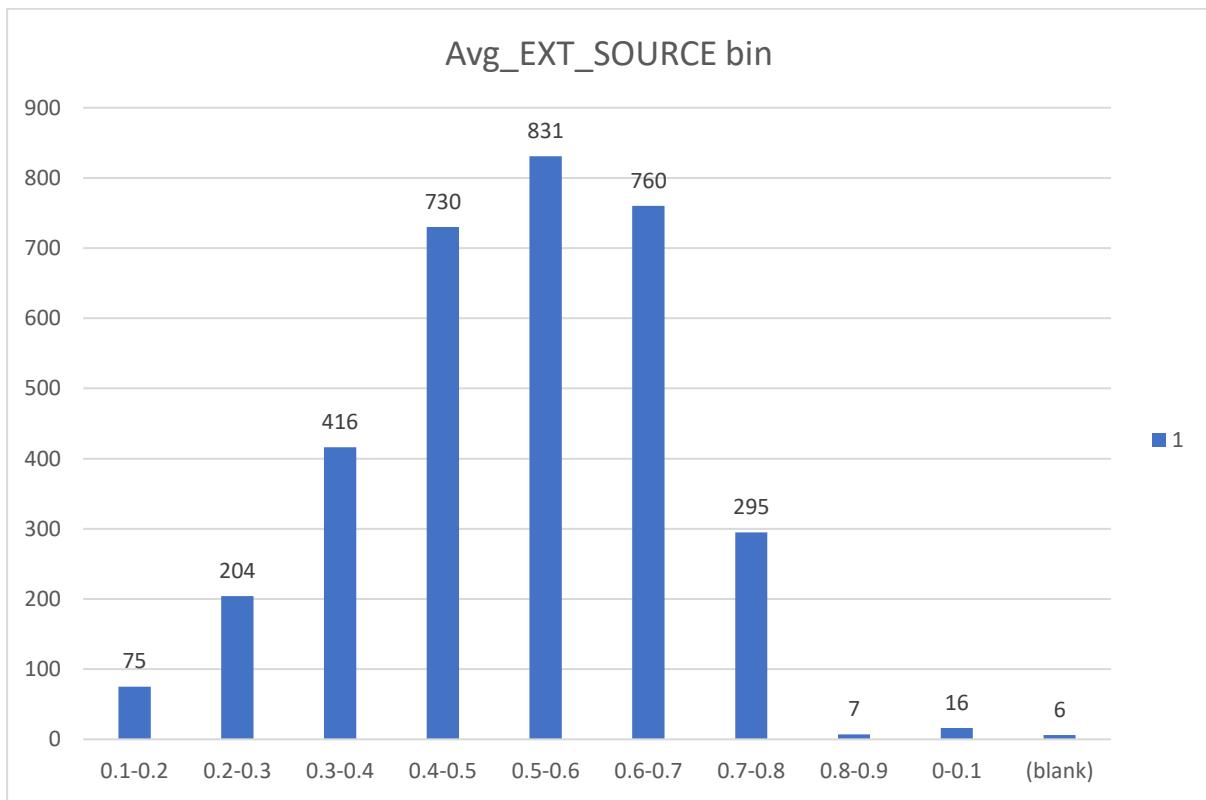
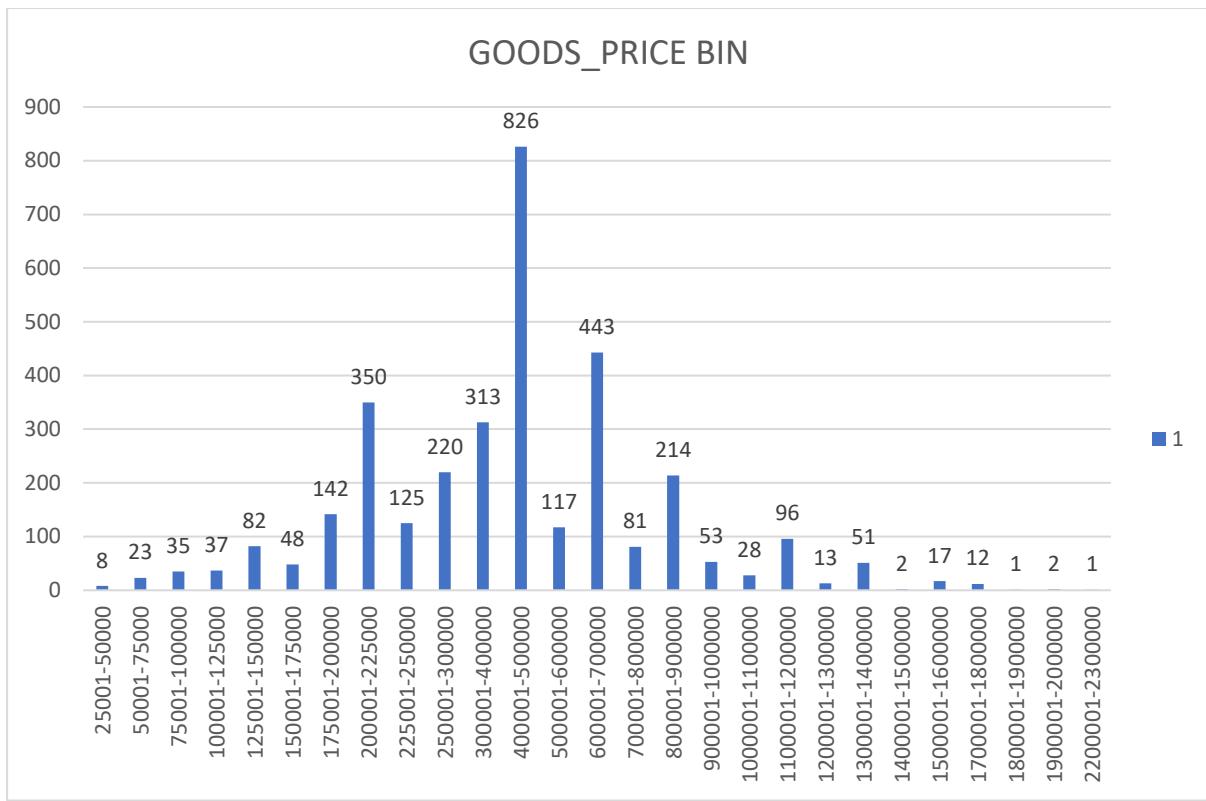
To conduct the analysis effectively, several key columns were selected, including AMT_INCOME_TOTAL, AMT_CREDIT, AMT_ANNUITY, AMT_GOODS_PRICE, Avg_EXT_SOURCE, and YEARS_EMPLOYED.

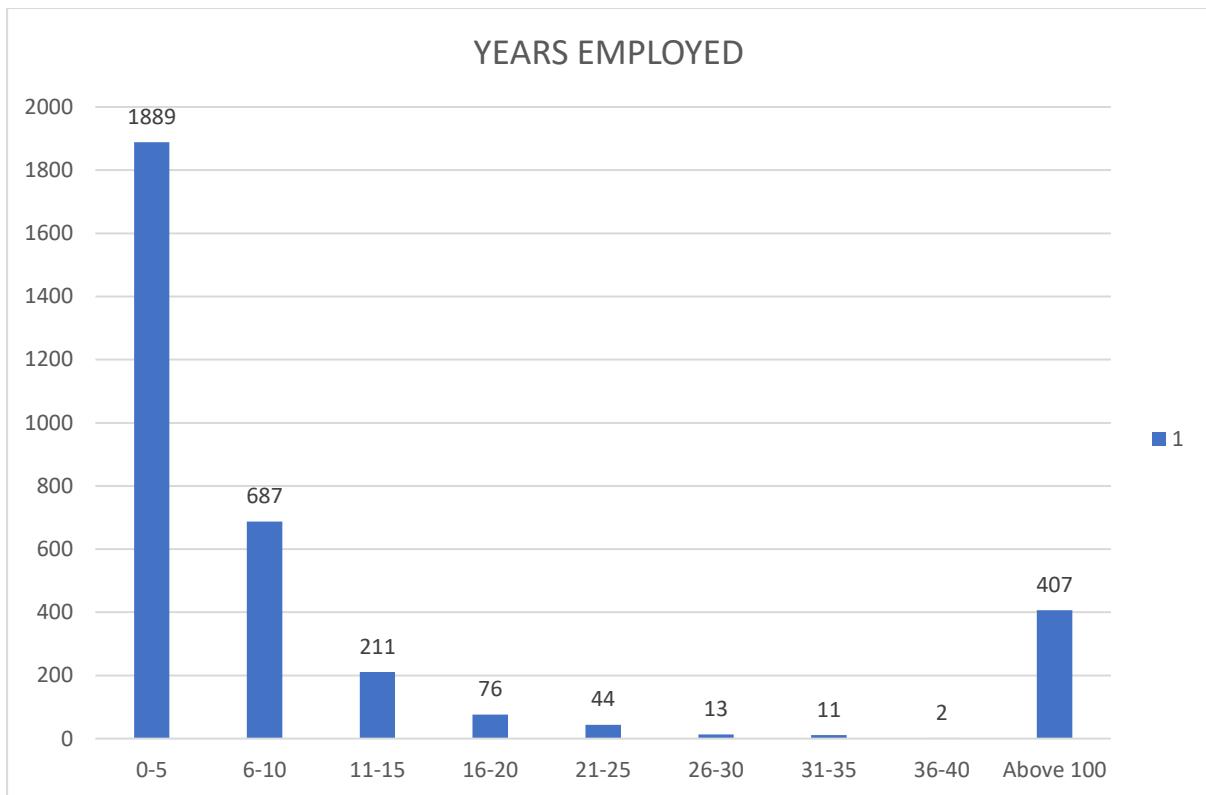
Next, the values in these columns were categorized into bins using nested IF and AND functions, resulting in the creation of six new columns representing the bin values.

Subsequently, pivot tables were generated from these columns to perform **univariate analysis**, comparing each of the six columns with the Target value of 1. This allowed for an understanding of how each column correlates with loan default.

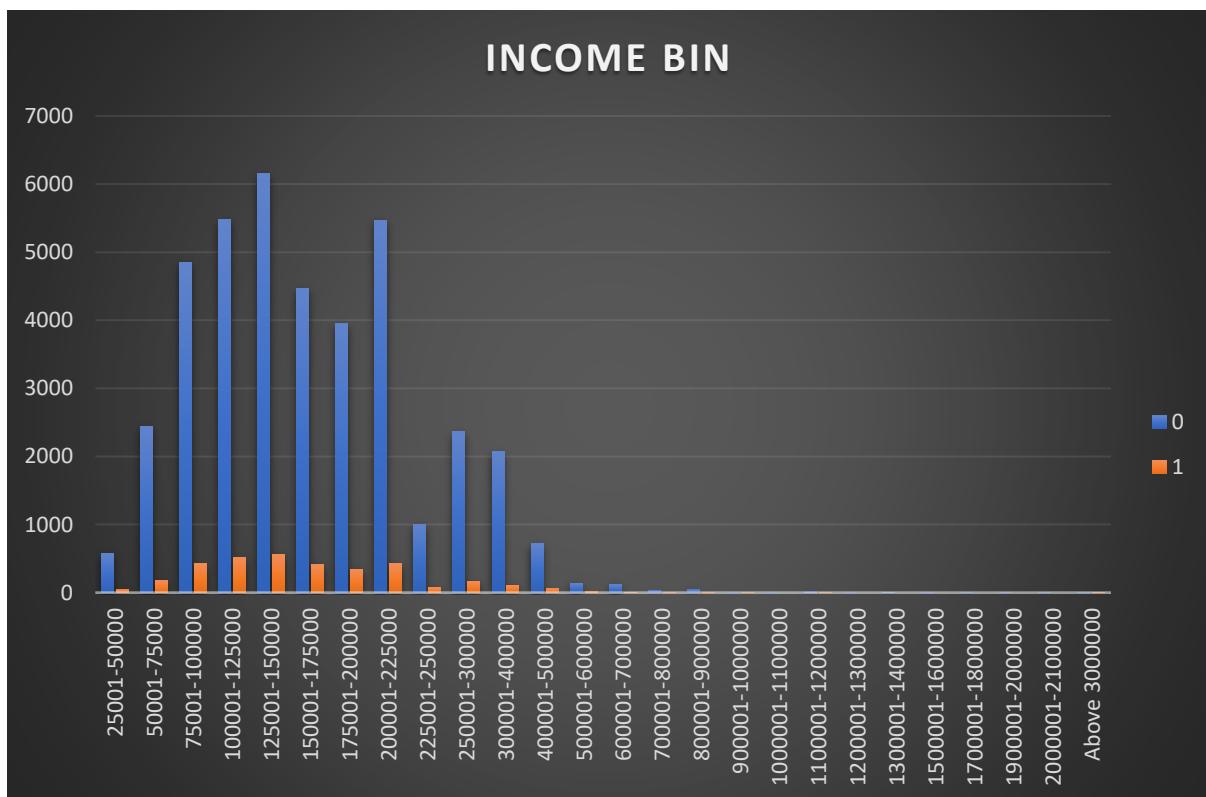


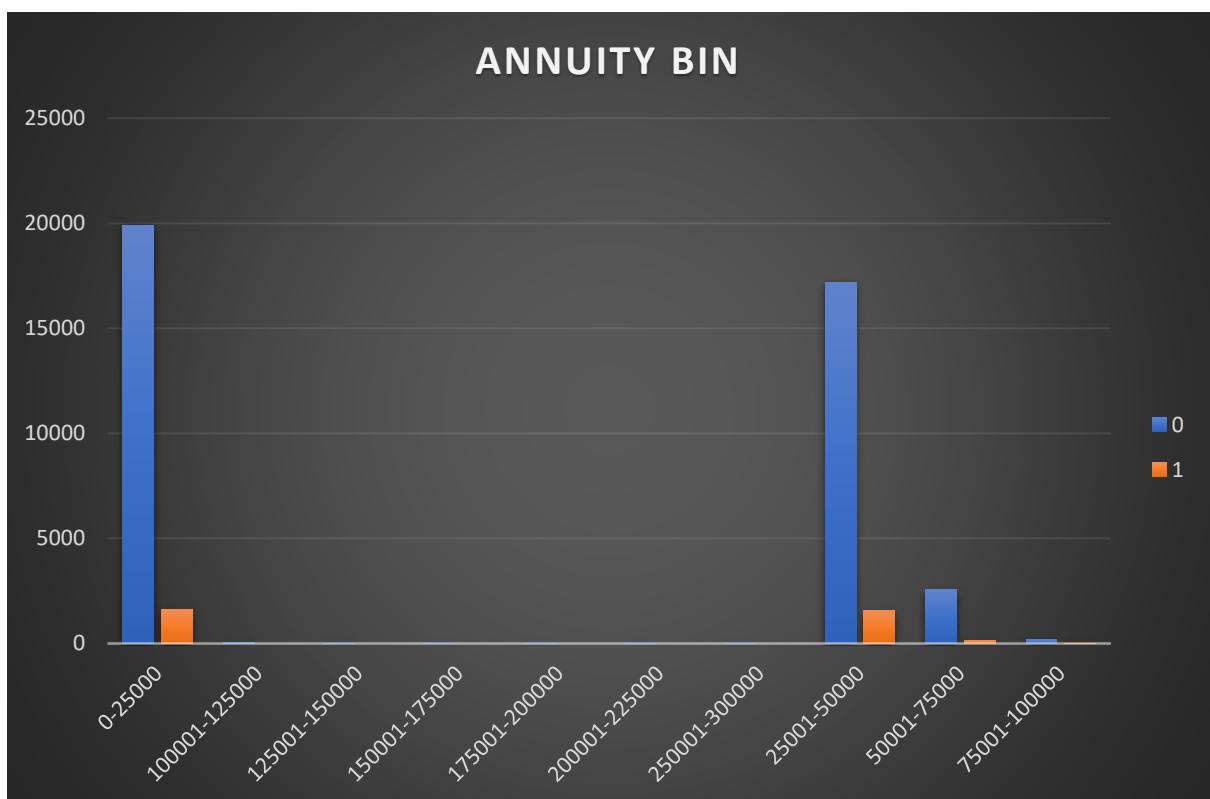
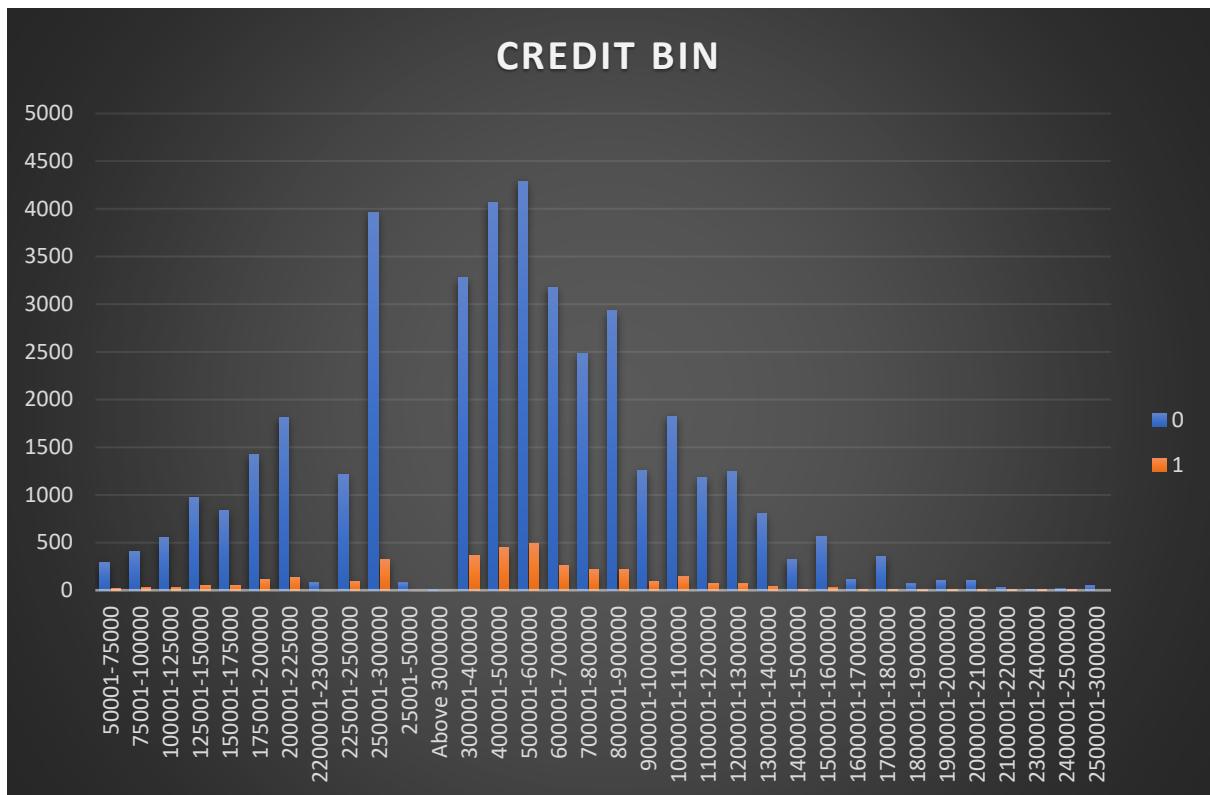


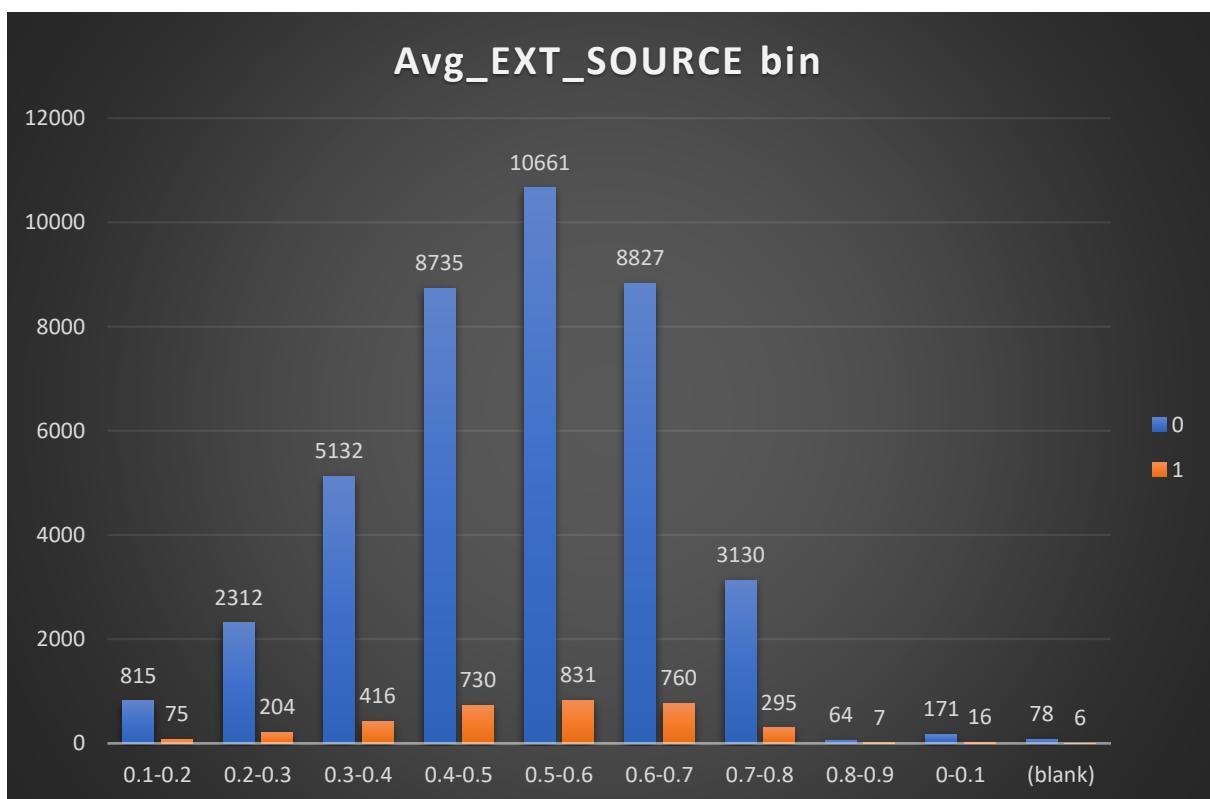


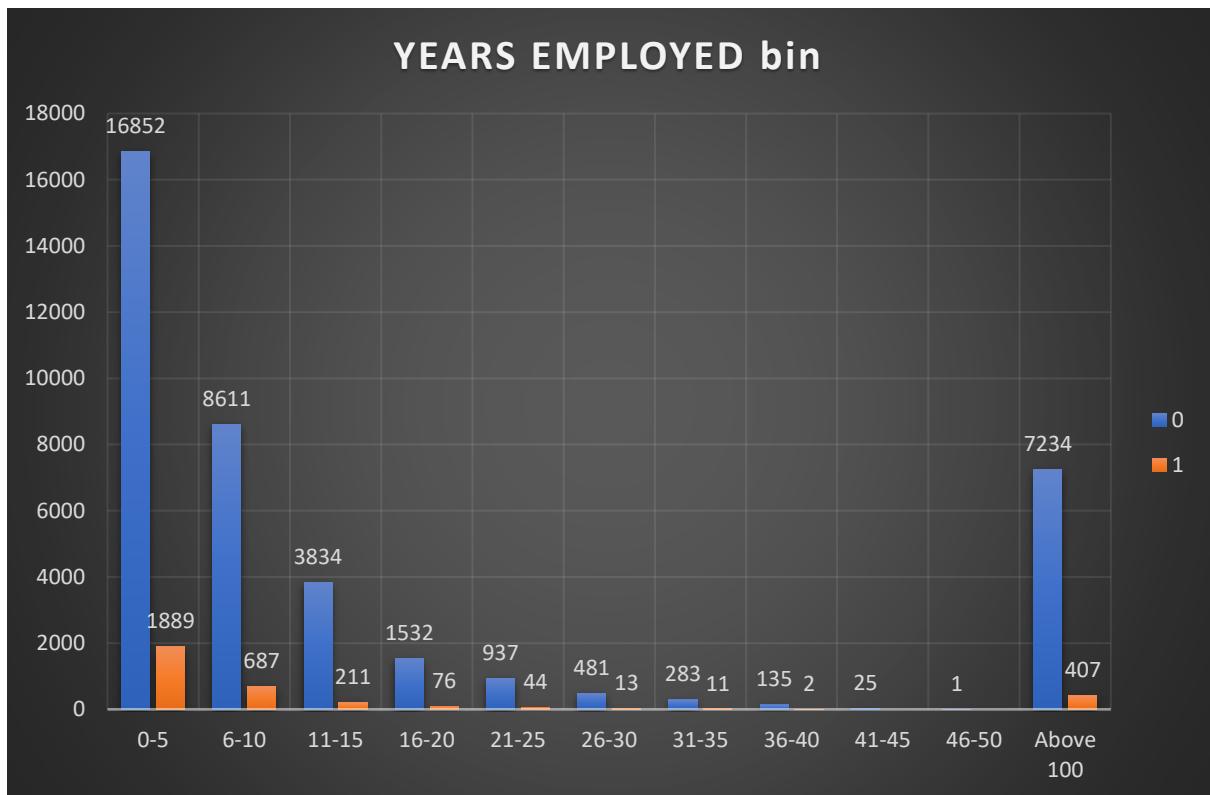


Following the univariate analysis, **segmented univariate analysis** was conducted, enabling both Target 0 and 1 against each of the six columns individually. This segmentation provided insights into how different values within each column impact loan default rates.

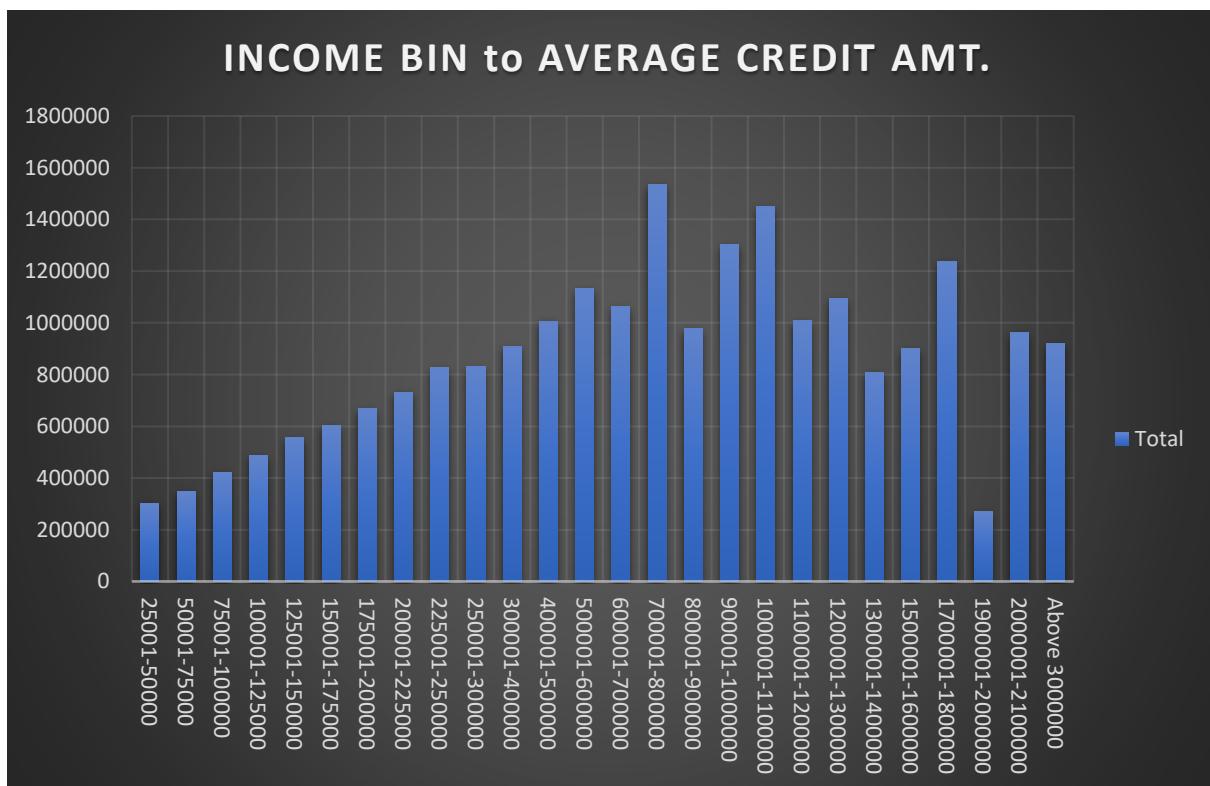


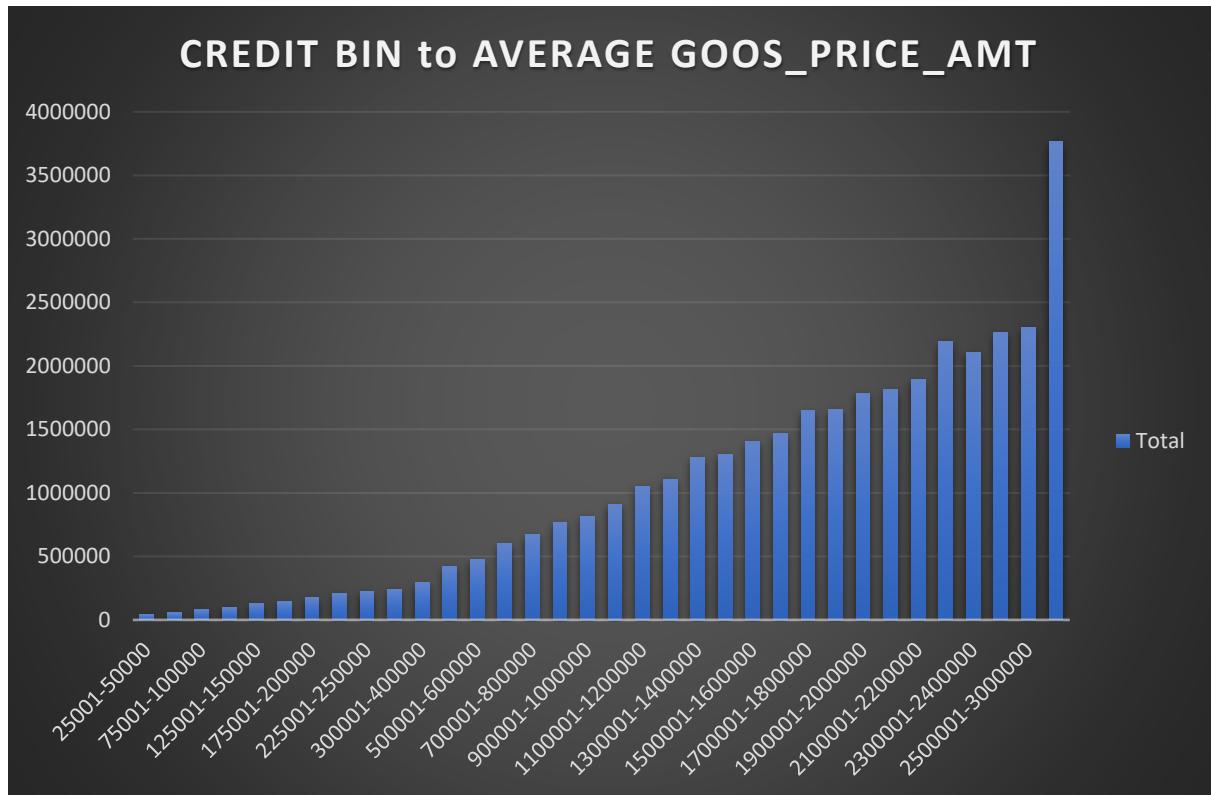
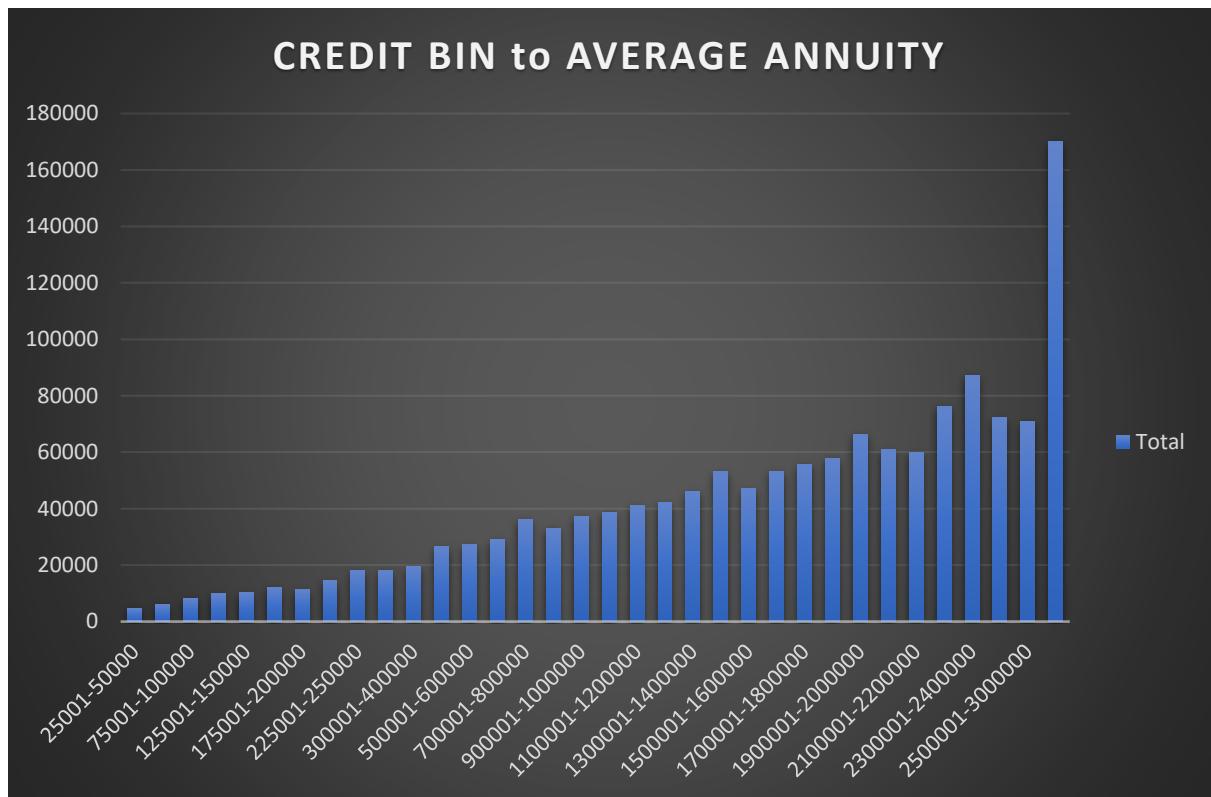






Finally, **bivariate analysis** was conducted by examining relationships between certain column pairs. For instance, the relationship between income bin and average credit amount, credit bin and average annuity, and credit bin and average goods price amount were analyzed.





This comprehensive approach helped gain deeper insights into the factors influencing loan default and enabled informed decisions regarding loan approval processes.

E. Identify Top Correlations for Different Scenarios:

Top correlations were identified through three iterations of analysis:

1. All Applicants:

- Avg_EXT_SOURCE: -0.226
- YEARS_BIRTH: -0.080
- REGION_RATING_CLIENT_W_CITY: 0.065
- REGION_RATING_CLIENT: 0.064
- DAYS_LAST_PHONE_CHANGE: 0.054

2. Applicants with Payment Difficulties:

- FLAG_EMP_PHONE vs. YEARS_EMPLOYED: -0.999
- OBS_60_CNT_SOCIAL_CIRCLE vs. OBS_30_CNT_SOCIAL_CIRCLE: 0.998
- AMT_GOODS_PRICE vs. AMT_CREDIT: 0.982
- REGION_RATING_CLIENT_W_CITY vs. REGION_RATING_CLIENT: 0.949
- FLAG_EMP_PHONE vs. YEARS_BIRTH: -0.586

3. All Other Cases:

- FLAG_EMP_PHONE vs. YEARS_EMPLOYED: -0.999
- OBS_60_CNT_SOCIAL_CIRCLE vs. OBS_30_CNT_SOCIAL_CIRCLE: 0.998
- AMT_GOODS_PRICE vs. AMT_CREDIT: 0.982
- REGION_RATING_CLIENT_W_CITY vs. REGION_RATING_CLIENT: 0.949
- FLAG_EMP_PHONE vs. YEARS_BIRTH: -0.586

These correlations provide valuable insights into the relationships between different variables under different scenarios, aiding in better decision-making processes.

All variable with complete dataset

	TARGET	T_CHILDREN	INCOME_TMT	CREDIT_ANNUIT	GOODS_PIPULATIONARS	BIRTRS_EMPLO	REGISTRAS_ID	PUB3_EMP_PHI	WORK_Pt	CONT_MCAT	PHONLAG	EMAILFAM	MEMI	RATING	TING	CLEN	NOT_LII	NOT_WCI	NOT_WCY	NOT_LIV	NOT_WC	NOT_WL	EXT_SOINT	SOCIAINT	SOCIAINT	SOCIAINT	PHONE_DOCUMENT	REDIT_BU	CREDIT_BU	REDIT_BU	CREDIT_BU	REDIT_BU						
TARGET	1	0.02883	0.013293	-0.03282	-0.01295	-0.04212	-0.03858	-0.07982	-0.04265	-0.03914	-0.04318	0.04156	0.017453	0.007608	-0.03515	-0.00083	0.016732	0.064078	0.065049	0.009161	-0.00098	0.0353	0.044985	0.031545	-0.22635	0.01749	0.043394	0.017325	0.04564	0.053953	0.021495	0.003258	0.011957	0.005731	-0.01136	-0.00081	0.02365	
CNT_CHIL	0.02883	1	0.009269	0.004098	0.026168	-0.00075	-0.02524	-0.33888	-0.24249	-0.1807	0.035493	0.21545	0.051013	0.002567	-0.03484	0.026711	0.881809	0.026272	0.023709	-0.00901	0.010091	0.01792	0.023021	0.073763	0.069149	-0.07276	0.016554	-0.00335	0.016485	-0.00435	0.008897	-0.01609	0.001963	-0.00213	0.000181	-0.01353	-0.00917	-0.04136
AMT_INCC	0.013293	0.009269	1	0.064063	0.07689	0.06447	0.026987	-0.01667	-0.02976	-0.00928	-0.00517	0.029787	-0.00899	-0.00368	-0.00238	0.01479	0.010436	-0.03559	-0.03812	0.01315	0.0262	0.024958	0.001132	0.000243	0.001244	0.003087	-0.00826	-0.00696	-0.00824	-0.00664	-0.00275	0.008366	0.000718	0.000982	0.000264	0.011521	-0.00044	0.01894
AMT_CREI	-0.03282	0.004098	0.064063	1	0.769379	0.986809	0.097991	0.04803	-0.07376	-0.00686	0.001141	0.074958	-0.00874	0.026308	0.021134	0.011384	0.06304	-0.10241	-0.11251	0.02562	0.057603	0.056951	-0.02326	-0.0148	0.004863	0.143936	0.000538	-0.01407	0.000845	-0.01853	-0.06808	0.226006	-0.00042	0.010845	0.001691	0.064058	0.016477	-0.04361
AMT_ANN	0.01295	0.026168	0.07689	0.769379	1	0.774177	0.115326	-0.01451	-0.11023	-0.0346	-0.01315	0.110971	-0.01752	0.026769	0.007234	0.069267	0.078187	-0.12315	-0.13868	0.043334	0.080867	0.074963	-0.00447	0.002011	0.011798	0.116622	-0.00838	-0.01869	-0.00823	-0.02264	-0.06375	0.201541	0.012316	0.007004	0.020182	0.040963	0.006668	-0.00819
AMT_GOO	-0.04212	-0.00075	0.06447	0.986809	0.774177	1	0.102311	0.046781	-0.07087	-0.01014	0.003004	0.072119	0.012768	-0.023667	0.036795	0.010744	0.060777	-0.10481	-0.11403	0.027334	0.05844	0.056979	-0.02285	-0.01595	0.003219	0.153125	-0.00021	-0.01511	3.76E-05	-0.01907	-0.07205	0.196892	0.000182	0.01139	0.002146	0.065871	0.017378	-0.04698
REGION_P	0.03858	-0.02524	0.026987	0.097991	0.115326	0.102311	1	0.030865	-0.00198	0.056775	0.004244	0.00199	-0.01641	-0.00584	0.095109	0.042182	-0.02367	-0.52864	-0.52581	-0.00569	0.061607	0.08702	-0.04358	-0.03752	-0.01118	0.140705	-0.1623	0.010843	-0.01514	0.007295	-0.04712	-0.00928	-0.0024	-0.001	0.00364	0.078138	0.006538	
YEARS_BIF	-0.07982	-0.33888	-0.01667	0.04803	-0.01451	0.046781	0.030865	1	0.62105	0.334463	0.248281	-0.61699	-0.17218	0.011267	0.044881	-0.09408	-0.28953	-0.01381	-0.01136	-0.05971	-0.09183	-0.06722	-0.17735	-0.2351	-0.15077	0.290203	-0.01216	5.89E-05	-0.01206	-0.00195	-0.07313	0.014673	-0.00508	-0.00003	0.016474	0.07074		
YEARS_EM	0.04265	-0.24249	-0.02976	0.07376	-0.11023	-0.07087	-0.00198	0.62105	1	0.207224	0.264799	-0.99972	-0.22997	0.015943	0.025083	-0.06728	-0.23154	0.035154	0.038299	-0.03901	-0.10635	-0.0945	-0.09245	-0.25354	-0.21652	0.103049	0.005845	0.016818	0.005814	0.014705	0.025667	0.009798	-0.0044	0.005308	-0.00413	-0.033	0.018456	0.047823
YEARS_REF	-0.03914	-0.1807	-0.00928	0.00686	-0.0346	-0.01014	0.005675	0.334463	0.207224	1	0.094513	-0.20491	-0.05669	-0.0291	0.072419	-0.03038	-0.17082	-0.08473	-0.07751	-0.02811	-0.03588	-0.02448	-0.06693	-0.09385	-0.06343	0.136549	-0.0118	-0.00264	0.01514	0.003363	0.022994	0.0006522	0.011099	-0.00305	0.023646			
YEARS_ID	-0.04318	0.035493	-0.00517	0.001141	-0.01315	0.003004	0.004244	0.248281	0.264799	0.094513	1	0.26391	-0.0474	0.007723	0.031069	-0.03515	0.026045	0.003261	0.007708	-0.02924	-0.04551	-0.03353	-0.06937	-0.09524	-0.05813	0.121314	0.010182	-0.00124	0.010286	-0.0007	-0.08062	0.035528	-0.00712	-0.00578	-0.00315	0.00662	0.012588	0.034643
FLAG_EMF	0.04156	0.241545	0.029787	0.074958	0.110971	0.072119	0.00199	-0.61699	-0.99972	-0.20491	-0.26391	1	0.23023	-0.01597	-0.02435	0.066847	0.230999	-0.03501	-0.03821	0.038215	0.105065	0.093505	0.090866	0.251638	0.215376	-0.09774	-0.058	-0.01701	-0.00577	-0.01491	-0.02752	-0.00934	0.036516	-0.00538	0.004214	0.03334	-0.01871	-0.04791
FLAG_WOI	0.017453	0.051013	-0.00899	-0.00874	-0.01752	0.012768	-0.01641	-0.17218	-0.22997	-0.05669	-0.0474	0.23023	1	0.028219	0.296595	-0.01067	0.064344	0.003479	0.088806	0.067907	0.040186	0.051481	0.118603	0.105173	-0.05929	-0.02383	-0.01698	-0.02405	-0.01373	-0.0449	0.028699	-0.00914	-0.00558	-0.00405	-0.00791	-0.02931	-0.07395	
FLAG_CON	0.007603	0.002567	-0.00368	0.026308	0.026769	0.023667	-0.00584	0.011267	0.015943	-0.00291	0.007723	-0.01597	0.028219	1	0.004005	-0.01057	0.000769	0.015722	0.015999	0.001251	-0.00128	0.000114	0.003804	0.000886	0.001553	0.007076	0.006594	0.002028	0.006527	0.002638	-0.027	-0.05122	-0.00204	-0.01549	-0.02093	0.006301	0.001453	0.026287
FLAG_PHO	-0.03515	-0.03484	-0.00238	0.021134	0.007234	0.036795	0.095109	0.044881	0.025083	0.072419	0.031069	-0.02435	0.296595	0.004005	1	0.012758	-0.02019	-0.08887	0.08333	0.010515	-0.0021	-0.00582	-0.04001	-0.04316	-0.02581	0.055597	-0.03604	-0.02802	-0.03569	-0.02476	-0.06621	0.036398	-0.01012	0.000327	0.001952	0.043439	-0.0152	-0.02114
FLAG_EMA	-0.00083	0.026711	0.01479	0.011384	0.069267	0.010744	0.042182	-0.09408	-0.06728	-0.03038	-0.03515	0.066847	-0.01067	-0.01057	0.012758	1	0.021989	-0.05947	-0.05691	0.020406	0.040355	0.039709	0.013192	0.0107027	0.001386	-0.01228	-0.00443	-0.0182	-0.00389	-0.0024	-0.01788	0.000796	0.006393	0.003322	0.025054	0.021161	0.015871	0.046885
CNT_FAM	0.016732	0.881809	0.010436	0.06304	0.078187	0.060777	-0.02367	-0.28953	-0.23154	-0.17082	0.026045	0.230999	0.064344	0.007698	-0.02019	-0.021989	1	0.027315	0.026921	-0.01135	0.005403	0.013531	0.016671	0.076503	0.078873	-0.03928	0.02488	-0.00271	0.024823	-0.00407	0.01832	-0.00473	0.003203	-0.00328	0.001576	-0.00676	-0.00888	-0.02827
REGION_R	0.064078	-0.03515	-0.010241	-0.12135	-0.10481	-0.152864	-0.01381	0.035154	-0.08473	0.003261	-0.03501	0.003479	0.015722	-0.08887	-0.05947	0.027315	1	0.950228	-0.03905	-0.14187	-0.14836	0.033944	0.005854	-0.01897	-0.211	0.032036	0.009372	0.032041	0.01047	0.024348	-0.00106	0.007373	0.004743	-0.0006	-0.07131	0.011512	0.009303	
REGION_D	0.065049	0.023709	-0.03812	-0.11251	-0.13868	-0.11403	-0.052581	-0.01136	0.038299	-0.07751	0.007708	-0.03821	0.008806	0.015999	-0.08333	-0.05691	0.026921	0.950228	1	0.03486	-0.1346	-0.14203	0.04325	0.026393	-0.00267	-0.20692	0.030306	0.007871	0.02381	-0.0022	0.006076	0.003837	-0.00391	-0.06897	0.009663	0.006915		
REG_REGH	0.009161	-0.00901	0.01315	0.02562	0.043334	0.027334	-0.00569	-0.05971	-0.03901	-0.02811	-0.02924	0.038215	0.067196	0.001251	0.010515	0.020406	-0.01135	-0.03905	-0.03486	1	0.45157	0.082467	0.339489	0.144701	0.003409	-0.02179	-0.0146	-0.00515	-0.01468	-0.00597	0.031512	-0.01828	-0.00304	-0.00474	8.57E-05	-0.00299	0.001448	-0.02008
REG_REGH	-0.00098	0.010091	0.0262	0.057603	0.080867	0.05844	0.061607	-0.09183	-0.10635	-0.03588	-0.04551	0.105065	0.067907	-0.01218	-0.0021	0.040355	0.005403	-0.14187	-0.45157	1	0.861246	0.155467	0.235913	0.189637	-0.01728	-0.02618	-0.00806	-0.026	-0.01087	0.035396	-0.004464	0.002417	0.002081	0.002455	0.004593	-0.00785	-0.02833	
LIVE_CITY	0.031545	0.069149	0.0101244	0.004863	0.011798	0.003219	-0.01118	-0.15077	-0.21652	-0.06343	-0.05813	0.215376	0.105173	0.001553	-0.02581	0.001386	0.078873	-0.01897	-0.00267	0.003409	0.189637	0.023124	0.026416	0.825196	1	-0.08486	-0.07778	-0.00311	-0.00749	0.000615	0.019078	0.004942	0.005808	-0.00043	0.004346	-0.00513	-0.00689	-0.01054
Avg_EXT	-0.22635</																																					

All variable with datatset for Target 0

	T_CHILDREN	INCOME	TMT_CREDIT_ANNUI	GOODS_PIPULATION	YEARS_BIRTHS_EMPL	REGISTRAS_ID	PUBS_EMP	PHI_WORK	PF_CONT	MCAG_PHONE	FAM_EMAIL	MEMI_RATING	TING_CLIEU	NOT_LIV	NOT_WC1	NOT_WC2	NOT_WC3	NOT_WC4	EXT_SOUINT	SOCIAINT_SOCIAINT	SOCIAINT_SOCIAINT	PHONE_OUMENT	CREDIT_BU	CREDIT_BU	CREDIT_BU	CREDIT_BU												
CNT_CHIL	1																																					
AMT_INCI	0.035346	1																																				
AMT_CREC	0.004527	0.374889	1																																			
AMT_ANN	0.025692	0.448659	0.770864	1																																		
AMT_GOO	0.000188	0.38102	0.98712	0.775993	1																																	
REGION_P	-0.02403	0.180381	0.097833	0.117294	0.101537	1																																
YEARS_BIF	-0.34524	-0.08059	0.04098	-0.01617	0.03926	0.028666	1																															
YEARS_EM	-0.24655	-0.16222	-0.07947	-0.11206	-0.07699	-0.00398	0.622563	1																														
YEARS_REL	-0.18299	-0.07164	-0.01109	-0.03607	-0.01484	0.055478	0.336008	0.207167	1																													
YEARS_ID	0.035247	-0.04268	-0.00205	-0.01604	-0.00078	0.002165	0.247598	0.267339	0.09486	1																												
FLAG_EMP	0.245576	0.162644	0.080624	0.112785	0.078176	0.003944	-0.61851	-0.99971	-0.20485	-0.26652	1																											
FLAG_WOI	0.050545	-0.03524	-0.00589	-0.01657	0.015518	-0.01476	-0.1722	-0.23209	-0.05648	-0.04928	0.232368	1																										
FLAG_CON	0.002145	-0.0207	0.026368	0.026294	0.023833	-0.00591	0.011067	0.016537	-0.00293	0.07952	-0.01657	0.023213	1																									
FLAG_PHO	-0.03551	0.002533	0.018588	0.007102	0.034135	0.094738	0.042455	0.023114	0.071612	0.029041	-0.02244	0.29813	0.003698	1																								
FLAG_EMA	0.02746	0.094201	0.012768	0.067695	0.012033	0.042039	-0.09479	-0.06799	-0.03295	-0.03482	0.067564	-0.01281	-0.01157	0.01164	1																							
CNT_FAM	0.880396	0.037989	0.063817	0.078321	0.061975	-0.02291	-0.29608	-0.2356	-0.1724	0.024613	0.23502	0.064609	0.000882	-0.02205	0.023117	1																						
REGION_R	0.020871	-0.20561	-0.10428	-0.1277	-0.10582	-0.53541	-0.00655	0.041515	-0.07943	0.008477	-0.04128	0.000353	0.014553	-0.08909	-0.06318	0.022247	1																					
REGION_R	0.017709	-0.22199	-0.11472	-0.1431	-0.11548	-0.53236	-0.00422	0.044526	-0.07217	0.012393	-0.04435	0.005843	0.014894	-0.08355	-0.06081	0.021399	0.950026	1																				
REG_REGI	-0.00851	0.081627	0.027491	0.04518	0.029363	-0.00602	-0.0606	-0.03834	-0.02977	-0.02984	0.037535	0.065287	0.001019	0.009249	0.018988	-0.01239	-0.04013	-0.03604	1																			
REG_REGI	0.011577	0.16039	0.059396	0.081917	0.061018	0.064733	-0.09349	-0.1081	-0.03777	-0.04617	0.106787	0.06168	-0.0017	-0.00373	0.041684	0.006007	-0.14527	-0.13788	0.443823	1																		
LIVE_REGI	0.01959	0.151375	0.057632	0.074483	0.057452	0.088528	-0.06901	-0.09644	-0.02571	-0.03405	0.09541	0.039436	-0.00144	-0.00777	0.042019	0.015104	-0.15	-0.14346	0.081361	0.865459	1																	
REG_CITY	0.023454	0.017598	-0.01929	-0.00189	-0.01864	-0.04348	-0.17816	-0.09225	-0.06816	-0.06848	0.090668	0.052961	0.003135	-0.04	0.013936	0.016078	0.031473	0.041009	0.339926	0.151591	0.023438	1																
REG_CITY	0.074853	0.018908	-0.0115	0.003175	-0.01204	-0.03511	-0.23283	-0.25349	-0.09163	-0.09515	0.251628	0.120094	-0.00055	-0.04408	0.006887	0.078387	0.0017	0.021934	0.144323	0.236879	0.184676	0.434711	1															
LIVE_CITY	0.069496	0.021462	0.006719	0.012029	0.005634	-0.00842	-0.1489	-0.21725	-0.06127	-0.05943	0.216148	0.106195	0.000474	-0.02789	0.001267	0.080721	-0.02295	-0.0069	0.00354	0.191522	0.233614	0.02893	0.828197	1														
Avg_EXT_S	-0.06996	0.069279	0.13991	0.11832	0.146997	0.137325	0.281918	0.094204	0.131781	0.115127	-0.09176	-0.05625	0.00908	0.051763	-0.01291	-0.03686	-0.20609	-0.20084	-0.02102	-0.01895	-0.00711	-0.09266	-0.12035	-0.07915	1													
OBS_30_C	0.015511	-0.03077	-2.8E-06	-0.00943	-0.00448	-0.01656	-0.01178	0.006982	-0.01182	0.010189	-0.00695	-0.02007	0.006231	-0.03379	-0.00308	0.02305	0.031709	0.029905	-0.01323	-0.0257	-0.02112	-0.00579	-0.00718	-0.00697	-0.01574	1												
DEF_30_CI	-0.00447	-0.02937	-0.01165	-0.01753	-0.01302	0.011648	0.003205	0.018611	0.000192	-0.00138	-0.01877	-0.01726	0.003628	-0.02659	0.000199	-0.00392	0.004988	0.003358	-0.00804	-0.00993	-0.00791	0.004627	0.001402	-0.00171	-0.03806	0.304685	1											
OBS_60_C	0.015666	-0.03083	0.000284	-0.00927	-0.00026	-0.01553	-0.01179	0.00688	-0.01213	0.010458	-0.00685	-0.02037	0.006167	-0.03335	-0.00265	0.023284	0.031453	0.029535	-0.0133	-0.02581	-0.02126	-0.00603	-0.00708	-0.00672	-0.0157	0.998353	0.307258	1										
DEF_60_CI	-0.00541	-0.02995	-0.01633	-0.02124	-0.01732	0.007188	0.000975	0.017422	-0.00477	-0.00098	-0.01758	-0.01418	0.004904	-0.02406	-0.00054	-0.00574	0.008413	0.006188	-0.00898	-0.01392	-0.01164	0.003113	0.003598	0.001833	-0.04266	0.226273	0.850926	0.228652	1									
DAYS_LAST	-0.00241	-0.04504	-0.06343	-0.06154	-0.05706	-0.04044	-0.05553	0.030701	-0.04789	-0.07488	-0.0325	-0.0487	-0.02797	-0.06412	-0.02133	-0.02139	0.021539	0.021045	0.032091	0.036566	0.027169	0.046981	0.038053	0.015931	-0.17671	-0.00953	0.003492	-0.01055	0.003306	1								
FLAG_DOC	-0.01614	0.04297	0.227758	0.2008	0.198871	-0.01099	0.042816	0.011298	0.004813	0.034885	-0.01082	0.028694	-0.05185	0.038076	-0.00017	-0.00445	-0.00095	-0.00164	-0.01892	-0.00437	0.005711	-0.02339	-0.00999	0.003877	0.007965	0.020392	0.005251	0.020803	0.001547	-0.00361	1							
AMT_REQ	0.002232	0.006923	-0.00165	0.01028	-0.00077	-0.0032	-0.00287	-0.00436	0.004017	-0.00585	0.003537	-0.01005	-0.00234	-0.01141	0.005116	0.003135	0.00861	0.007465	-0.00219	0.000375	0.002859	0.001304	0.00533	0.004693	0.000349	0.002496	-0.00462	0.00274	-0.00329	0.000663	-0.00022	1						
AMT_REQ	0.000814	0.008689	0.012978	0.00933	0.013188	-0.00021	-0.00317	0.001964	0.003569	-0.00604	-0.00203	-0.00641	-0.0168	-0.00071	0.003877	1.5E-05	0.002275	0.00138	-0.00591	0.001123	0.00329	0.000731	0.000343	-0.00103	-0.00203	0.0007	0.004075	0.000588	0.003133	0.001196	0.012761	0.230225	1					
AMT_REQ	0.003469	0.006488	0.002263	0.01907	0.00276	0.003241	0.000121	-0.00594	-0.00054	-0.006036	-0.00228	-0.02254	0.003582	0.024825	0.004483	-0.00111	-0.00508	-0.00088	0.004366	0.006266	0.000433	0.003825	0.003382	-0.00926	-0.00541	-0.00603	0.00586	-0.00222	0.012826	0.010403	0.248195	1						
AMT_REQ	-0.0145	0.07436	0.062465	0.038804	0.064487	0.077484	-0.00174	-0.03355	0.011411	0.004738	0.033867	-0.00642	0.006265	0.045366	0.021048	-0.00783	-0.07019	-0.06769	-0.00765	0.005954	0.011455	-0.01206	-0.01088	-0.00363	0.033064	0.007577	0.00877	0.007547	0.004895	-0.04338	0.014626	0.006538	-0.0032	-0.01699	1			
AMT_REQ	-0.00815	0.007437	0.019833	0.007476	0.02073	-0.00961	0.01704	0.018085	-0.00363	0.012049	-0.01834	-0.02579	0.001747	-0.01437	0.015012	-0.00907	0.012381	0.011067	0.002505	-0.00762	-0.0125	0.004149	-0.00047	-0.00369	-0.01561	0.007892	0.006493	0.007733	0.009944	-0.00295	0.015595	-0.001	-0.01175	-0.0242	-0.00437	1		
AMT_REQ	-0.042	0.023605	-0.004395	-0.004689	0.006171	0.070461	0.050749	0.024679	0.032745	-0.05081	-0.07298	0.026023	-0.0192	0.046646	-0.02974	0.007021	0.004616	-0.01878	-0.0281	-0.02351	-0.00284	-0.00894	-0.01197	-0.03838	0.035296	0.016565	0.035772	0.017587	-0.11525	0.000703	-0.00079	0.00497	0.01495	0.00184	0.098412	1		

All variable with datatset for Target 1

Insights

- Within the application data file, a significant majority of applications, accounting for 90.68%, were registered for cash loans, whereas 9.32% of applications were earmarked for revolving loans. This distribution underscores the prevalent preference for cash loans among applicants.
- The Average External Sources (Avg EXT Sources) emerged as a pivotal feature in assessing loan defaults. Introducing additional features akin to Avg EXT Sources could enhance the model's robustness, enabling a more comprehensive analysis of loan repayment probabilities.
- An intriguing insight from the dataset's 'Name Income Type' column indicates that 17.66% of applicants identify as pensioners. This demographic segment represents a substantial portion of the applicant pool, warranting targeted strategies tailored to their unique financial circumstances.
- In loan borrower analysis, the region rating assumes paramount importance. Regions with lower ratings often face heightened scrutiny during the loan approval process, potentially leading to application rejections. Juxtaposing region rating with metrics like OBS_30_CNT_SOCIAL_CIRCLE and OBS_60_CNT_SOCIAL_CIRCLE could yield valuable insights. Such an approach could help identify applicants associated with groups of defaulters, thereby safeguarding against potential exploitation of financial institution resources.
- An intriguing observation stems from the 'DAYS_LAST_PHONE_CHANGE' column in the dataset. Among defaulted applications, a noteworthy 14.79% of individuals altered their phone numbers on the same day as their application, as indicated by the unique value '0' in this column. This finding suggests a potential correlation between immediate phone number changes and loan default, warranting further investigation into the underlying causes.

Results

- A. The original dataset comprised 122 features. However, a significant portion of these features contained less than 50% of the data. Following a thorough evaluation, 49 columns were retained for in-depth analysis.
- B. Outliers were identified primarily in the "Income Amount" and "Years Employed" features. Notably, an anomaly was detected in the "Years Employed" feature, with instances where the reported employment duration exceeded plausible limits, such as 1001 years, suggesting potential data entry errors.
- C. The feature "FLAG_CONT_MOBILE" exhibited extreme data imbalance, with a minority proportion of only 0.21%.
- D. Univariate, segmented univariate, and bivariate analyses yielded invaluable insights. These analyses involved the creation of new columns containing bin values derived from principal features. This approach facilitated a comprehensive examination of various data attributes and their relationships.
- E. Correlation analysis with the target variable proved instrumental in gauging feature significance. Notably, the "Avg Ext Source" feature demonstrated noteworthy correlation performance, indicating its potential significance in predicting loan default probabilities.

For Excel File Please click [here](#)

For Loom Video Please click [here](#)