

# Bank Loan Case Study

Final Project-2

Yashi Gupta

# Project Description

As a data analyst at a finance company specializing in lending to urban customers, my role is to conduct **Exploratory Data Analysis (EDA)** to understand customer behavior and loan repayment patterns. The company faces a challenge where applicants with **insufficient credit history** sometimes exploit the system, leading to **loan defaults**. This increases financial risk and affects overall profitability. The goal of this project is to analyze loan application data to identify key indicators of default and help the company make better lending decisions.

## What We Are Going to Do

### *1. Data Cleaning & Preparation*

- Load and inspect the dataset to understand its structure.
- Handle missing values, duplicates, and inconsistencies.
- Perform data transformations to ensure usability.

### *2. Exploratory Data Analysis (EDA)*

- Analyze distributions of **customer demographics** (age, income, employment type, etc.).
- Examine **loan characteristics** (amount, interest rates, tenure).
- Compare **payment behavior** across different customer segments.
- Identify patterns in customers who defaulted vs. those who paid on time.
- Explore correlations between **loan approval outcomes** and customer attributes.

### *3. Segmentation of Loan Applicants*

- Classify customers into **high-risk** and **low-risk** categories.

- Evaluate how different loan types and amounts impact repayment behavior.
- Identify key risk factors that contribute to **late payments and defaults**.

#### 4. *Visualization & Insights*

- Use **charts, graphs, and heatmaps** to present findings.
- Highlight trends in **loan defaults, approval rates, and customer profiles**.
- Provide recommendations based on data-driven insights.

#### 5. Business Impact & Decision Support

- Develop strategies to **reduce loan default risk** (e.g., adjusting approval criteria, modifying loan terms, or implementing risk-based pricing).
- Ensure the company does not reject **creditworthy applicants unnecessarily**.
- Assist in optimizing lending policies to balance **profitability and risk management**.

By the end of this analysis, we will have a clearer understanding of **what factors influence loan defaults** and how the company can make more informed decisions to enhance its loan approval process.

## Approach

The dataset contains detailed records of loan applications, including customer demographics, financial background, and loan repayment history. Customers fall into two main categories:

1. **Customers with payment difficulties** – Those who have missed or delayed payments beyond a certain threshold.
2. **Customers without payment issues** – Those who have consistently paid on time.

Each loan application results in one of four outcomes:

- **Approved** – The loan is granted.
- **Cancelled** – The customer withdraws the application during the process.
- **Refused** – The application is rejected.
- **Unused Offer** – The loan is approved, but the customer does not utilize it.

The workflow consists of three main stages: **Data Cleaning, Data Analysis, and Data Visualization.**

## 1. Data Cleaning

Before performing any analysis, it is crucial to clean the dataset to ensure accuracy and reliability.

### Step 1: Identifying and Handling Missing Data

- Used COUNTBLANK(range) to identify missing values in each column.
- Applied IF(ISBLANK(cell), "Missing", "Present") to flag missing data points.
- Dealt with missing values based on data type:
  - **Numerical data:** Replaced missing values with the **median** (MEDIAN(range)) to avoid skewness.
  - **Categorical data:** Filled missing values with the most frequently occurring category (MODE(range)).
- **Visualization:** Created a **bar chart** to show the proportion of missing values across different columns.

### Step 2: Identifying and Handling Outliers

- Used **Interquartile Range (IQR) Method:**

- $Q1 = \text{QUARTILE}(\text{range}, 1)$ ,  $Q3 = \text{QUARTILE}(\text{range}, 3)$ ,  $\text{IQR} = Q3 - Q1$ .
- Defined outliers as values outside  $Q1 - 1.5 * \text{IQR}$  or  $Q3 + 1.5 * \text{IQR}$ .
- Highlighted outliers using **Conditional Formatting** in Excel.
- **Visualization:** Created **box plots** to identify extreme values in loan amounts, applicant income, and credit history.

## 2. Data Analysis

This stage focuses on identifying trends, relationships, and patterns in loan default scenarios.

### Step 3: Analyzing Data Imbalance

- Used  $\text{COUNTIF}(\text{range}, \text{condition})$  to count occurrences of each loan status (Approved, Refused, etc.).
- Calculated the **class distribution percentage** to check for imbalance.
- **Visualization:**
  - **Pie chart** showing the distribution of loan approval statuses.
  - **Bar chart** comparing the proportion of defaulters vs. non-defaulters.

## Step 4: Univariate Analysis

- Examined individual variable distributions using:
  - AVERAGE(range), MEDIAN(range), and STDEV.P(range) for numerical data.
  - **Pivot Tables** to summarize categorical variables.
- **Visualization:**
  - **Histograms** for income, loan amount, and credit history.
  - **Bar charts** for categorical variables like employment type and loan type.

## Step 5: Segmented Univariate Analysis

- Compared variable distributions across different loan scenarios (approved vs. refused applicants).
- Used **Pivot Tables and Sorting** to segment data by customer profile.
- **Visualization:**
  - **Stacked bar charts** for comparing different applicant segments.
  - **Grouped bar charts** to analyze how income levels differ across loan statuses.

## Step 6: Bivariate Analysis

- Explored relationships between two variables (e.g., income vs. default rate, loan amount vs. approval status).
- Used CORREL(range1, range2) to calculate correlation coefficients.
- **Visualization:**
  - **Scatter plots** for identifying trends in numerical data.
  - **Heatmaps (Conditional Formatting)** to highlight high correlations.

## Step 7: Identifying Top Correlations for Loan Default

- Segmented dataset into two groups:
  - **Customers with payment difficulties**
  - **All other cases**
- Used CORREL(range1, range2) to rank top indicators of loan default.
- **Visualization:**
  - **Correlation matrix (Heatmap)** to display strongest predictors of default.
  - **Box plots** showing variations in top correlated features.

### 3. Data Visualization

To effectively communicate findings, multiple Excel visualization techniques were used:

- **Bar Charts & Pie Charts** → For categorical data distributions.
- **Box Plots** → For identifying outliers in income and loan amounts.
- **Histograms** → For numerical variable distributions.
- **Scatter Plots** → For relationships between two numerical variables.
- **Stacked Bar Charts** → For comparing different segments (approved vs. refused).

## Tech-Stack Used

**Microsoft Excel 2022:** Excel was used to perform the data analysis and calculations. The software's vast array of built-in functions allowed me to calculate descriptive statistics, correlation coefficients, and percentiles with ease. I also used its charting capabilities to visualize the relationships between different variables.

# Insights

## Loan Application Analysis

### 1. Loan Approval & Rejection Trends

- Out of 43,344 loan applications, about 38.6% were approved, while 58.9% were rejected due to strict credit policies.
- The most common loan types were cash loans and repair loans, but large loans (like home loans) faced more rejections.
- Consumer loans had the highest approval rate, while loans for paying off other loans had the lowest.

### 2. Who is More Likely to Default?

The data reveals key factors that influence whether a client will default on a loan:

Lower risk (less likely to default):

- Older clients with more experience → They tend to manage finances better.
- Higher education levels → College-educated clients default less than those with only secondary education.
- Female clients → Women default less than men.
- Corporate employees → Safer bet than labor workers.
- Clients from Region 1 → Lowest default rate, safest for banks.
- Older clients take larger loans but default less, making them highly profitable and low-risk.

Higher risk (more likely to default):

- Younger clients with fewer years of experience.
- Lower education levels (secondary or lower).
- Male clients default more than females.

- Labor-class clients.
- Clients from Region 3 → Highest percentage of defaulters. Banks should apply stricter loan policies for this region.

### 3. General Loan Insights

- Loan approval rates depend on purpose – loans for repairs and urgent needs get approved more than home loans.
- Higher-income individuals apply for fewer loans.
- Married individuals take out more loans than unmarried ones.
- Most clients are aged 35-50, but younger clients (26-31) default the most.
- Clients asking for loans above ₹3,500 are more likely to get denied.

### 4. What Should Banks Do?

- ◆ Prioritize older, educated, and corporate clients – they are safer and bring better profits.
- ◆ Be stricter with younger, less educated, and labor-class applicants – they have a higher risk of default.
- ◆ Region-based policies – stricter rules for Region 3 clients, while clients from Region 1 are the safest bet.
- ◆ Encourage responsible borrowing – since younger clients default more, financial literacy programs could help reduce risk.

Older, educated, and corporate clients are the safest bets. Younger, less-educated, and labor-class clients are riskier. Women default less than men. Banks should be stricter with Region 3 applicants and offer more loans to experienced clients.

# CONCLUSION

## PROJECT SUMMARY

In this project, we analyzed two datasets: the **application dataset** and the **previous application dataset**. These not only helped us understand the current status of individuals applying for loans but also provided insights into their previous loan history. The dataset was very large, making it challenging to manipulate and perform queries on, but it provided valuable insights. We learned that **data imbalance** can cause serious problems while analyzing data—for example, since the **labor profession** had a higher representation, it dominated both **default and non-default categories**, affecting analysis accuracy. We also explored **correlations between different columns** and identified key factors influencing loan defaults and approvals. **Older, educated, and corporate clients are the safest bets**, while **younger, less-experienced individuals from Region 3 have a higher default risk**. Women tend to default less than men. Based on these findings, the bank should implement **stricter credit policies for high-risk clients, prioritize profitable, low-risk clients**, and **enhance data management** to improve decision-making and reduce default rates.

## 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.
- **Hint:** Utilize Excel functions like COUNT, ISBLANK, and IF to identify missing data. Consider using functions like AVERAGE or MEDIAN for imputation or other appropriate methods available in Excel.
- **Graph suggestion:** Create a bar chart or column chart to visualize the proportion of missing values for each variable.

1)WE WILL COUNT THE NO OF BLANKS CELLS PER COLUMN USING COUNT OF BLANK AS SHOWN IN THE IMAGE

2)ALL COLUMNS HAVING BLANKS GREATER THAN 30 ARE HIGHLIGHTED IN RED

3)THIS WILL HELP US ELIMINATE MISSING DATA,IF MORE AMOUNT OF CELLS ARE EMPTY AND FILLED WE WILL REMOVE THAT SAID COLUMN.

A	B	C	D	E	F	G	H	I	J	K	L	M	N	
1	CountOfBlank	0	0	0	0	0	0	0	0	1	38	192	0	
2	Blank %	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	0%	
3	Columns	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	NAME_INCOME_TYPE
4	100002	1	Cash loans	M	N	Y		0	202500	406597.5	24700.5	351000	Unaccompanied	Working
5	100003	0	Cash loans	F	N	N	0	270000	1293502.5	35698.5	1129500	Family	State servant	
6	100004	0	Revolving loans	M	Y	Y	0	67500	135000	6750	135000	Unaccompanied	Working	
7	100006	0	Cash loans	F	N	Y	0	135000	312682.5	29686.5	297000	Unaccompanied	Working	
8	100007	0	Cash loans	M	N	Y	0	121500	513000	21865.5	513000	Unaccompanied	Working	
9	100008	0	Cash loans	M	N	Y	0	99000	490495.5	27517.5	454500	Spouse, partner	State servant	
10	100009	0	Cash loans	F	Y	Y	1	171000	1560726	41301	1395000	Unaccompanied	Commercial associate	
11	100010	0	Cash loans	M	Y	Y	0	360000	1530000	42075	1530000	Unaccompanied	State servant	
12	100011	0	Cash loans	F	N	Y	0	112500	1019610	33826.5	913500	Children	Pensioner	
13	100012	0	Revolving loans	N	Y	0	135000	405000	20250	405000	Unaccompanied	Working		
14	100014	0	Cash loans	N	Y	1	112500	652500	21177	652500	Unaccompanied	Working		
15	100015	0	Cash loans	F	N	Y	0	38419.155	148365	10678.5	135000	Children	Pensioner	
16	100016	0	Cash loans	F	N	Y	0	67500	80865	5881.5	67500	Unaccompanied	Working	
17	100017	0	Cash loans	M	Y	N	1	225000	918468	28966.5	697500	Unaccompanied	Working	
18	100018	0	Cash loans	F	N	Y	0	189000	773680.5	32778	679500	Unaccompanied	Working	
19	100019	0	Cash loans	M	Y	Y	0	157500	299772	20160	247500	Family	Working	
20	100020	0	Cash loans	M	N	N	0	108000	509602.5	26149.5	387000	Unaccompanied	Working	
21	100021	0	Revolving loans	F	N	Y	1	81000	270000	13500	270000	Unaccompanied	Working	
22	100022	0	Revolving loans	F	N	Y	0	112500	157500	7875	157500	Other_A	Working	
23	100023	0	Cash loans	F	N	Y	1	90000	544491	17563.5	454500	Unaccompanied	State servant	
24	100024	0	Revolving loans	M	Y	Y	0	135000	427500	21375	427500	Unaccompanied	Working	
25	100025	0	Cash loans	F	Y	Y	1	202500	1132573.5	37561.5	927000	Unaccompanied	Commercial associate	
26	100026	0	Cash loans	F	N	N	1	450000	497520	32521.5	450000	Unaccompanied	Working	
27	100027	0	Cash loans	F	N	Y	0	83250	239850	23850	225000	Unaccompanied	Pensioner	
28	100029	0	Cash loans	M	Y	N	2	135000	247500	12703.5	247500	Unaccompanied	Working	
29	100030	0	Cash loans	F	N	Y	0	90000	225000	11074.5	225000	Unaccompanied	Working	
30	100031	1	Cash loans	F	N	Y	0	112500	979992	27076.5	702000	Unaccompanied	Working	
31	100032	0	Cash loans	M	N	Y	1	12500	327024	23827.5	270000	Family	Working	
32	100033	0	Cash loans	M	Y	Y	0	270000	790830	57676.5	675000	Unaccompanied	State servant	
33	100034	0	Revolving loans	M	N	Y	0	90000	180000	9000	180000	Unaccompanied	Working	
34	100035	0	Cash loans	F	N	Y	0	292500	665892	24592.5	477000	Unaccompanied	Commercial associate	
35	100036	0	Cash loans	F	N	Y	0	112500	512064	25033.5	360000	Family	Working	
36	100037	0	Cash loans	F	N	N	0	90000	199008	20893.5	180000	Unaccompanied	Working	
37	100039	0	Cash loans	M	Y	N	1	360000	733315.5	39069	679500	Unaccompanied	Commercial associate	
38	100040	0	Cash loans	F	N	Y	0	135000	1125000	32895	1125000	Unaccompanied	State servant	

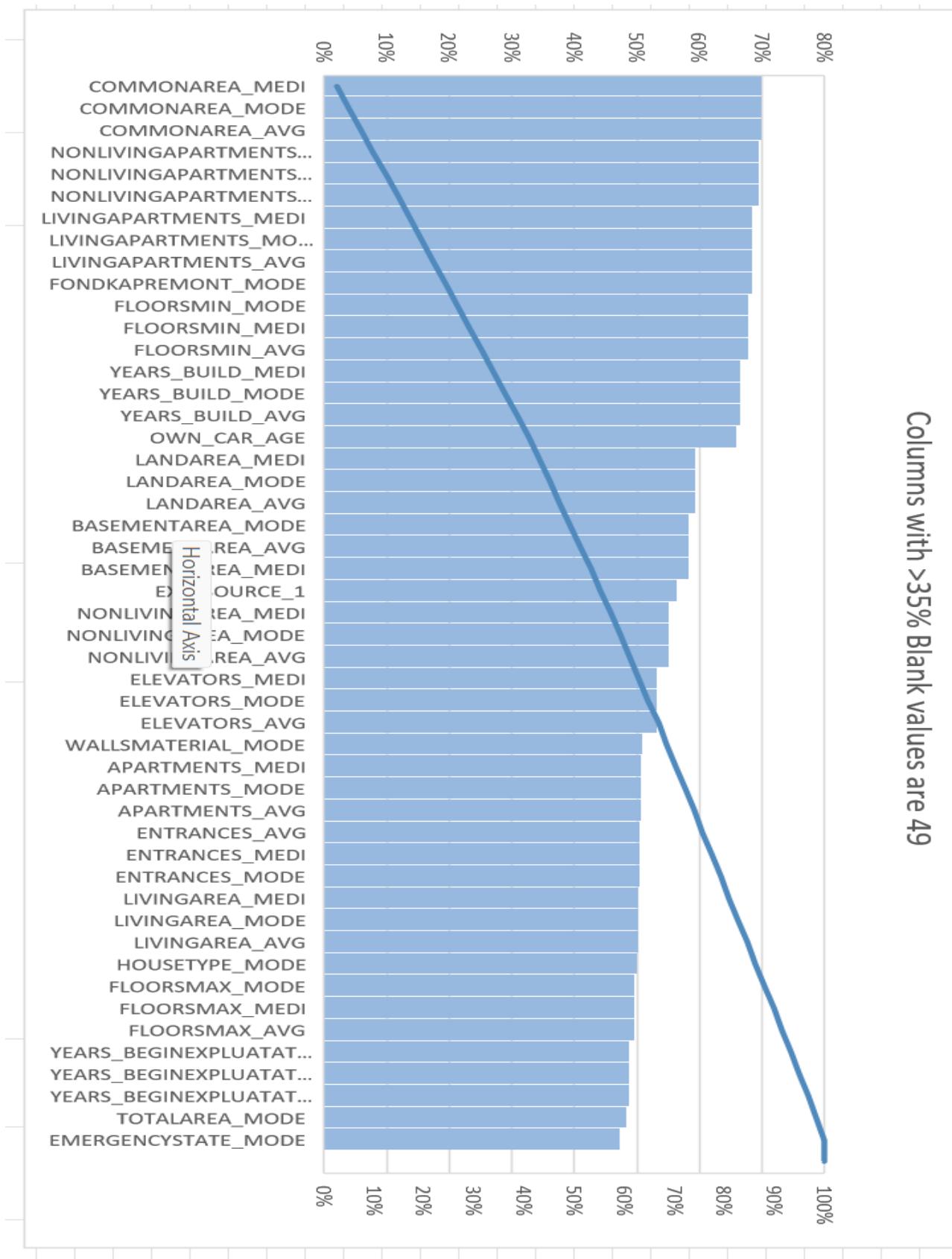
0	0	0	0	0	0	0	0	0	0	0	32950	0	0
NAME_INCOME_TYPE	NAME_EDUCATION_TYPE	NAME_FAMILY_STATUS	NAME_HOUSING_TYPE	REGION_POPULATION_RELATIVE	DAYS_BIRTH	DAYS_EMPLOYED	DAYS_REGISTRATION	DAYS_ID_PUBLISH	OWN_CAR_AGE	FLAG_MOBIL	FLAG_EMP_PHONE	FLAG_DOCUMENT_1	FLAG_DOCUMENT_2
Working	Secondary / secondary special	Single / not married	House / apartment	0.018801	-9461	-637	-3648	-2120	1	1	1	1	1
State servant	Higher education	Married	House / apartment	0.003541	-16765	-1188	-1186	-291	1	1	1	1	1
Working	Secondary / secondary special	Single / not married	House / apartment	0.010032	-19046	-225	-4260	-2531	26	1	1	1	1
Working	Secondary / secondary special	Civil marriage	House / apartment	0.008019	-19005	-3039	-9833	-2437	1	1	1	1	1
Working	Secondary / secondary special	Single / not married	House / apartment	0.028663	-19932	-3038	-4311	-2458	1	1	1	1	1
State servant	Secondary / secondary special	Married	House / apartment	0.035792	-16941	-1588	-4970	-477	1	1	1	1	1
Commercial associate	Higher education	Married	House / apartment	0.035792	-13778	-3130	-1213	-619	17	1	1	1	1
State servant	Higher education	Married	House / apartment	0.003122	-18850	-449	-4597	-2379	8	1	1	1	1
Pensioner	Secondary / secondary special	Married	House / apartment	0.018634	-20999	365243	-7427	-3514	1	0	1	1	1
Working	Secondary / secondary special	Single / not married	House / apartment	0.019689	-14469	-2019	-14437	-3992	1	1	1	1	1
Working	Higher education	Married	House / apartment	0.0228	-10197	-679	-4427	-738	1	1	1	1	1
Pensioner	Secondary / secondary special	Married	House / apartment	0.015221	-20417	365243	-5246	-2512	1	0	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.031329	-13439	-2717	-311	-3227	1	1	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.016612	-14086	-3028	-643	-4911	23	1	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.010006	-14583	-203	-615	-2056	1	1	1	1	1
Working	Secondary / secondary special	Single / not married	Rented apartment	0.020713	-8728	-1157	-3494	-1368	17	1	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.018634	-12931	-1317	-6392	-3866	1	1	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.010966	-9776	-191	-4143	-2427	1	1	1	1	1
Working	Secondary / secondary special	Widow	House / apartment	0.046422	-17718	-7804	-8751	-1259	1	1	1	1	1
State servant	Higher education	Single / not married	House / apartment	0.015221	-11348	-2038	-1021	-3964	1	1	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.015221	-18252	-4286	-298	-1800	7	1	1	1	1
Commercial associate	Secondary / secondary special	Married	House / apartment	0.025164	-14815	-1652	-2299	-2299	14	1	1	1	1
Working	Secondary / secondary special	Married	Rented apartment	0.020713	-11146	-4306	-114	-2518	1	1	1	1	1
Pensioner	Secondary / secondary special	Married	House / apartment	0.006296	-24827	365243	-9012	-3684	1	0	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.026392	-11286	-746	-108	-3729	7	1	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.028663	-19334	-3494	-2419	-2893	1	1	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.018029	-18724	-2628	-6573	-1827	1	1	1	1	1
Working	Secondary / secondary special	Married	House / apartment	0.019101	-15948	-1234	-5782	-3153	1	1	1	1	1
State servant	Higher education	Single / not married	House / apartment	0.046422	-9994	-1796	-4668	-2661	1	1	1	1	1
Working	Higher education	Single / not married	With parents	0.030755	-10341	-1010	-4799	-3015	1	1	1	1	1
Commercial associate	Secondary / secondary special	Civil marriage	House / apartment	0.025164	-15280	-2668	-5266	-3787	1	1	1	1	1
Working	Secondary / secondary special	Civil marriage	House / apartment	0.008575	-11144	-1104	-7846	-2904	1	1	1	1	1
Working	Secondary / secondary special	Civil marriage	House / apartment	0.010032	-12974	-4404	-7123	-4464	1	1	1	1	1
Commercial associate	Secondary / secondary special	Married	House / apartment	0.015221	-11694	-2060	-3557	-3557	3	1	1	1	1
State servant	Higher education	Married	House / apartment	0.019689	-15997	-4585	-5735	-4067	1	1	1	1	1

Columns	Blanks	Blank %	Columns with >35% blanks are 49	
			Column Name	Blank %
SK_ID_CURR	0	0%	EMERGENCYSTATE_MODE	47%
TARGET	0	0%	TOTALAREA_MODE	48%
NAME_CONTRACT_TYPE	0	0%	YEARS_BEGINEXPLUATATION_AVG	49%
CODE_GENDER	0	0%	YEARS_BEGINEXPLUATATION_MODE	49%
FLAG_OWN_CAR	0	0%	YEARS_BEGINEXPLUATATION_MEDI	49%
FLAG_OWN_REALTY	0	0%	FLOORSMAX_AVG	50%
CNT_CHILDREN	0	0%	FLOORSMAX_MODE	50%
AMT_INCOME_TOTAL	0	0%	FLOORSMAX_MEDI	50%
AMT_CREDIT	0	0%	HOUSETYPE_MODE	50%
AMT_ANNUITY	1	0%	LIVINGAREA_AVG	50%
AMT_GOODS_PRICE	38	0%	LIVINGAREA_MODE	50%
NAME_TYPE_SUITE	192	0%	LIVINGAREA_MEDI	50%
NAME_INCOME_TYPE	0	0%	ENTRANCES_AVG	50%
NAME_EDUCATION_TYPE	0	0%	ENTRANCES_MODE	50%
NAME_FAMILY_STATUS	0	0%	ENTRANCES_MEDI	50%
NAME_HOUSING_TYPE	0	0%	APARTMENTS_AVG	51%
REGION_POPULATION_RELATIVE	0	0%	APARTMENTS_MODE	51%
DAYS_BIRTH	0	0%	APARTMENTS_MEDI	51%
DAYS_EMPLOYED	0	0%	WALLSMATERIAL_MODE	51%
DAYS_REGISTRATION	0	0%	ELEVATORS_AVG	53%
DAYS_ID_PUBLISH	0	0%	ELEVATORS_MODE	53%
EMERGENCYSTATE_MODE	23698	47%	ELEVATORS_MEDI	53%
FLAG_MOBIL	0	0%	NONLIVINGAREA_AVG	55%
FLAG_EMP_PHONE	0	0%	NONLIVINGAREA_MODE	55%
FLAG_WORK_PHONE	0	0%	NONLIVINGAREA_MEDI	55%
FLAG_CONT_MOBILE	0	0%	EXT_SOURCE_1	56%
FLAG_PHONE	0	0%	BASEMENTAREA_AVG	58%
FLAG_EMAIL	0	0%	BASEMENTAREA_MODE	58%
OCCUPATION_TYPE	15654	31%	BASEMENTAREA_MEDI	58%
CNT_FAM_MEMBERS	1	0%	LANDAREA_AVG	59%
REGION_RATING_CLIENT	0	0%	LANDAREA_MODE	59%
REGION_RATING_CLIENT_W_CITY	0	0%	LANDAREA_MEDI	59%
WEEKDAY_APPR_PROCESS_START	0	0%	OWN_CAR_AGE	66%
HOUR_APPR_PROCESS_START	0	0%	YEARS_BUILD_AVG	66%
REG_REGION_NOT_LIVE_REGION	0	0%	YEARS_BUILD_MODE	66%
REG_REGION_NOT_WORK_REGION	0	0%	YEARS_BUILD_MEDI	66%
LIVE_REGION_NOT_WORK_REGION	0	0%	MOBILITY_AVG	66%

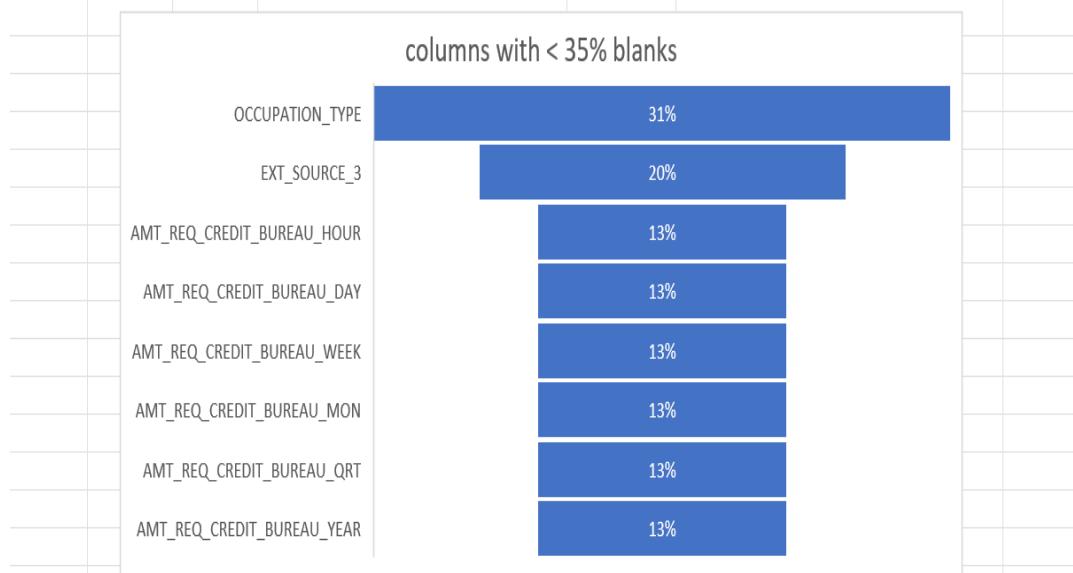
< > ... application\_data | columns\_description | EDA\_APP | [Proportion of Missing Values](#) | [Handling Missing Values](#) | [Model](#)

Adv. [Accessibility](#) [Investigate](#)

Columns with >35% Blank values are 49



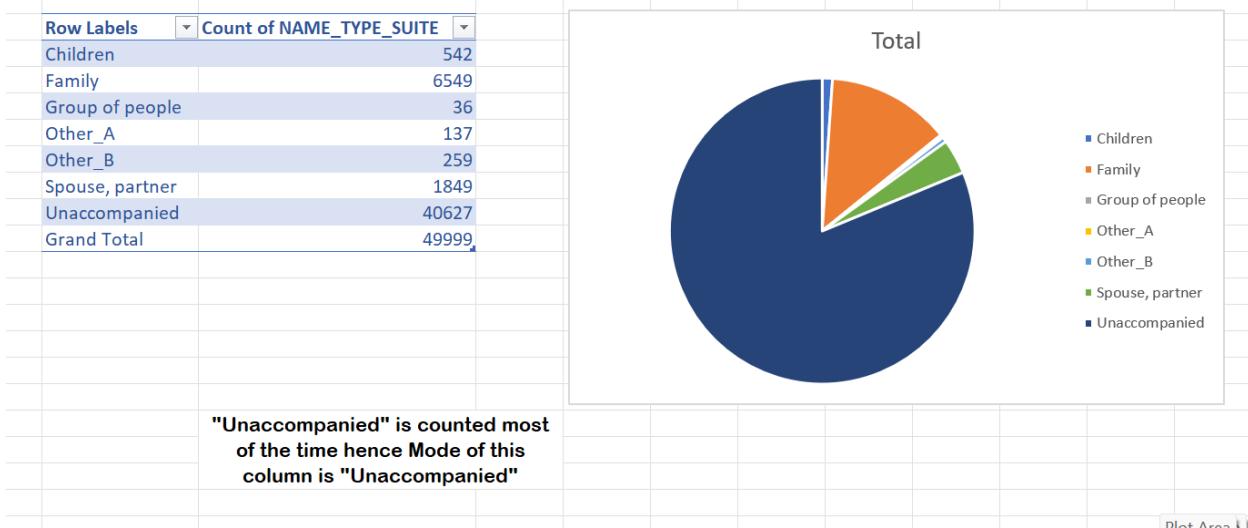
H	I	J	K	L	M	N
Columns with < 35% blank values are 8						
			Column Name	Blank %		
			OCCUPATION_TYPE	31%		
			EXT_SOURCE_3	20%		
			AMT_REQ_CREDIT_BUREAU_HOUR	13%		
			AMT_REQ_CREDIT_BUREAU_DAY	13%		
			AMT_REQ_CREDIT_BUREAU_WEEK	13%		
			AMT_REQ_CREDIT_BUREAU_MON	13%		
			AMT_REQ_CREDIT_BUREAU_QRT	13%		
			AMT_REQ_CREDIT_BUREAU_YEAR	13%		

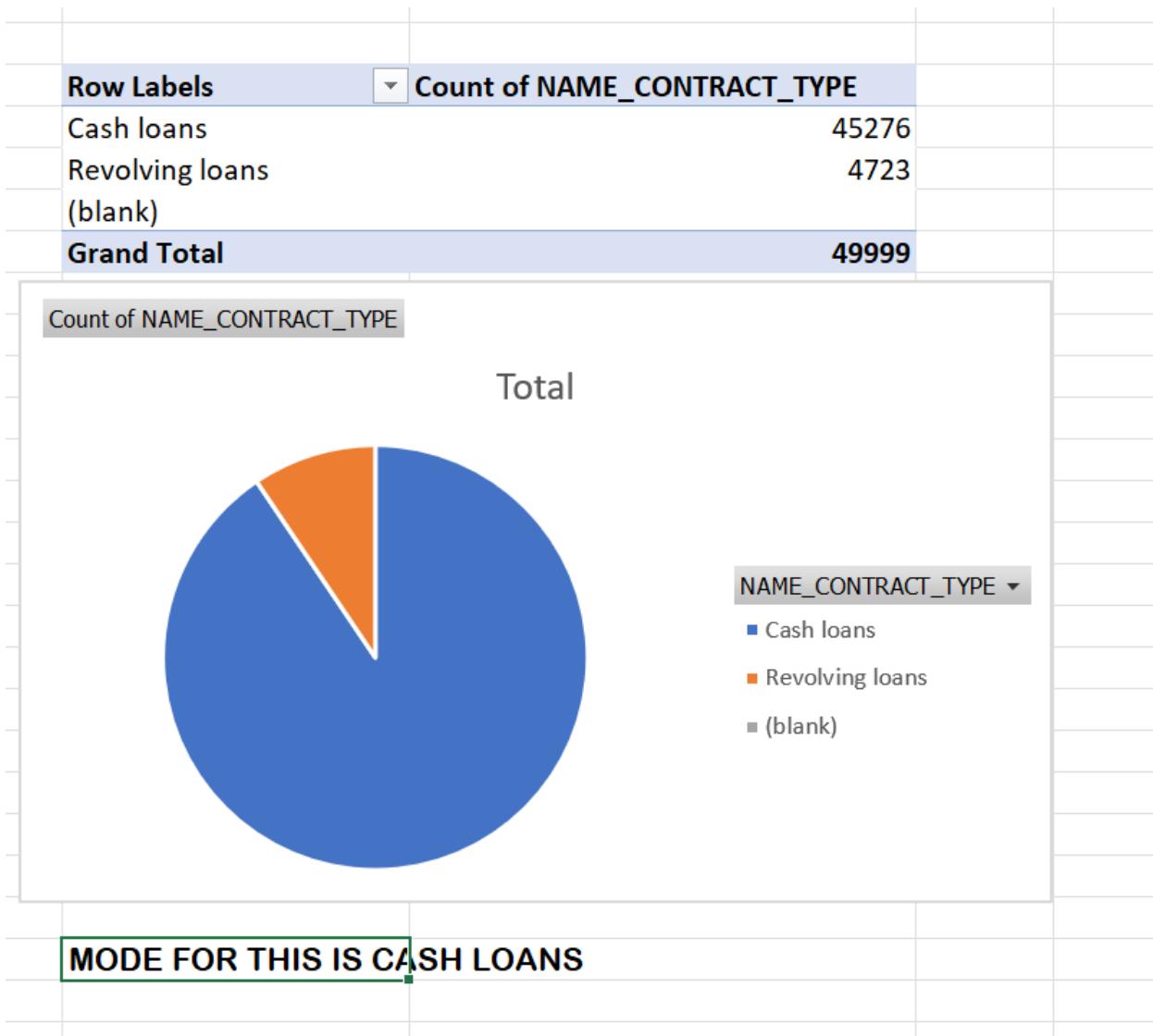


## HANDLING MISSING VALUES

- 1) WE FIND MEAN, MEDIAN AND MODE FOR THE COLUMNS AND FILL THE BLANK CELLS WITH THE RESPECTIVE VALUES
- 2) IF NUMEICAL VALUE- MEDIAN
- 3) IF CATEGORICAL VALUE – MODE

MEDIAN	127652	0	#NUM!	#NUM!	#NUM!	#NUM!	0	144000	513531	24939	450000	#NUM!	#NUM!	#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	NAME_INCOME_TYPE	NAME_EDUCATION_TYPE	NAME_FAMILY_STATUS
100002	1	Cash loans	M	N	Y	0	202500	406597.5	24700.5	351000	Unaccompanied	Working	Secondary / secondary spec	Married
100003	0	Cash loans	F	N	N	0	270000	1293502.5	35698.5	1129500	Family	State servant	Higher education	Married
100004	0	Revolving loans	M	Y	Y	0	6750	135000	6750	135000	Unaccompanied	Working	Secondary / secondary spec	Married
100006	0	Cash loans	F	N	Y	0	135000	312682.5	29686.5	297000	Unaccompanied	Working	Secondary / secondary spec	Married
100007	0	Cash loans	M	N	Y	0	121500	513000	21865.5	513000	Unaccompanied	Working	Secondary / secondary spec	Married
100008	0	Cash loans	M	N	Y	0	99000	490495.5	27517.5	454500	Spouse, partner	State servant	Secondary / secondary spec	Married
100009	0	Cash loans	F	Y	Y	1	171000	1560726	41301	1395000	Unaccompanied	Commercial associate	Higher education	Married
100010	0	Cash loans	M	Y	Y	0	360000	1530000	42075	1530000	Unaccompanied	State servant	Higher education	Married
100011	0	Cash loans	F	N	Y	0	112500	1019610	33826.5	913500	Children	Pensioner	Secondary / secondary spec	Married
100012	0	Revolving loans	M	N	Y	0	135000	405000	20250	405000	Unaccompanied	Working	Secondary / secondary spec	Married
100014	0	Cash loans	F	N	Y	1	112500	652500	21177	652500	Unaccompanied	Working	Higher education	Married





**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.
- **Hint:** Utilize Excel functions like QUARTILE, IQR, and conditional formatting to identify potential outliers. Consider applying thresholds or business rules to determine if the outliers are valid data points or require further investigation.
- **Graph suggestion:** Create box plots or scatter plots to visualize the distribution of numerical variables and highlight the outliers.

THE FOLLOWING STEPS WERE EXECUTED TO PERFORM THE GIVEN TASK-

Step 1: Calculate Quartiles and IQR

1. Calculate Q1 (25th percentile) using the formula:

=QUARTILE.INC(B2:B1000, 1)

2. Calculate Q3 (75th percentile) using the formula:

=QUARTILE.INC(B2:B1000, 3)

3. Calculate the Interquartile Range (IQR):

=Q3 - Q1

4. Calculate the Lower Bound:

=Q1 - 1.5 \* IQR

5. Calculate the Upper Bound:

=Q3 + 1.5 \* IQR

## Step 2: Flag Outliers

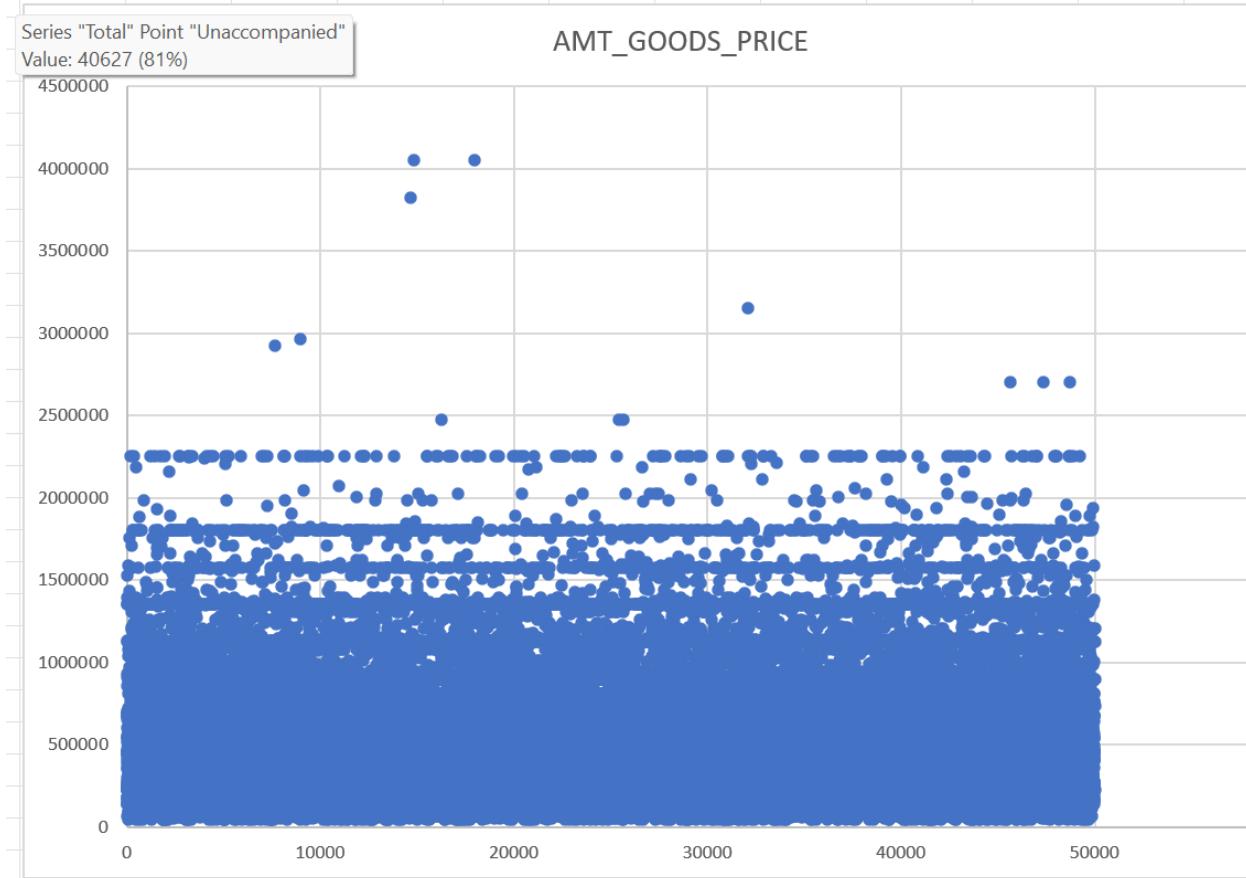
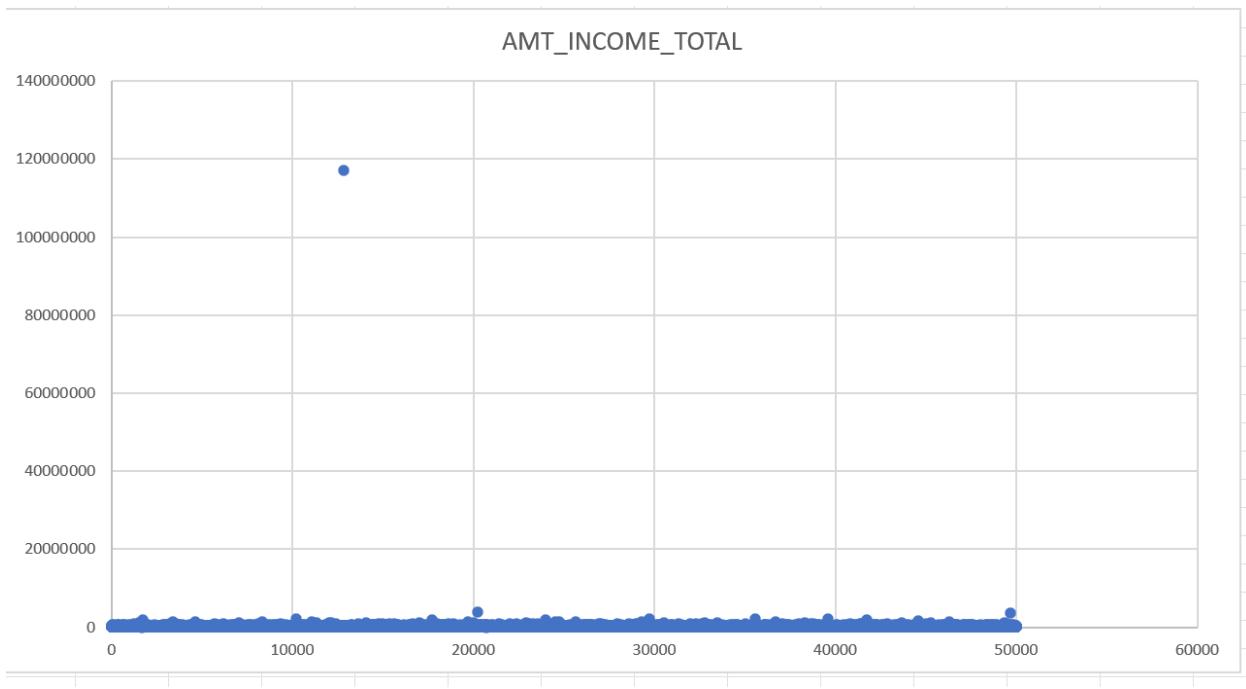
### 1. Apply Conditional Formatting:

- Select the data range (e.g., B2:B1000).
- Go to Home → Conditional Formatting → New Rule.
- Choose “Use a formula to determine which cells to format.”
- Enter the formula: `=OR(B2 < Lower_Bound, B2 > Upper_Bound)`
- Click Format, choose a color, and press OK.

The screenshot shows the Microsoft Excel ribbon with the 'Home' tab selected. A table is open with columns: AMT\_INCOME\_TOTAL, AMT\_CREDIT, AMT\_ANNUITY, AMT\_GOODS\_PRICE, REGION\_POP, YEARS\_REGISTRATION, and AMT\_INCOME\_TOTAL. The YEARS\_REGISTRATION column is currently selected. A 'Conditional Formatting' dialog box is open, showing the formula `=OR(B2 < Lower_Bound, B2 > Upper_Bound)` in the 'Format values where this formula is true' field. The 'Format' button is highlighted. The background shows the table data with various numerical values.

AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE	YEARS_BIRTH	YEARS_EMPLOYED	YEARS_REGISTRATION
202500	406597.5	24700.5	351000	0.018801	25.9	1.7	10.0
270000	1293502.5	35698.5	1129500	0.003541	45.9	3.3	3.2
67500	135000	6750	135000	0.010032	52.2	0.6	11.7
135000	312682.5	29686.5	297000	0.008019	52.1	8.3	26.9
121500	513000	21865.5	513000	0.028663	54.6	8.3	11.8
99000	490495.5	27517.5	454500	0.035792	46.4	4.4	13.6
171000	1560726	41301	1395000	0.035792	37.7	8.6	3.3
360000	1530000	42075	1530000	0.003122	51.6	1.2	12.6
112500	1019610	33826.5	913500	0.018634	55.1	1000.7	20.3
135000	405000	20250	405000	0.019689	39.6	5.5	39.6
112500	652500	21177	652500	0.0228	27.9	1.9	12.1
38419.155	148365	10678.5	135000	0.015221	55.9	1000.7	14.4
67500	80865	5881.5	67500	0.031329	36.8	7.4	0.9
225000	918468	28966.5	697500	0.016612	38.6	8.3	1.8
189000	773680.5	32778	679500	0.010006	40.0	0.6	1.7
157500	299772	20160	247500	0.020713	23.9	3.2	9.6
108000	509602.5	26149.5	387000	0.018634	35.4	3.6	17.5
81000	270000	13500	270000	0.010966	26.8	0.5	11.4
112500	157500	7875	157500	0.04622	48.5	21.4	24.0
90000	544491	17563.5	454500	0.015221	31.1	5.6	2.8
135000	427500	21375	427500	0.015221	50.0	11.7	0.8
202500	1132573.5	37561.5	927000	0.025164	40.6	4.5	6.3
450000	497520	32521.5	450000	0.020713	30.5	11.8	0.3
83250	239850	23850	225000	0.006296	68.0	1000.7	24.7
135000	247500	12703.5	247500	0.026392	30.9	2.0	0.3
90000	225000	11074.5	225000	0.028663	53.0	9.6	6.6
112500	979992	27076.5	702000	0.018029	51.3	7.2	18.0
112500	327024	23827.5	270000	0.019101	43.7	3.4	15.8
270000	790830	57676.5	675000	0.04622	27.4	4.9	12.8
90000	180000	9000	180000	0.030755	28.3	2.8	13.1
292500	665892	24592.5	477000	0.025164	41.9	7.3	14.4
112500	512064	25033.5	360000	0.008575	30.5	3.0	21.5
90000	199008	20893.5	180000	0.010032	35.5	12.1	19.5
360000	733315.5	39069	679500	0.015221	32.0	5.6	9.7
135000	1125000	32895	1125000	0.019689	43.8	12.6	15.7
112500	450000	44509.5	450000	0.008575	33.3	3.5	17.2
198000	641173.5	23157	553500	0.01885	47.1	2.1	0.2

Columns	AMT_INCOME_TOTAL	AMT_CREDI	AMT_ANNUITY	AMT_GOODS_PRIC	REGION_POPULATION_RELATIV	YEARS_BIRTH	YEARS_EMPLOYE	YEARS_REGISTRATIO
Q1	112500	270000	16456.5	238500	0.010006	33.91232877	2.556164384	5.473972603
Q3	202500	808650	34596	679500	0.028663	53.81917808	15.66575342	20.44931507
IQR	90000	538650	18139.5	441000	0.018657	19.90684932	13.10958904	14.97534247
Upper Limit	337500	1616625	61805.25	1341000	0.0566485	83.67945205	35.33013699	42.91232877
Lower Limit	-22500	-537975	-10752.75	-423000	-0.017975	4.052054795	-17.10821918	-16.9890411





## VISUAL REPRESENTATION OF 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.
- **Hint:** Utilize Excel functions like COUNTIF and SUM to calculate the proportions of each class. Compare the class frequencies to assess data imbalance.
- **Graph suggestion:** Create a pie chart or bar chart to visualize the distribution of the target variable and highlight the class imbalance.

Class imbalance occurs when there are more observations in one class than in another. This is a common problem in machine learning.

### Explanation

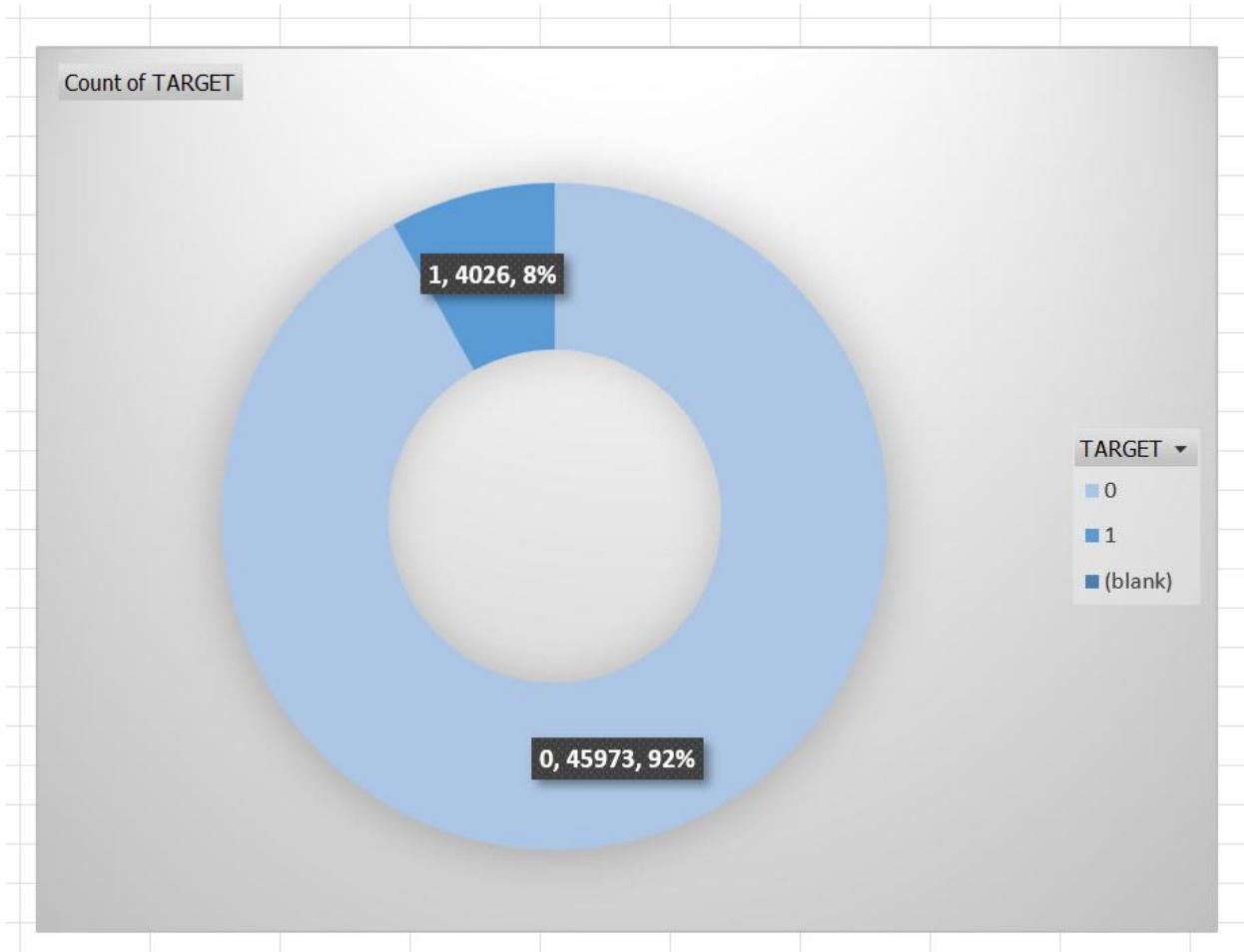
The class with more observations is called the majority class.

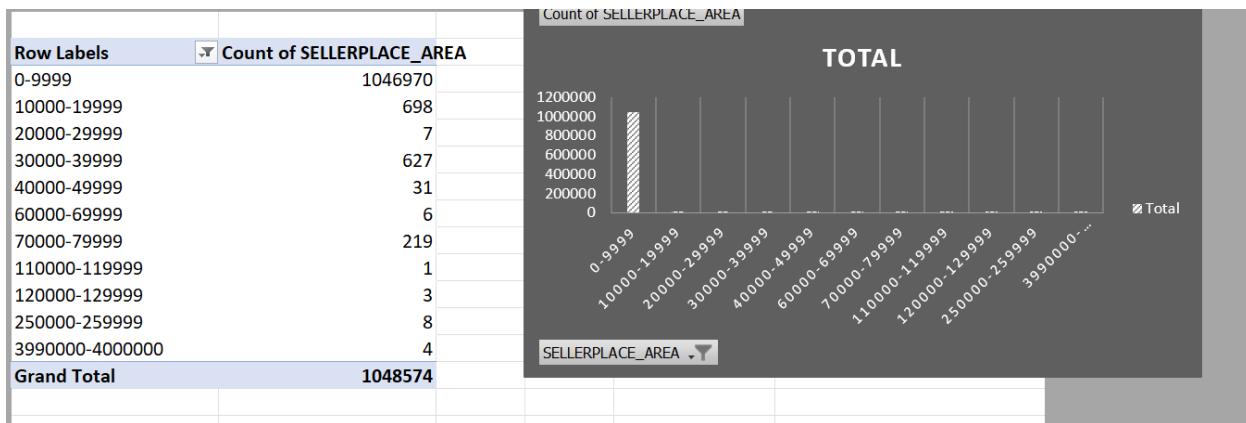
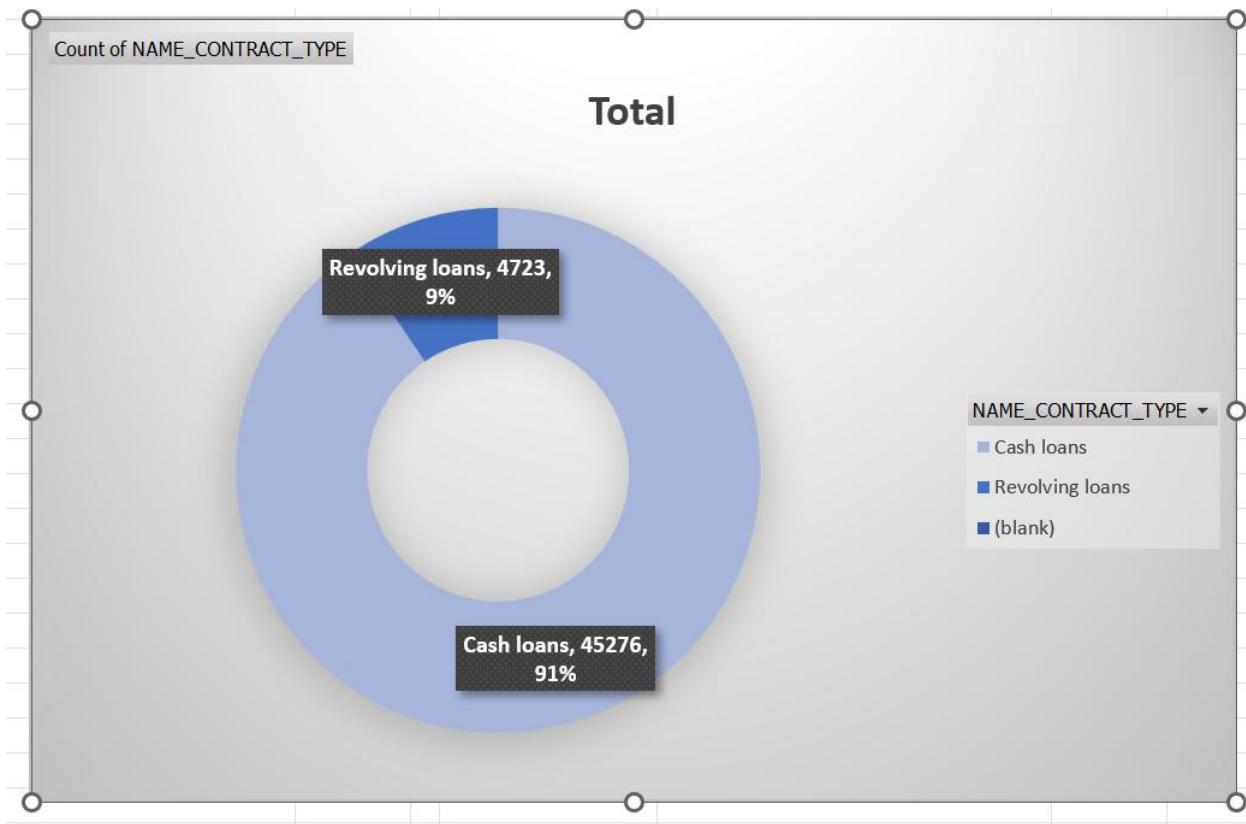
The class with fewer observations is called the minority class.

In a balanced dataset, the number of positive and negative labels is roughly equal.

Class imbalance can affect model accuracy.

We have calculated data imbalance and its ratio using pivot tables on columns then visualizing then using pie chart and bar graphs.





Row Labels	Count of AMT_APPLICATION	
0-500000	957615	
500000-1000000	58212	
1000000-1500000	25535	
1500000-2000000	4314	
2000000-2500000	2395	
2500000-3000000	211	
3000000-3500000	270	
3500000-4000000	12	
4000000-4500000	9	
6500000-7000000	1	
<b>Grand Total</b>	<b>1048574</b>	

**TOTAL**

AMT\_APPLICATION

Row Labels	Count of FLAG_LAST_APPL_PER_CONTRACT	
N	5372	
Y	1043202	
<b>Grand Total</b>	<b>1048574</b>	

**TOTAL**

FLAG\_LAST\_APPL\_PER\_CONTRACT

Row Labels	Count of AMT_ANNUITY	
0-50000	786610	
50000-100000	27254	
100000-150000	1493	
150000-200000	141	
200000-250000	54	
250000-300000	7	
350000-400000	2	
400000-450000	4	
<b>Grand Total</b>	<b>815565</b>	

**TOTAL**

AMT\_ANNUITY

Row Labels	Count of AMT_CREDIT	
0-500000	935511	
500000-1000000	73644	
1000000-1500000	28609	
1500000-2000000	6851	
2000000-2500000	2909	
2500000-3000000	677	
3000000-3500000	235	
3500000-4000000	126	
4000000-4500000	11	
6500000-7000000	1	
<b>Grand Total</b>	<b>1048574</b>	

**TOTAL**

AMT\_CREDIT

**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.
- **Hint:** Utilize Excel functions like COUNT, AVERAGE, MEDIAN, and statistical functions for descriptive analysis. Utilize Excel features like filters, sorting, and pivot tables for segmented and bivariate analysis.
- **Graph suggestion:** Create histograms, bar charts, or box plots to visualize the distributions of variables. Create stacked bar charts or grouped bar charts to compare variable distributions across different scenarios. Create scatter plots or heatmaps to visualize the relationships between variables and the target variable.

In data analysis, "univariate" means analyzing a single variable at a time, "segmented univariate" refers to analyzing a single variable but broken down into different groups or segments within the data, and "bivariate" means examining the relationship between two variables simultaneously.

## UNIVARIATE ANALYSIS

WE ARE SUPPOSE TO CALCULATE MEAN MEDIAN STD DEVIATION MIN AND MAX VALUES FOR THE COLUMNS

1. **Objective:** Summarize each variable's statistics.
2. **Formula:** Use Excel functions like *AVERAGE*, *STDEV*, *MIN*, *MAX*, and *COUNT*.

	A	B	C	D	E
	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	
MEAN	170767.5905	599700.5815	27107.33399	538992.3491	
MEDIAN	145800	514777.5	24939	450000	
STD DEV	531813.7768	402411.4096	14562.6564	369717.1252	
MIN	25650	45000	2052	45000	
MAX	117000000	4050000	258025.5	4050000	
COUNT	49999	49999	49999	49999	
0	427500	814041	23800.5	679500	
1	99000	808650	26217	675000	
2	90000	545040	26640	450000	
3	225000	720000	21051	720000	
4	67500	254700	14220	225000	
5	90000	592560	25105.5	450000	
6	54000	71955	7245	67500	
7	157500	225000	26703	225000	
8	81000	539100	27652.5	450000	
	58500	675000	21775.5	675000	

## Segmented Univariate Analysis

### What We Do in Segmented Univariate Analysis

- **Group by a Category:**

We divide the data into segments based on a categorical variable. In your case, the TARGET variable (with values 0 for non-defaults and 1 for defaults) is used to segment the data.

- **Analyze One Variable at a Time:**

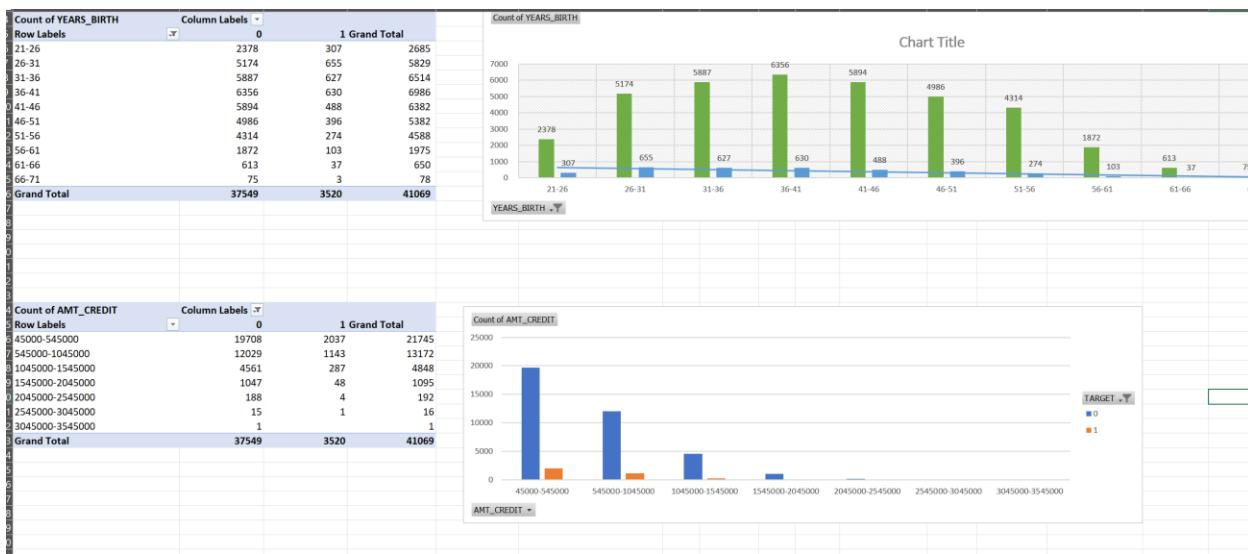
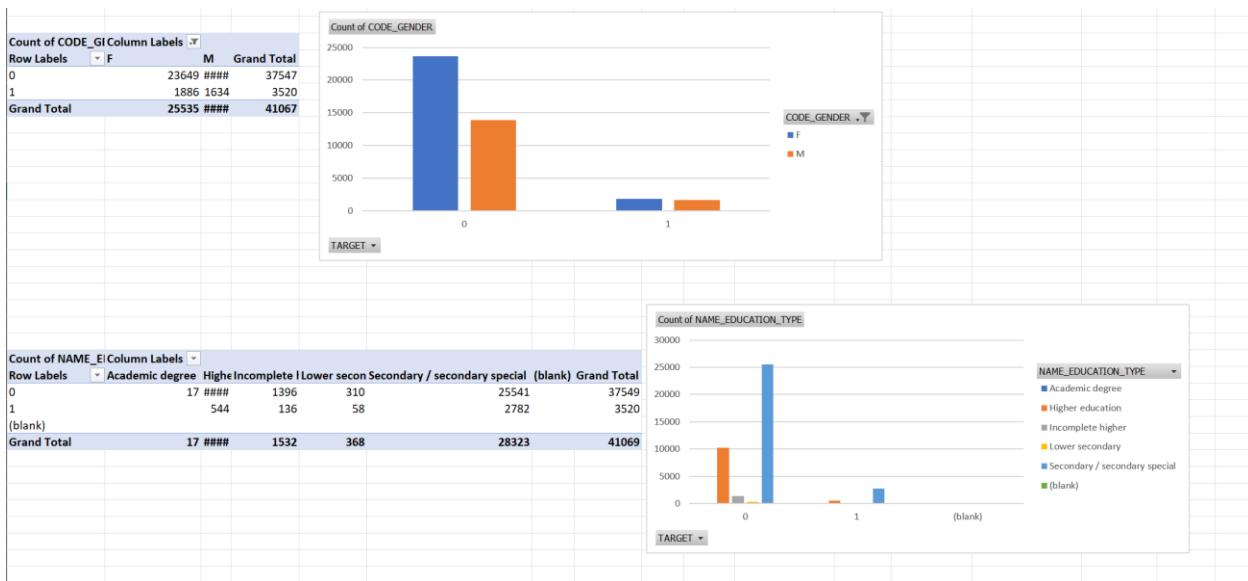
For each segment, we examine individual variables—such as

**AMT\_INCOME\_TOTAL, AMT\_CREDIT, AMT\_ANNUITY, and**

**YEARS\_BIRTH**—to understand their distribution, central tendency (mean, median), variability (standard deviation), and range (min, max).

- **Compare Statistics Across Groups:**

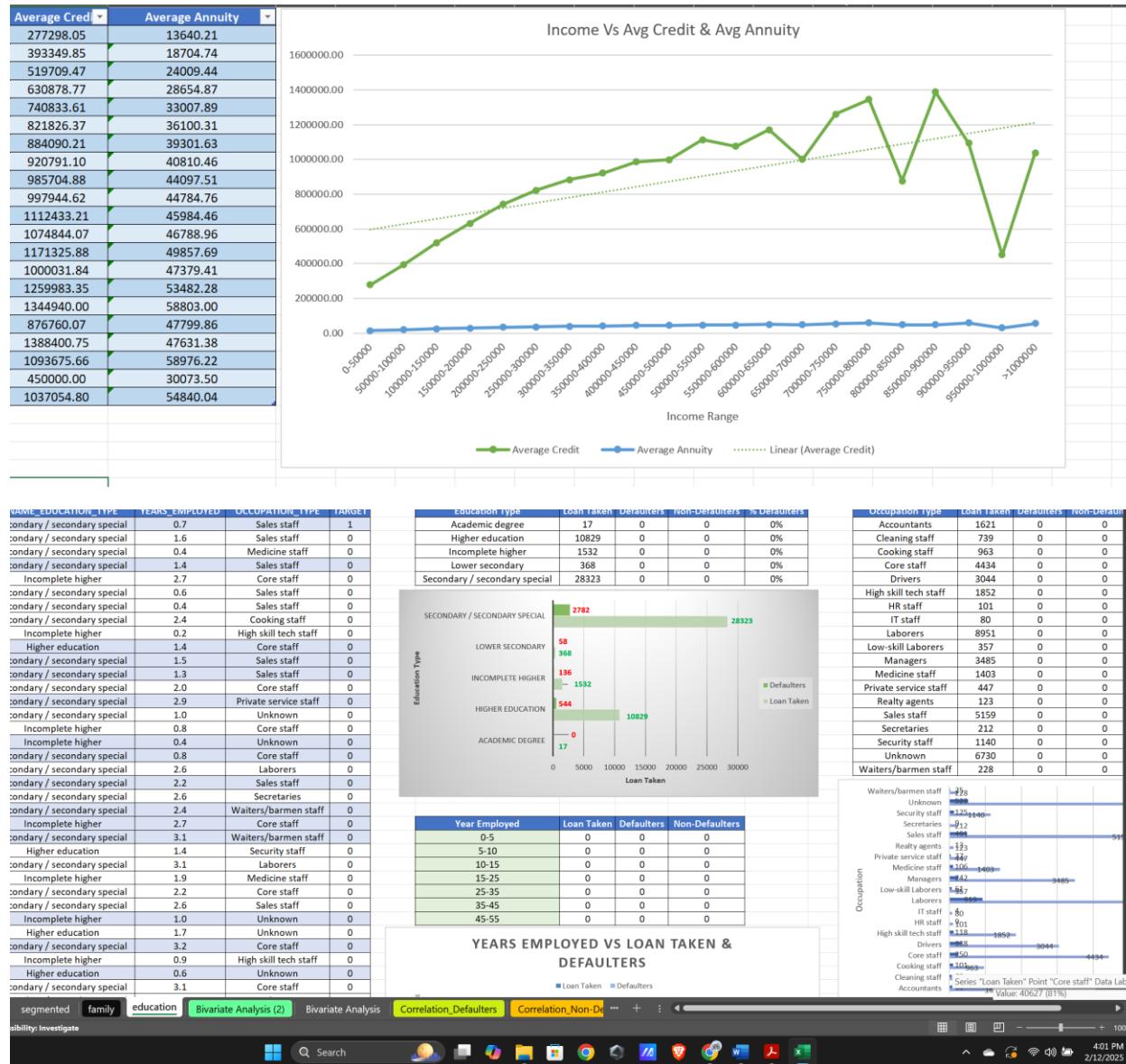
By looking at the descriptive statistics for each segment separately, you can compare how a variable behaves across different groups. For example, you might compare the average income between those who defaulted and those who did not.

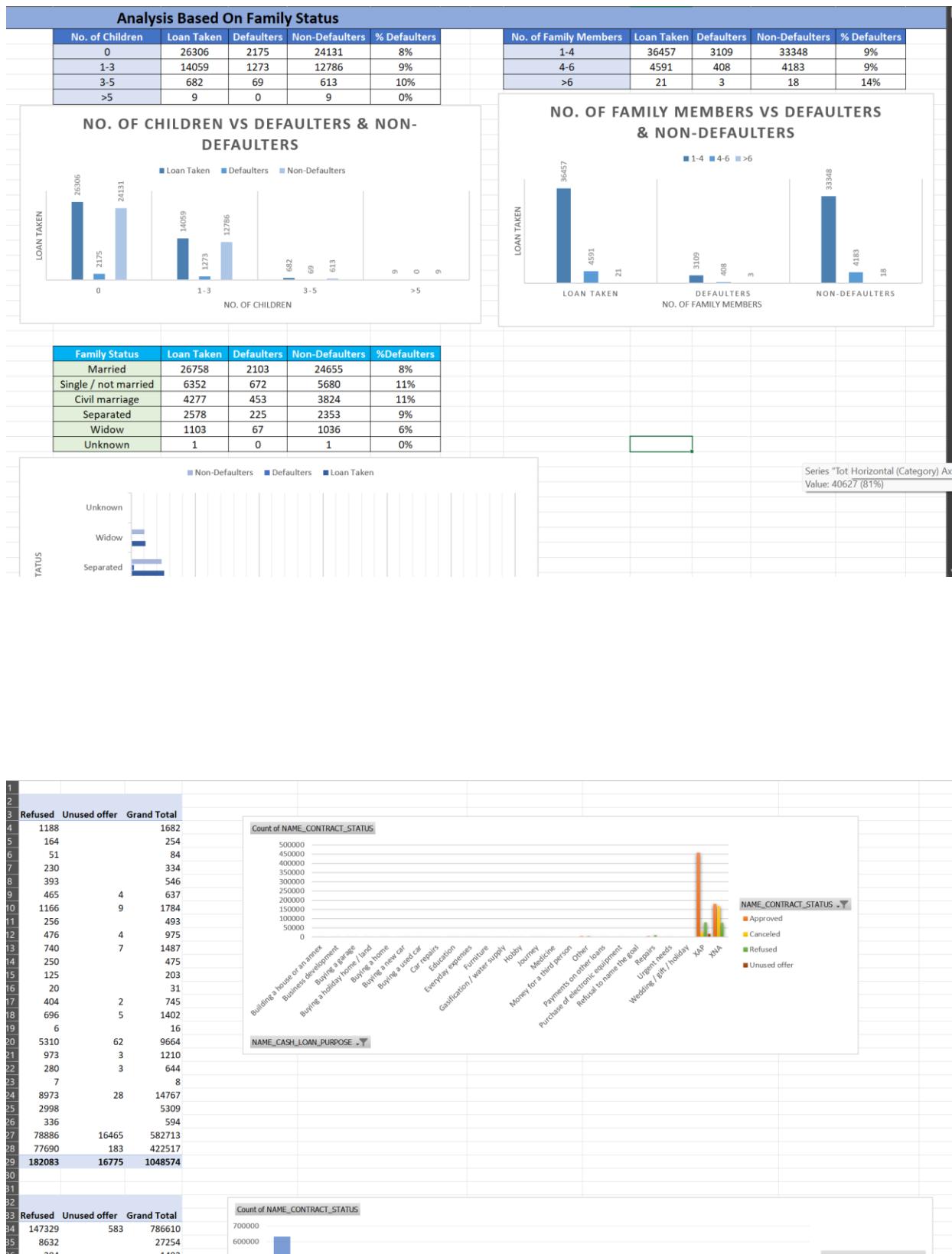


We have compare multiple columns with target column to check and compare payment difficulties.

## Bivariate analysis

Bivariate analysis helps us understand the relationship between two variables. Since our dataset involves loan application data, this analysis will help us explore relationships between numeric and categorical variables to identify trends or risk factors.





**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.
- **Hint:** Utilize Excel functions like CORREL to calculate correlation coefficients between variables and the target variable within each segment. Rank the correlations to identify the top indicators of loan default for each scenario.
- **Graph suggestion:** Create correlation matrices or heatmaps to visualize the correlations between variables within each segment. Highlight the top correlated variables for each scenario using different colors or shading.

The CORREL function in Excel is used to calculate the correlation coefficient between two data sets. The correlation coefficient measures the strength and direction of the linear relationship between two variables.

Formula-

```
=CORREL(array1, array2)
```

## Steps-

1. Separate the data based on defaulters and no defaulters that is 1 and 0 respectively.

2. Then use =CORREL(array1, array2) function to get the correlation between different columns in the data set for both defaulters and non defaulters.

	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE	YEARS_BIRTH	YEARS_EMPLOYED	YEARS_REGISTRATION	YEARS_ID_PUBLISH	FLAG_MOBIL	FLAG_EMP_PHONE
1		1									
2	AMT_INCOME_TOTAL	0.312173644	0.745132112								
3	AMT_CREDIT	0.022601082	0.04109669	0.054426768							
4	AMT_ANNUITY	0.022601082	0.04109669	0.054426768							
5	AMT_GOODS_PRICE	0.313726831	0.981928143	0.746422447	1						
6	REGION_POPULATION_RELATIVE	0.096758897	0.055597704	0.06586731	0.061151451	1					
7	YEARS_BIRTH	0.087629898	0.194437334	0.086228175	0.18810843	0.013409073	1				
8	YEARS_EMPLOYED	0.022601082	0.04109669	0.054426768	0.113070145	-0.001640893	0.305741728	1			
9	YEARS_REGISTRATION	-0.00293807	0.042506313	-0.021558946	0.041100465	0.046966518	0.23986934	0.1506959			
10	YEARS_ID_PUBLISH	0.037532601	0.054409392	0.050271371	0.058966989	0.008905666	0.125405421	0.099252606	0.043760239	1	
11	FLAG_MOBIL	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	
12	FLAG_EMP_PHONE	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	
13	FLAG_WORK_PHONE	-0.099296959	-0.05398719	-0.056408023	-0.026160664	-0.027465876	-0.059625901	0.010872911	-0.021512932	0.011798852	#DIV/0!
14	FLAG_CONT_MOBILE	-0.055294138	0.030523104	0.027004036	0.020027334	0.025658022	0.008828731	0.001112875	0.006168352	0.020707896	#DIV/0!
15	FLAG_PHONE	0.007844835	0.031813347	0.0074097	0.046284121	0.062164434	0.035470603	0.058136973	0.039390682	0.021780211	#DIV/0!
16	FLAG_EMAIL	0.079615096	0.050626381	0.098213289	-0.04526882	0.052260116	-0.05972844	0.032831631	-0.018327464	0.021620140	#DIV/0!
17	CNT_FAM_MEMBERS	-0.02862832	0.048692162	0.04693887	0.044282726	0.016834882	-0.10836037	0.011257766	-0.13215779	0.111308835	#DIV/0!
18	REGION_RATING_CLIENT_W_CITY	-0.160225589	-0.04618178	-0.046652059	-0.038562059	-0.036699303	-0.051303568	-0.036513733	-0.126769392	-0.28399284	#DIV/0!
19	REGION_RATING_CLIENT	-0.168885987	0.059358113	-0.085320242	-0.06228612	-0.439771312	-0.047188741	0.00028359	-0.15800314	-0.19625809	#DIV/0!
20	HOUR_APPR_PROCESS_START	0.060162205	0.047723296	0.038563626	0.060384921	0.157891399	-0.037057941	0.016151137	0.065340507	0.002561693	#DIV/0!
21	REG_REGION_NOT_LIVE_REGION	0.053977901	0.006100489	0.029105568	0.007101455	-1.773456-06	0.025388406	-0.039793874	-0.001284827	-0.01881892	#DIV/0!
22	REG_REGION_NOT_WORK_REGION	0.103557050	0.03889956	0.062818511	0.029248955	0.021144356	0.032366508	-0.070872835	-4.255866-05	-0.023017339	#DIV/0!
23	LIVE_REGION_NOT_WORK_REGION	0.10570591	0.03889956	0.072919945	0.04182111	0.064372174	0.051085566	-0.049838788	0.000373368	0.013720931	#DIV/0!
24	REG_CITY_NOT_LIVE_CITY	-0.015606392	0.050684801	-0.026168688	-0.051317501	0.035041409	0.125601404	-0.096813278	-0.032794847	-0.046355785	#DIV/0!
25	REG_CITY_NOT_WORK_CITY	-0.032745594	0.039205058	-0.01935929	-0.036992095	0.045820766	0.107686126	-0.130516219	-0.062359874	-0.027848522	#DIV/0!
26	LIVE_CITY_NOT_WORK_CITY	-0.018222562	-0.040987208	-0.009392877	-0.010328456	-0.052860959	-0.033905166	-0.070086942	-0.036728491	-0.009833436	#DIV/0!
27	EXT_SOURCE_2	0.116520912	0.113463669	0.103480509	0.127359953	0.157475877	0.162842312	0.10075063	0.084750692	0.047820366	#DIV/0!
28	EXT_SOURCE_3	-0.058902187	0.04753462	0.021729952	0.0502147563	-0.028018518	0.11561263	0.076331438	0.02953675	0.057080185	#DIV/0!
29	OBS_30_CNT_SOCIAL_CIRCLE	0.020308899	0.031064636	0.017327962	0.028834846	0.012352115	0.006111492	0.044734843	0.010210586	0.027148224	#DIV/0!
30	DEF_30_CNT_SOCIAL_CIRCLE	0.045671551	0.036799595	0.0305075618	-0.031456725	0.018822572	0.01408454	-0.002967611	0.000934258	0.024895657	#DIV/0!
31	OBS_60_CNT_SOCIAL_CIRCLE	0.018606977	0.03199129	0.031550337	0.029962662	-0.018322232	0.007623365	0.045377158	0.003967488	0.02578489	#DIV/0!
32	DEF_60_CNT_SOCIAL_CIRCLE	0.033290582	0.030934721	-0.031117108	-0.031737368	0.017898089	0.003982514	-0.003643277	0.007986126	0.029607764	#DIV/0!
33	DAY_LAST_PHONE_CHANGE	0.091538841	0.11445685	0.094456961	0.115732447	0.064243239	0.158157576	0.137960847	0.086530723	0.136204331	#DIV/0!
34	FLAG_DOCUMENT_2	0.001436883	0.0147693	0.046382214	-0.018123198	0.03546858	-0.006297068	0.021935177	0.02340501	0.02416477	#DIV/0!
35	FLAG_DOCUMENT_3	-0.096028928	0.027796234	0.049492474	0.00086598	-0.036526013	0.026021803	0.046172849	0.011742752	0.043196477	#DIV/0!
36	FLAG_DOCUMENT_4	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!	#DIV/0!
37	FLAG_DOCUMENT_5	-0.004620618	0.00805165	0.0317324504	0.000534489	-0.005401576	0.00047853	0.061050354	-0.022764648	0.007421727	#DIV/0!
38	FLAG_DOCUMENT_6	-0.00471022	0.031691958	0.021665664	0.078617093	0.014553277	0.07059202	-0.017522591	0.006912122	0.018534261	#DIV/0!
< > ...	family	education	bivariate	contract_Status	Correlation_Non-Defaulters	Correlation	+				
	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE	YEARS_BIRTH	YEARS_EMPLOYED	YEARS_REGISTRATION	YEARS_ID_PUBLISH	FLAG_MOBIL	FLAG_EMP_PHONE
1		1									
2	AMT_INCOME_TOTAL	0.360011781	1								
3	AMT_CREDIT	0.431026919	0.760827873	1							
4	AMT_ANNUITY	0.367729007	0.98635817	0.765201743	1						
5	AMT_GOODS_PRICE	0.188785687	0.09564213	0.116108426	0.099641681	1					
6	REGION_POPULATION_RELATIVE	0.049352629	0.16087890	0.099056566	0.155143369	0.048987416	1				
7	YEARS_BIRTH	0.036221572	0.094943177	0.062134245	0.059518555	-0.005606334	0.352389434	1			
8	YEARS_EMPLOYED	-0.033500961	0.024294688	0.000851932	0.020392259	0.064621189	0.304733514	0.175692735	1		
9	YEARS_REGISTRATION	0.023115928	0.044246818	0.03222207	0.044495693	0.040352194	0.107692262	0.08250215	0.03702339	1	
10	YEARS_ID_PUBLISH	0.00254311	0.004263478	0.00070066	0.004159253	0.003826644	-0.00816817	0.004667144	-0.000181687	0.005949554	1
11	FLAG_MOBIL	0.00254311	0.004263478	0.00070066	0.004159253	0.003826644	-0.00816817	0.004667144	-0.000181687	0.005949554	1
12	FLAG_EMP_PHONE	0.001840188	-0.003321741	-0.000809968	-0.002120288	0.0395230	-0.008240593	0.000482146	0.001072225	-0.009423469	-3.76647E-05
13	FLAG_WORK_PHONE	-0.080193477	0.034281389	-0.00995972	-0.01965972	-0.019062959	0.040974482	0.013544813	-0.013564483	0.014735981	0.02096452
14	FLAG_CONT_MOBILE	-0.016152058	0.027480204	0.026055469	0.024835936	-0.006171892	0.00180144	0.012321114	-0.005343434	0.000519242	-0.000205715
15	FLAG_PHONE	0.011981609	0.022918689	0.018811979	0.030816451	0.019760336	0.0358450329	0.0560680072	0.056055531	0.028924793	0.003197907
16	FLAG_EMAIL	0.079189884	0.00476983	0.057534045	0.050506054	0.041495511	-0.067389587	-0.03088492	-0.022046867	-0.031549504	0.00133979
17	CNT_FAM_MEMBERS	-0.000847483	0.036117532	0.042839942	0.034143737	-0.010306544	0.180331567	-0.035963433	-0.151212131	0.114272484	0.001406748
18	REGION_RATING_CLIENT_W_CITY	-0.22990386	-0.1145176	-0.144021079	-0.11616212	-0.540963979	-0.048821639	-0.011631131	-0.093685548	-0.030329496	0.000107595
19	REGION_RATING_CLIENT	-0.22990386	-0.1145176	-0.144021079	-0.11616212	-0.540963979	-0.048821639	-0.011631131	-0.093685548	-0.030329496	0.000107595
20	HOUR_APPR_PROCESS_START	0.061241520	0.023978892	0.050214074	0.169741541	-0.05249584	-0.022972212	-0.022382729	-0.011940947	-0.01218541	-0.00283
21	REG_REGION_NOT_LIVE_REGION	0.070705756	0.024088354	0.042248024	0.0268503	-0.010125917	0.050584554	-0.026168535	-0.017619762	0.009077835	0.003188915
22	REG_REGION_NOT_WORK_REGION	0.150620706	0.052029517	0.076135793	0.053710538	0.069312685	-0.037096108	-0.084713495	-0.016582688	-0.020195986	0.001317332
23	LIVE_REGION_NOT_WORK_REGION	0.142728268	0.052033499	0.070788918	0.052444038	0.090678969	-0.031270406	-0.067876551	-0.004016713	-0.0089165	0.001170453
24	REG_CITY_NOT_WORK_CITY	-0.005249486	-0.032441358	-0.017797182	-0.031003034	-0.05272645	-0.169778772	-0.110955527	-0.052450252	-0.054359507	0.001610151
25	LIVE_CITY_NOT_WORK_CITY	-0.029395219	-0.038075471	-0.031084111	-0.038105486	-0.04546779	-0.105586448	-0.128027403	-0.048038224	-0.036453026	0.00319281
26	EXT_SOURCE_2	0.157592079	0.14329988	0.131789222	0.150147083	0.21040741	0.145947751	0.079985787	0.075749806	0.061074607	-0.00112609
27	EXT_SOURCE_3	-0.055446702	0.048696373	0.036034041	0.048434991	-0.015946362	0.157059546	0.112506111	0.086551796	0.086764284	-0.00059366
28	OBS_30_CNT_SOCIAL_CIRCLE	0.034628883	-0.003083088	-0.012532838	-0.003651531	-0.01943995	-0.022913038	-0.002616499	-0.017619762	0.009077835	0.003188915
29	DEF_30_CNT_SOCIAL_CIRCLE	-0.032942435	-0.016637802	-0.022537261	-0.0188242423	0.008881307	-0.016976929	-0.013442007	-0.008043937	-0.009174922	0.001610554
30	OBS_60_CNT_SOCIAL_CIRCLE	-0.03475411	-0.027585959	-0.012217567	-0.00340683	-0.018358055	-0.022686332	-0.002316044	-0.017610276	0.009490502	0.003177755
31	DEF_60_CNT_SOCIAL_CIRCLE	-0.033893862	-0.020419982	-0.024979398	-0.021464651	0.00209051	-0.0190726	-0.014226927	-0.011579482	-0.00959006	0.001370527
32	DAY_LAST_PHONE_CHANGE	0.03									

B	C	D	E
Top 5 Correlation (Non-Defaulters)			
Variable 1	Variable 2	Correlation	
OBS_60_CNT_SOCIAL_CIRCLE	OBS_30_CNT_SOCIAL_CIRCLE	0.998	
AMT_GOODS_PRICE	AMT_CREDIT	0.986	
REGION_RATING_CLIENT_W_CITY	REGION_RATING_CLIENT	0.948	
LIVE_REGION_NOT_WORK_REGION	REG_REGION_NOT_WORK_REGION	0.861	
DEF_60_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	0.853	

Top 5 Correlation (Defaulters)			
Variable 1	Variable 2	Correlation	
OBS_60_CNT_SOCIAL_CIRCLE	OBS_30_CNT_SOCIAL_CIRCLE	0.998	
AMT_GOODS_PRICE	AMT_CREDIT	0.982	
REGION_RATING_CLIENT_W_CITY	REGION_RATING_CLIENT	0.951	
DEF_60_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	0.891	
LIVE_REGION_NOT_WORK_REGION	REG_REGION_NOT_WORK_REGION	0.806	

## LINKS

### EXCEL SHEET

<https://docs.google.com/spreadsheets/d/1v9-xxUrDCIQv0AKGV6XI7dska5lKu5Y3/edit?usp=sharing&ouid=108286913145936487778&rtpof=true&sd=true>

### LOOM VIDEO

<https://www.loom.com/share/0bb4d6886c954c4081a76613f925af74?sid=23bc2763-2581-4188-9c76-86a824dc7a84>