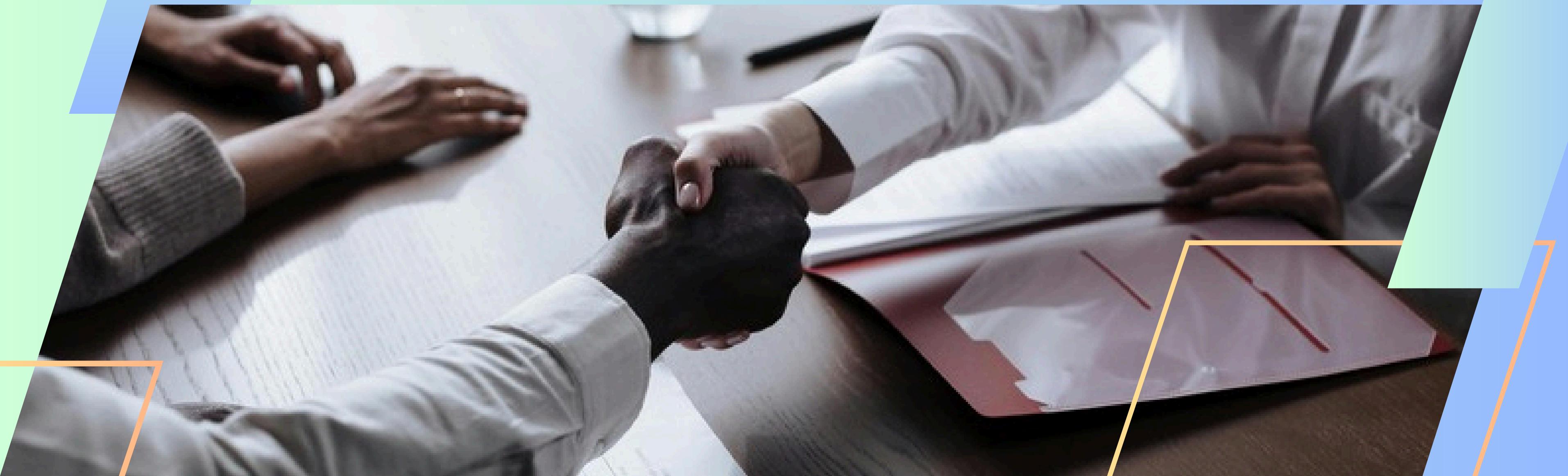


# Welcome To Presentation

I'M ANALYST, AND I'LL BE SHARING WITH YOU MY POINT OF VIEW.



A professional photograph of a woman with long dark hair, wearing an orange blazer over a teal top, sitting at a desk and looking directly at the camera with a neutral expression. The background is a blurred office environment.

# Bank Loan Case Study

ANKITA TANEJA

# Table Of Contents

1

## Introduction

### PROJECT DESCRIPTION

2

## Tasks

### TASKS DESCRIPTION

3

## Outputs

### TASKS OUTCOME

# I Introduction

PROJECT DESCRIPTION



# Strategic Perspectives



## Description

AS A DATA ANALYST AT A FINANCE COMPANY FOCUSED ON URBAN LOANS, YOUR TASK IS TO LEVERAGE EXPLORATORY DATA ANALYSIS (EDA) TO UNCOVER PATTERNS IN CUSTOMER AND LOAN DATA. THE GOAL IS TO IDENTIFY HIGH-RISK APPLICANTS WITH INSUFFICIENT CREDIT HISTORY WHO MAY DEFAULT WHILE ENSURING QUALIFIED APPLICANTS ARE APPROVED, STRIKING A BALANCE BETWEEN MINIMIZING LOSSES AND MAXIMIZING BUSINESS OPPORTUNITIES. THE GOAL OF THIS PROJECT IS TO USE EXPLORATORY Data Analysis (EDA) TO INVESTIGATE HOW CUSTOMER AND LOAN ATTRIBUTES IMPACT THE LIKELIHOOD OF LOAN DEFAULT. THIS ANALYSIS AIMS TO HELP THE COMPANY MINIMIZE FINANCIAL LOSSES BY IDENTIFYING HIGH-RISK APPLICANTS WHILE ENSURING ELIGIBLE APPLICANTS ARE NOT UNFAIRLY REJECTED.

# Mission And Vision



## OBJECTIVE

THE OBJECTIVE OF THIS PROJECT IS TO ANALYZE AND IDENTIFY PATTERNS THAT PREDICT A CUSTOMER'S LIKELIHOOD OF DEFAULTING ON INSTALLMENTS. THIS INSIGHT WILL HELP THE COMPANY MAKE INFORMED DECISIONS, SUCH AS DECLINING LOANS, ADJUSTING LOAN AMOUNTS, OR INCREASING INTEREST RATES FOR HIGH-RISK APPLICANTS, WHILE UNDERSTANDING THE KEY FACTORS INFLUENCING LOAN DEFAULTS TO IMPROVE APPROVAL STRATEGIES.

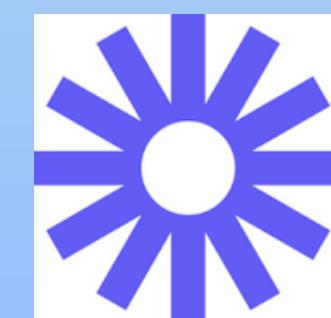

# Tech Stack



Microsoft Excel



Canva - Presentation



Loom - Video



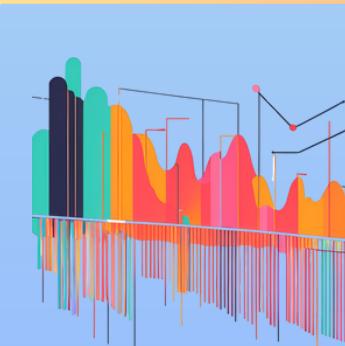
# Approach Used



Data Cleaning



Data Analyzing



Visualization





# 2 Tasks

TASKS DESCRIPTION

# Tasks



## One

IDENTIFY MISSING DATA AND  
DEAL WITH IT APPROPRIATELY



## Three

ANALYZE DATA IMBALANCE

PERFORM UNIVARIATE,  
SEGMENTED UNIVARIATE, AND  
BIVARIATE ANALYSIS



## Two

IDENTIFY OUTLIERS  
IN THE DATASET



## Four

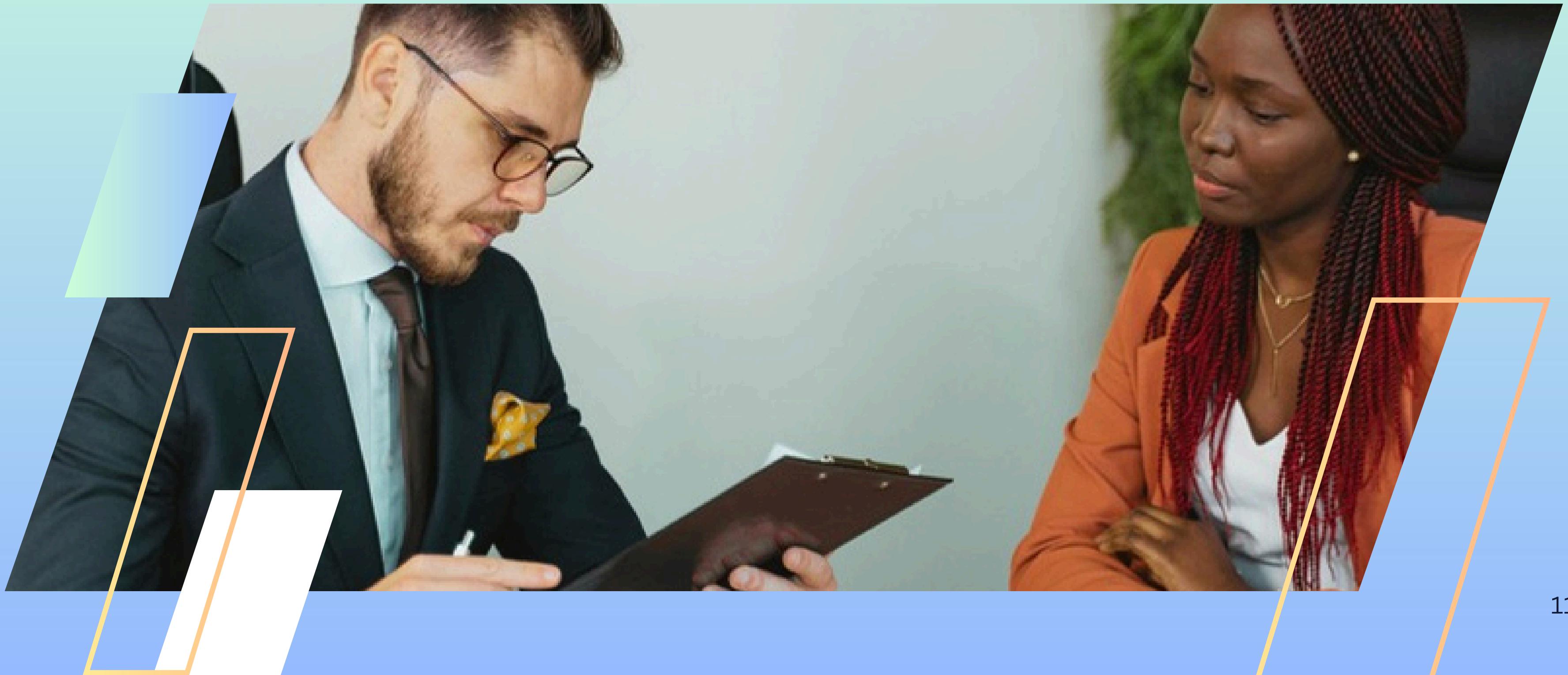


## Five

IDENTIFY TOP  
CORRELATIONS FOR  
DIFFERENT SCENARIOS



# A picture is worth a thousand words



# Tasks Description



## Task One

IDENTIFY THE MISSING DATA IN THE DATASET & DECIDE ON AN APPROPRIATE METHOD TO DEAL WITH IT USING EXCEL BUILT-IN FUNCTIONS & FEATURES.

## Task Two

DETECT & IDENTIFY OUTLIERS IN THE DATASET USING EXCEL STATS FUNCTION & FEATURE, FOCUSING ON NUMERICAL VARIABLES.

## Task Three

DETERMINE IF THERE IS DATA IMBALANCE IN THE LOAN APPLICATION DATASET & CALCULATE THE RATIO OF IMP. DATA IMBALANCE USING VARIOUS ANALYSES ON EXCEL FUNCTIONS..

## Task Four

To GAIN INSIGHTS INTO THE DRIVING FACTORS OF LOAN DEFAULT, IT IS CONDUCT VARIOUS ANALYSES ON CONSUMER & LOAN ATTRIBUTES.

## Task Five

SEGMENT THE DATASET BASED ON DIFFERENT SCENARIOS & IDENTIFY THE TOP CORRELATIONS FOR EACH SEGMENTED DATA USING EXCEL FUNCTIONS.



# 3 Output

OUTCOME RESULTS

# Task I

## IDENTIFY MISSING DATA AND DEAL WITH IT APPROPRIATELY

WHEN HANDLING MISSING DATA IN A LOAN APPLICATION DATASET USING EXCEL, THE PROCESS CAN BE SUMMARIZED INTO THE FOLLOWING STEPS:

### 1. IDENTIFY MISSING DATA

- USE THE ISBLANK FUNCTION TO FLAG MISSING CELLS.
  - FORMULA: =ISBLANK(CELL\_REFERENCE)
  - THIS WILL RETURN TRUE FOR MISSING CELLS AND FALSE OTHERWISE.
  - FORMULA: =BLANKCOUNT(RANGE)
  - THIS WILL RETURN COUNT OF BLANK CELLS.
- USE THE COUNTIF OR COUNT FUNCTION TO COUNT MISSING VALUES IN EACH COLUMN.
  - FORMULA: =COUNTIF(RANGE, "")
  - THIS COUNTS ALL BLANK CELLS IN THE SPECIFIED RANGE.

### 2. SUMMARIZE MISSING DATA

- CREATE A SUMMARY TABLE WITH COLUMN NAMES AND THE COUNT OF MISSING VALUES FOR EACH COLUMN.
  - USE THE ABOVE COUNTIF FORMULAS FOR EACH COLUMN.
- CALCULATE THE PROPORTION OF MISSING DATA FOR EACH COLUMN:
  - FORMULA: =COUNTIF(RANGE, "")/COUNTA(RANGE)
  - THIS PROVIDES THE PERCENTAGE OF MISSING VALUES.

### **3. HANDLE MISSING DATA**

DEPENDING ON THE TYPE OF DATA AND THE PROPORTION OF MISSING VALUES:

IMPUTATION FOR NUMERICAL DATA:

USE THE AVERAGE FUNCTION TO REPLACE MISSING VALUES WITH THE COLUMN MEAN.

FORMULA: =IF(ISBLANK(CELL\_REFERENCE), AVERAGE(RANGE), CELL\_REFERENCE)

ALTERNATIVELY, USE THE MEDIAN FUNCTION FOR SKEWED DATA.

IMPUTATION FOR CATEGORICAL DATA:

USE THE MODE FUNCTION TO REPLACE MISSING VALUES WITH THE MOST FREQUENT CATEGORY.

FORMULA: =IF(ISBLANK(CELL\_REFERENCE), MODE(RANGE), CELL\_REFERENCE)

DELETION:

IF A COLUMN OR ROW HAS A HIGH PERCENTAGE OF MISSING DATA, CONSIDER REMOVING IT, BUT ONLY IF IT WON'T HARM THE ANALYSIS.

### **4. VISUALIZE MISSING DATA**

CREATE A SUMMARY CHART TO VISUALIZE THE MISSING DATA DISTRIBUTION:

HIGHLIGHT THE SUMMARY TABLE WITH COLUMN NAMES AND PROPORTIONS.

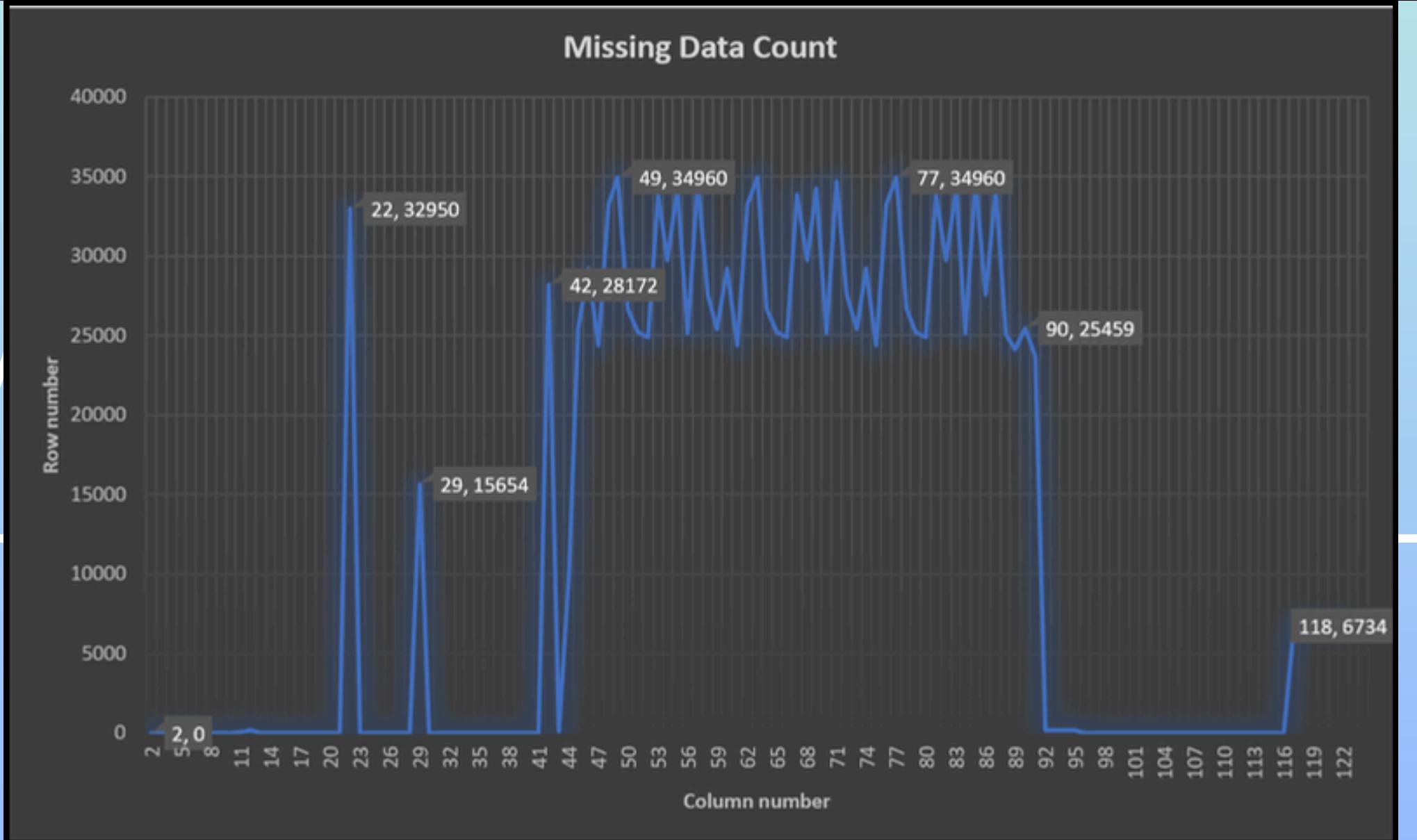
INSERT A BAR CHART OR COLUMN CHART:

GO TO THE INSERT TAB, SELECT BAR CHART OR COLUMN CHART.

CUSTOMIZE THE CHART TO LABEL AXES AND ADD A TITLE, E.G., "PROPORTION OF MISSING DATA BY VARIABLE".

BY FOLLOWING THESE STEPS, YOU CAN IDENTIFY, SUMMARIZE, AND HANDLE MISSING DATA EFFECTIVELY IN EXCEL.

PROPORTION
COUNTBLANK
Total Missing Values
1488212
Column with NO missing data
SK_ID_CURR, TARGET, NAME_CONTRACT_TYPE, CODE_GENDER, etc.
Column with missing data
AMT_ANNUITY, AMT_GOODS_PRICE, NAME_TYPE_SUITE etc.
Orange Color Column represents the least number of Missing Data Count. - Repaired
Red Color Column represents the number of Missing Data Count Greater 35% - Deleted



## INSIGHTS

TO CALCULATE THE TOTAL NUMBER OF MISSING VALUES FOR EACH COLUMN, USE:  
EXCEL

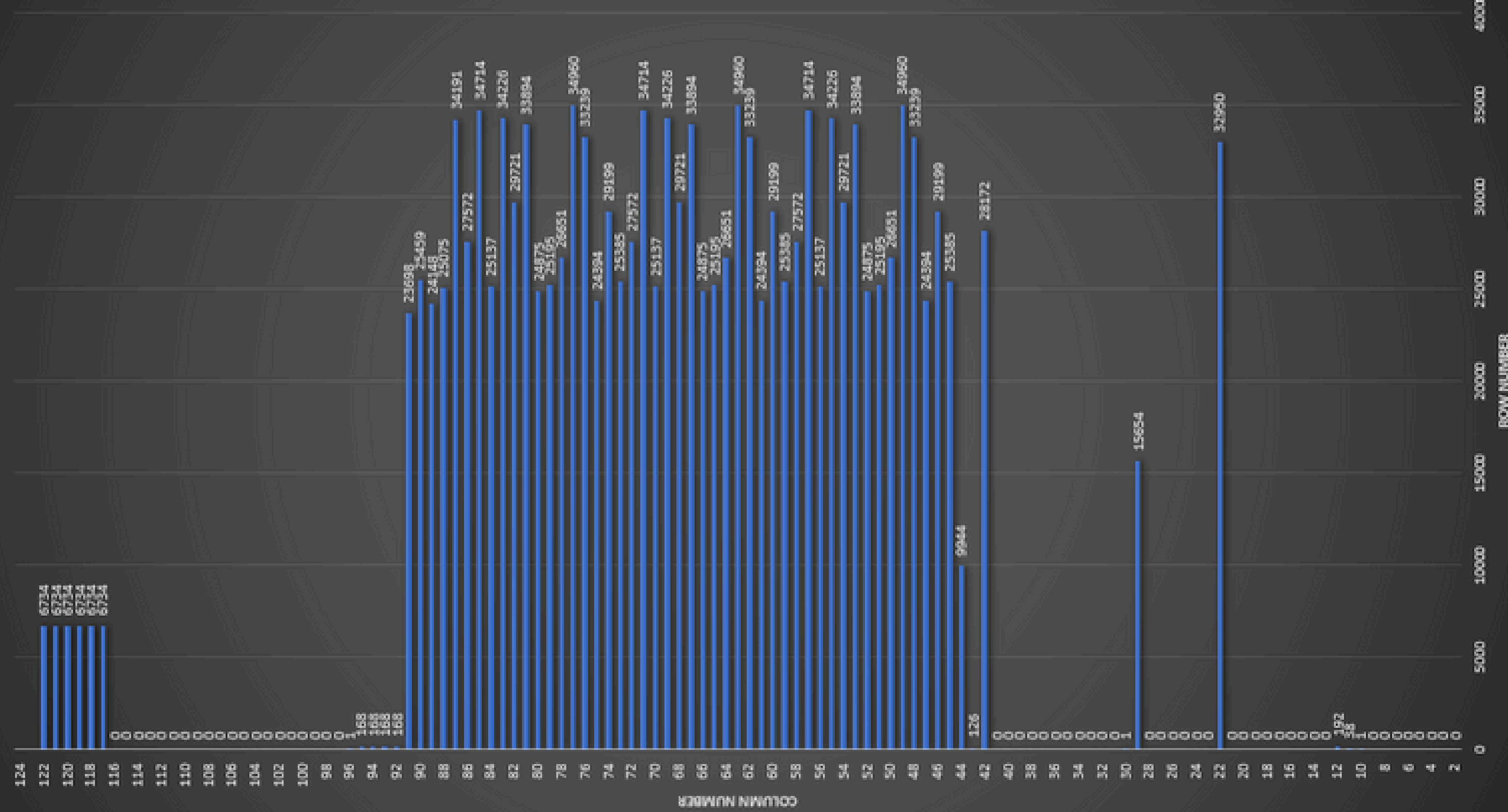
=COUNTBLANK(COLUMN\_RANGE)

**TOTAL MISSING VALUES: 1488212**

TO CALCULATE THE PROPORTION:  
=COUNTBLANK(COUNT\_RANGE)/Rows Count

B1		=COUNTBLANK(B4:B50003) / (50000-1)					
1	PROPORTION	A					
2	COUNTBLANK	B 0.00% 0					
B2		=COUNTBLANK(B4:B50002)					
1	PROPORTION	A					
2	COUNTBLANK	B 0.00% 0					
Days_Birth	Days_Birth	Days_Employed	Days_Employed	Days_Registration	Days_Registered	Days_ID_PU	Days_ID
-9461	25.92054795	-637	1.745205479	-3648	9.994520548	-2120	5.808219178
-19046	52.18082192	-225	0.616438356	-4260	11.67123288	-2531	6.934246575
-16941	46.41369863	-1588	4.350684932	-4970	13.61643836	-477	1.306849315
-13778	37.74794521	-3130	8.575342466	-1213	3.323287671	-619	1.695890411
-18850	51.64383562	-449	1.230136986	-4597	12.59452055	-2379	6.517808219
-20099	55.06575342	365243	1000.665753	-7427	20.34794521	-3514	9.62739726
-10197	27.9369863	-679	1.860273973	-4427	12.12876712	-738	2.021917808
-20417	55.9369863	365243	1000.665753	-5246	14.37260274	-2512	6.882191781
-13439	36.81917808	-2717	7.443835616	-311	0.852054795	-3227	8.84109589
-14086	38.59178082	-3028	8.295890411	-643	1.761643836	-4911	13.45479452
-8728	23.91232877	-1157	3.169863014	-3494	9.57260274	-1368	3.747945205
-12931	35.42739726	-1317	3.608219178	-6392	17.51232877	-3866	10.59178082
-17718	48.54246575	-7804	21.38082192	-8751	23.97534247	-1259	3.449315068
-11348	31.09041096	-2038	5.583561644	-1021	2.797260274	-3964	10.86027397
-14815	40.5890411	-1652	4.526027397	-2299	6.298630137	-2299	6.298630137
-11146	30.5369863	-4306	11.79726027	-114	0.312328767	-2518	6.898630137
-24827	68.01917808	365243	1000.665753	-9012	24.69041096	-3684	10.09315068
-11286	30.92054795	-746	2.043835616	-108	0.295890411	-3729	10.21643836
-19334	52.96986301	-3494	9.57260274	-2419	6.62739726	-2893	7.926027397
-18724	51.29863014	-2628	7.2	-6573	18.00821918	-1827	5.005479452
-15948	43.69315068	-1234	3.380821918	-5782	15.84109589	-3153	8.638356164
-9994	27.38082192	-1796	4.920547945	-4668	12.7890411	-2661	7.290410959
-15280	41.8630137	-2668	7.309589041	-5266	14.42739726	-3787	10.37534247
-12974	35.54520548	-4404	12.06575342	-7123	19.51506849	-4464	12.23013699
-11694	32.03835616	-2060	5.643835616	-3557	9.745205479	-3557	9.745205479
-12158	33.30958904	-1275	3.493150685	-6265	17.16438356	-2009	5.504109589
-17199	47.12054795	-768	2.104109589	-63	0.17260274	-735	2.01369863
-21077	57.74520548	-1288	3.528767123	-5474	14.99726027	-4270	11.69863014
22020	22.02054795	365243	1000.665753	-9817	2.9890001	-4969	12.61369863

## Missing Data Count



# PROPORTION

250.00%

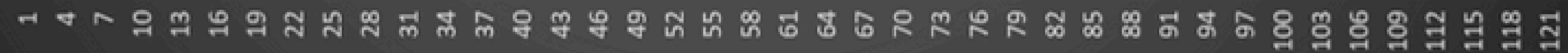
200.00%

150.00%

100.00%

50.00%

0.00%



# Task 2

## IDENTIFY OUTLIERS IN THE DATASET

TO IDENTIFY AND VISUALIZE OUTLIERS IN A LOAN APPLICATION DATASET USING EXCEL, FOLLOW THESE STEPS:

### 1. IDENTIFY OUTLIERS

#### STEP 1: CALCULATE THE INTERQUARTILE RANGE (IQR)

- IQR FORMULA:  $IQR = Q_3 - Q_1$ 
  - USE THE QUARTILE FUNCTION TO FIND  $Q_1$  (25TH PERCENTILE) AND  $Q_3$  (75TH PERCENTILE):
    - FORMULA FOR  $Q_1$ : =QUARTILE(RANGE, 1)
    - FORMULA FOR  $Q_3$ : =QUARTILE(RANGE, 3)
  - CALCULATE THE IQR:
    - FORMULA: = $Q_3 - Q_1$

#### STEP 2: DEFINE OUTLIER THRESHOLDS

- LOWER BOUND:  $Q_1 - (1.5 \times IQR)$
- UPPER BOUND:  $Q_3 + (1.5 \times IQR)$ 
  - FORMULA FOR LOWER BOUND: = $Q_1 - 1.5 * IQR$
  - FORMULA FOR UPPER BOUND: = $Q_3 + 1.5 * IQR$

#### STEP 3: IDENTIFY OUTLIERS

- USE AN IF STATEMENT TO FLAG OUTLIERS:
  - FORMULA: =IF(OR(CELL\_REFERENCE < LOWER\_BOUND, CELL\_REFERENCE > UPPER\_BOUND), "OUTLIER", "NORMAL")

## 2. VISUALIZE OUTLIERS

### Box Plot

- SELECT THE NUMERICAL DATA RANGE.
- GO TO THE INSERT TAB → INSERT STATISTIC CHART → BOX AND WHISKER.
- CUSTOMIZE THE BOX PLOT TO DISPLAY THE DISTRIBUTION AND HIGHLIGHT OUTLIERS.

### SCATTER PLOT

- USE THE NUMERICAL VARIABLE AS THE Y-AXIS AND AN IDENTIFIER (E.G., ROW NUMBER) AS THE X-AXIS.
- HIGHLIGHT OUTLIERS WITH A DIFFERENT COLOR:
  - ADD A NEW COLUMN INDICATING "OUTLIER" OR "NORMAL".
  - USE CONDITIONAL FORMATTING TO COLOR CELLS OR POINTS IN THE SCATTER PLOT FOR OUTLIERS.

## 3. INVESTIGATE OUTLIERS

- THRESHOLDS: EVALUATE IF THE OUTLIERS ARE VALID BASED ON BUSINESS RULES OR DOMAIN KNOWLEDGE.
  - FOR EXAMPLE, IN A LOAN DATASET, UNUSUALLY HIGH-INCOME VALUES COULD BE VALID FOR WEALTHY APPLICANTS.
- TREATMENT:
  - RETAIN VALID DATA: IF THE OUTLIER IS A LEGITIMATE VALUE, INCLUDE IT IN THE ANALYSIS.
  - CORRECT ERRORS: REPLACE INCORRECT DATA WITH APPROPRIATE VALUES (E.G., MEAN, MEDIAN).
  - EXCLUDE INVALID DATA: REMOVE OUTLIERS ONLY IF THEY DISTORT THE ANALYSIS.

QUARTILE 1		0	0	112500	274779	16456.5
QUARTILE 3		0	1	202500	820638	34587
IQR		0	1	90000	545859	18130.5
UL		0	2.5	337500	1639426.5	61782.75
LL		0	-1.5	-22500	-544009.5	-10739.25

*fx* =QUARTILE.EXC(B7:B56341,1)

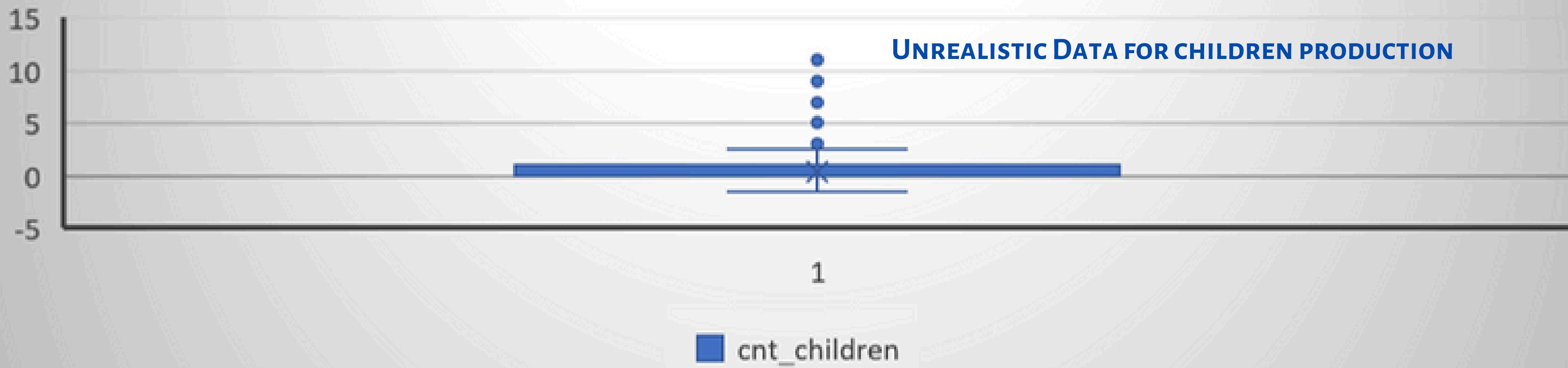
*fx* =QUARTILE.EXC(B7:B56341,3)

*fx* =B2-B1

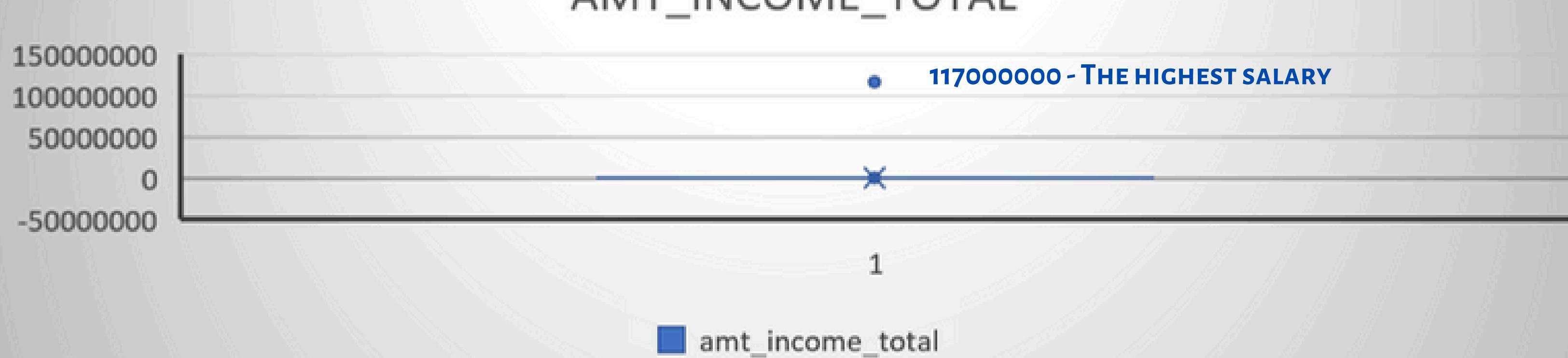
*fx* =B2+(1.5\*B3)

*fx* =B1-(1.5\*B3)

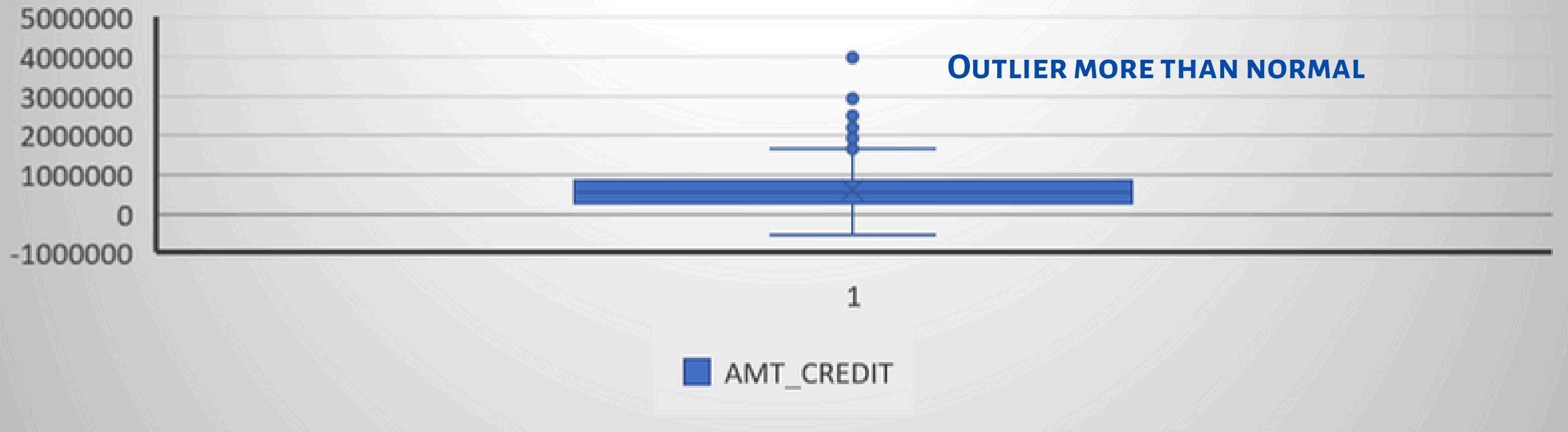
## CNT\_Children



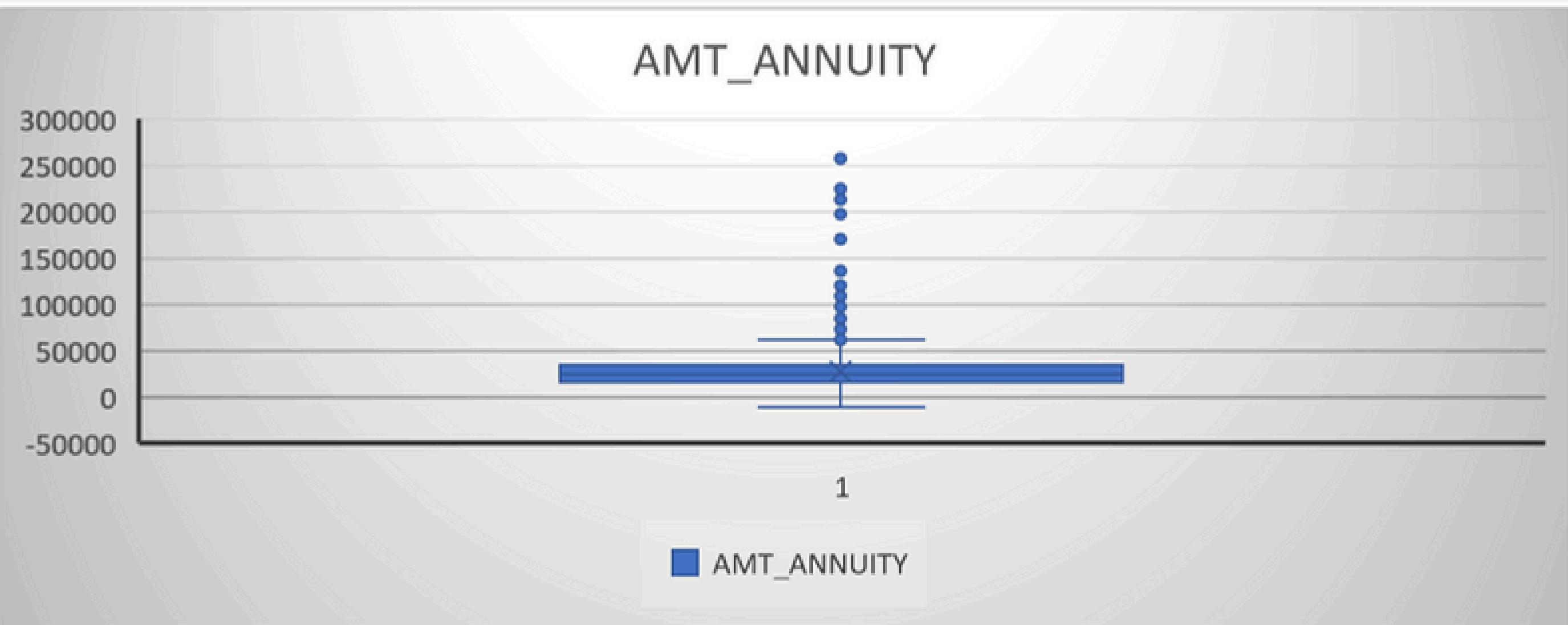
## AMT\_INCOME\_TOTAL



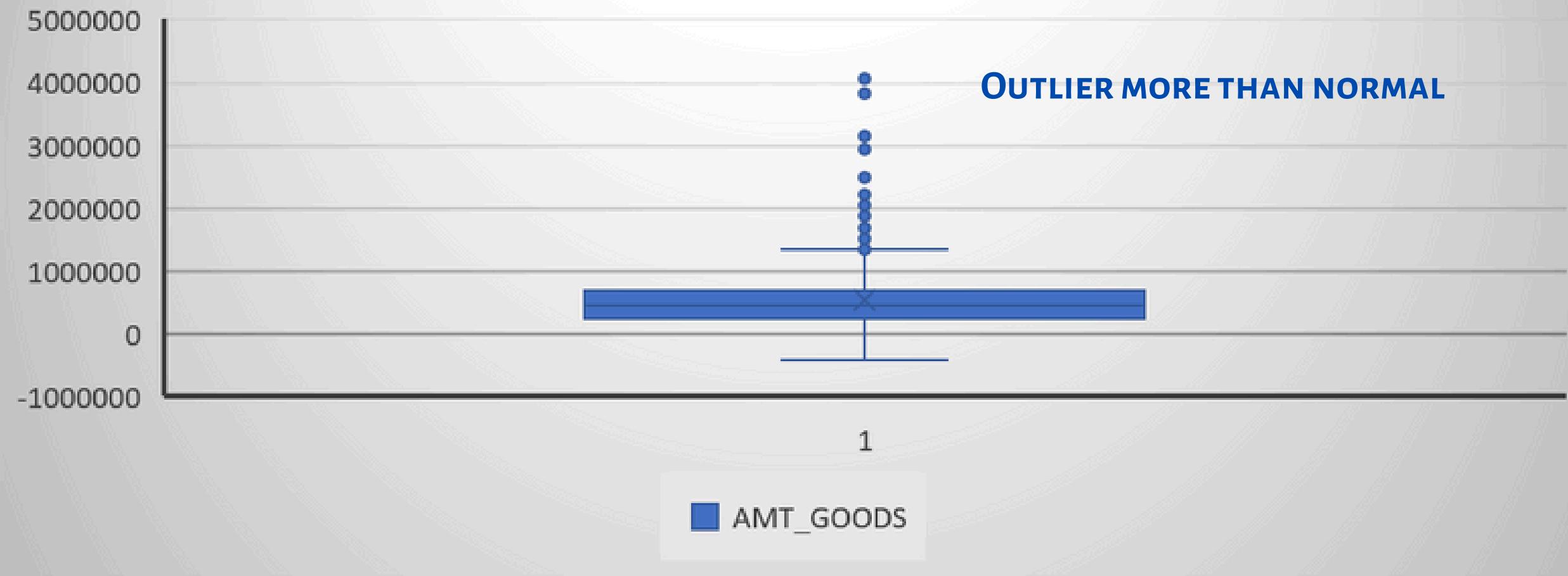
## AMT\_CREDIT



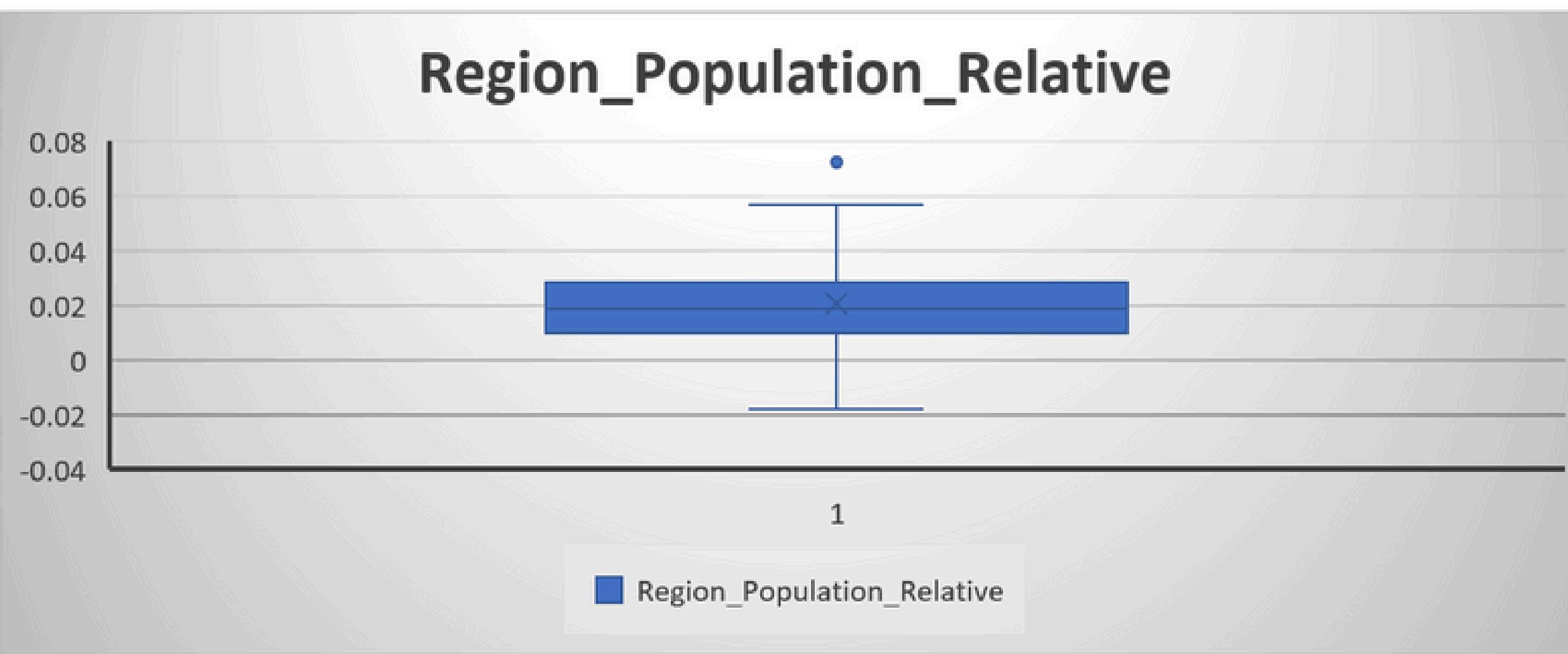
## AMT\_ANNUITY



### AMT\_GOODS

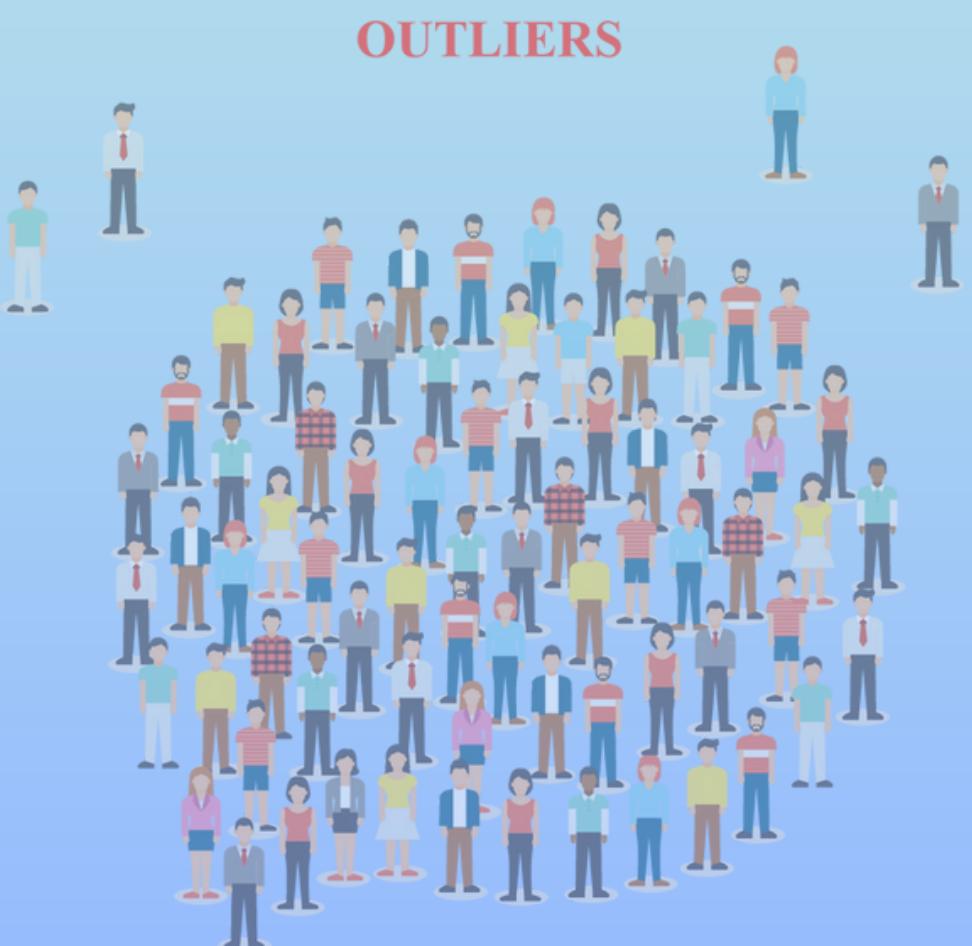


### Region\_Population\_Relative





DOTS IN THE BOX GRAPH ARE THE OUTLIERS LIKE YEARS\_EMPLOYED HAVE THE OUTLIER AT 7100 WHICH IS NOT POSSIBLE, AMT\_GOODS HAVE THE ABNORMAL OUTLIER SHOWN IN THE GRAPH, CNT\_CHILDREN CAN BE 11 WHICH IS AN OUTLIER AGAIN, TOTAL INCOME CANT BE LESS THAN ZERO ETC.



# Task 3

## ANALYZE DATA IMBALANCE

DATA IMBALANCE CAN SIGNIFICANTLY INFLUENCE THE ACCURACY OF MACHINE LEARNING MODELS, PARTICULARLY IN BINARY CLASSIFICATION TASKS. A KEY STEP IN UNDERSTANDING THIS IS ANALYZING THE DISTRIBUTION OF YOUR TARGET VARIABLE. THIS INVOLVES DETERMINING THE FREQUENCY OF EACH CLASS AND EVALUATING WHETHER ONE CLASS SIGNIFICANTLY OUTWEIGHS THE OTHER.

### STEPS TO ASSESS DATA IMBALANCE IN EXCEL

#### 1. ANALYZE CLASS DISTRIBUTION:

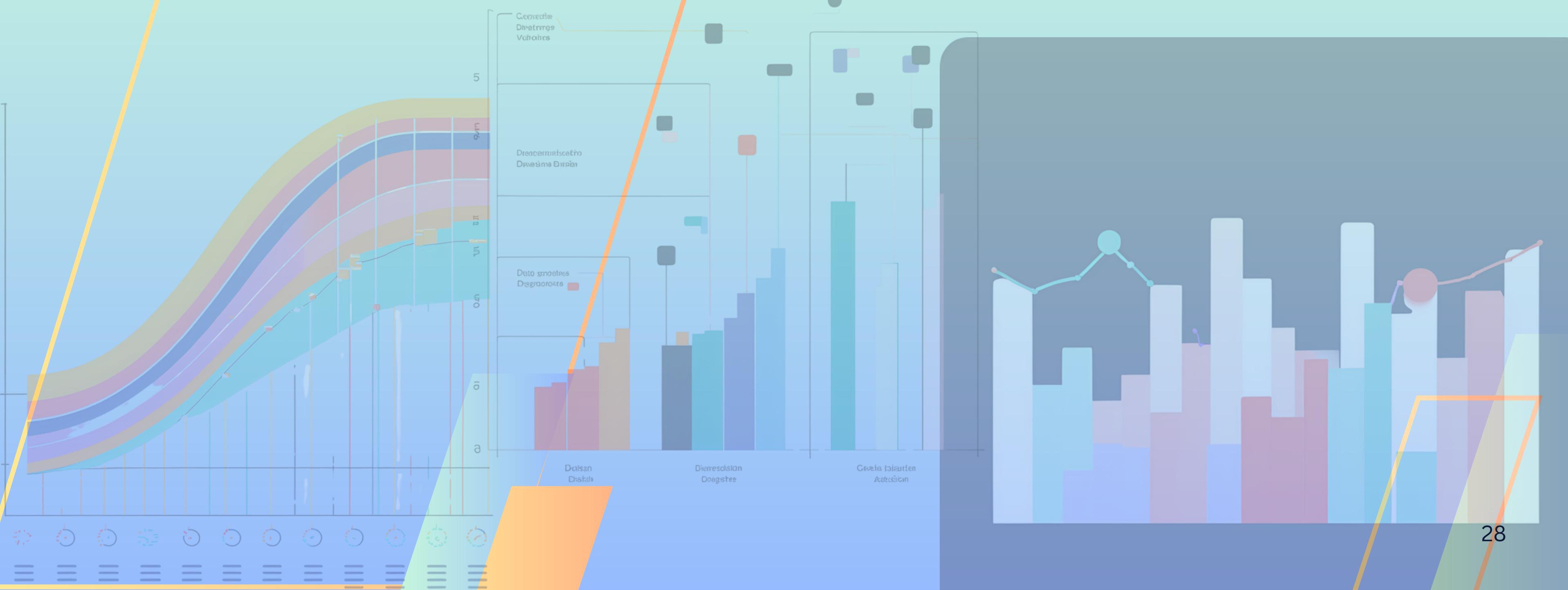
- USE THE COUNTIF FUNCTION TO CALCULATE THE NUMBER OF INSTANCES FOR EACH CLASS IN THE TARGET VARIABLE:
  - EXAMPLE: IF THE TARGET VARIABLE IS IN COLUMN A, USE:
    - =COUNTIF(A:A, "CLASS1") TO COUNT OCCURRENCES OF CLASS 1.
    - =COUNTIF(A:A, "CLASS2") TO COUNT OCCURRENCES OF CLASS 2.
- COMPUTE THE TOTAL NUMBER OF SAMPLES USING THE SUM FUNCTION:
  - =SUM(COUNTIF(A:A, "CLASS1"), COUNTIF(A:A, "CLASS2")).

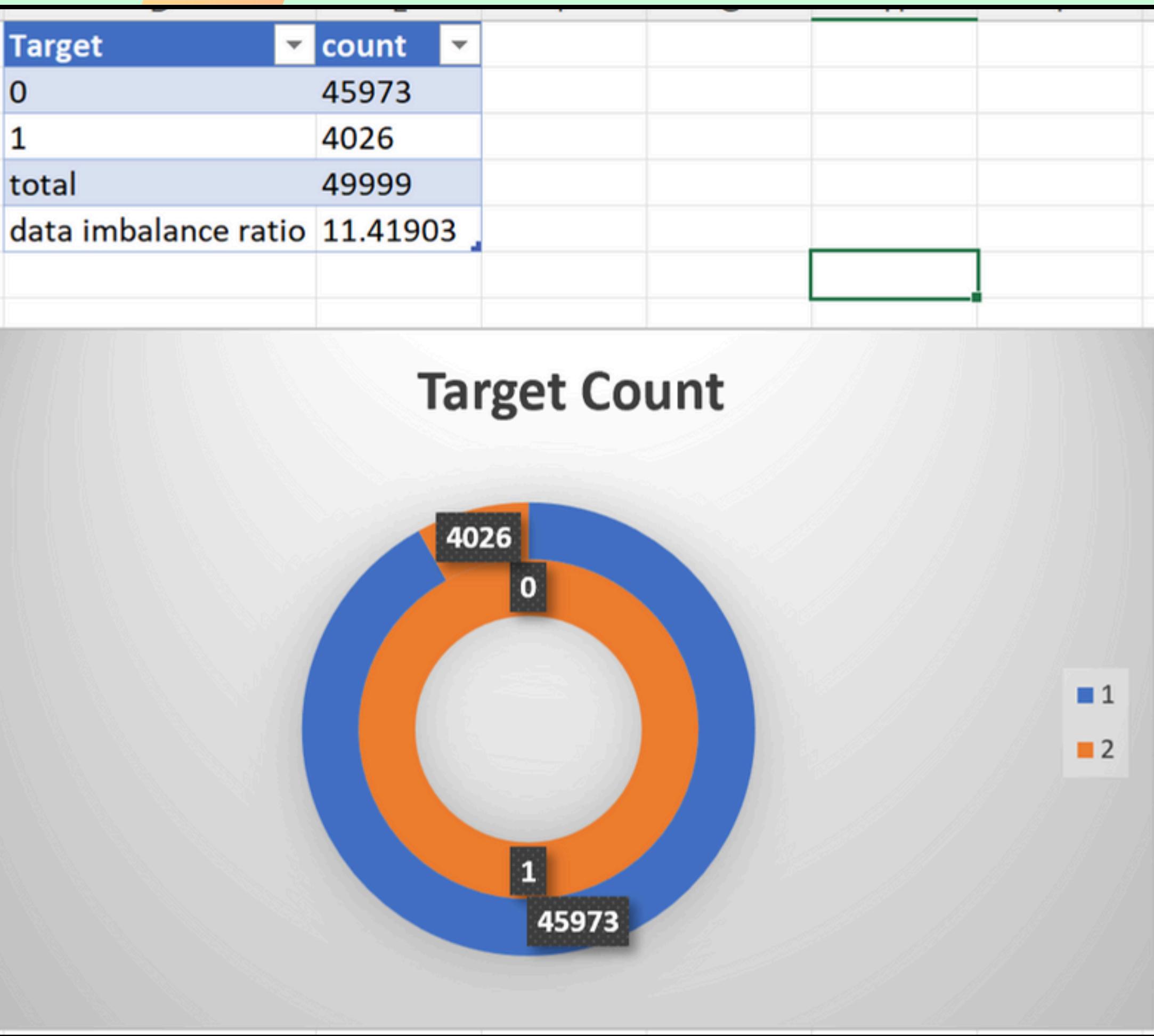
#### 2. CALCULATE THE IMBALANCE RATIO:

- COMPUTE THE PROPORTION OF EACH CLASS:
  - PROPORTION OF CLASS 1: =COUNTIF(A:A, "CLASS1") / TOTAL\_SAMPLES.
  - PROPORTION OF CLASS 2: =COUNTIF(A:A, "CLASS2") / TOTAL\_SAMPLES.
- DETERMINE THE RATIO:
  - IMBALANCE RATIO = PROPORTION OF MINORITY CLASS / PROPORTION OF MAJORITY CLASS.

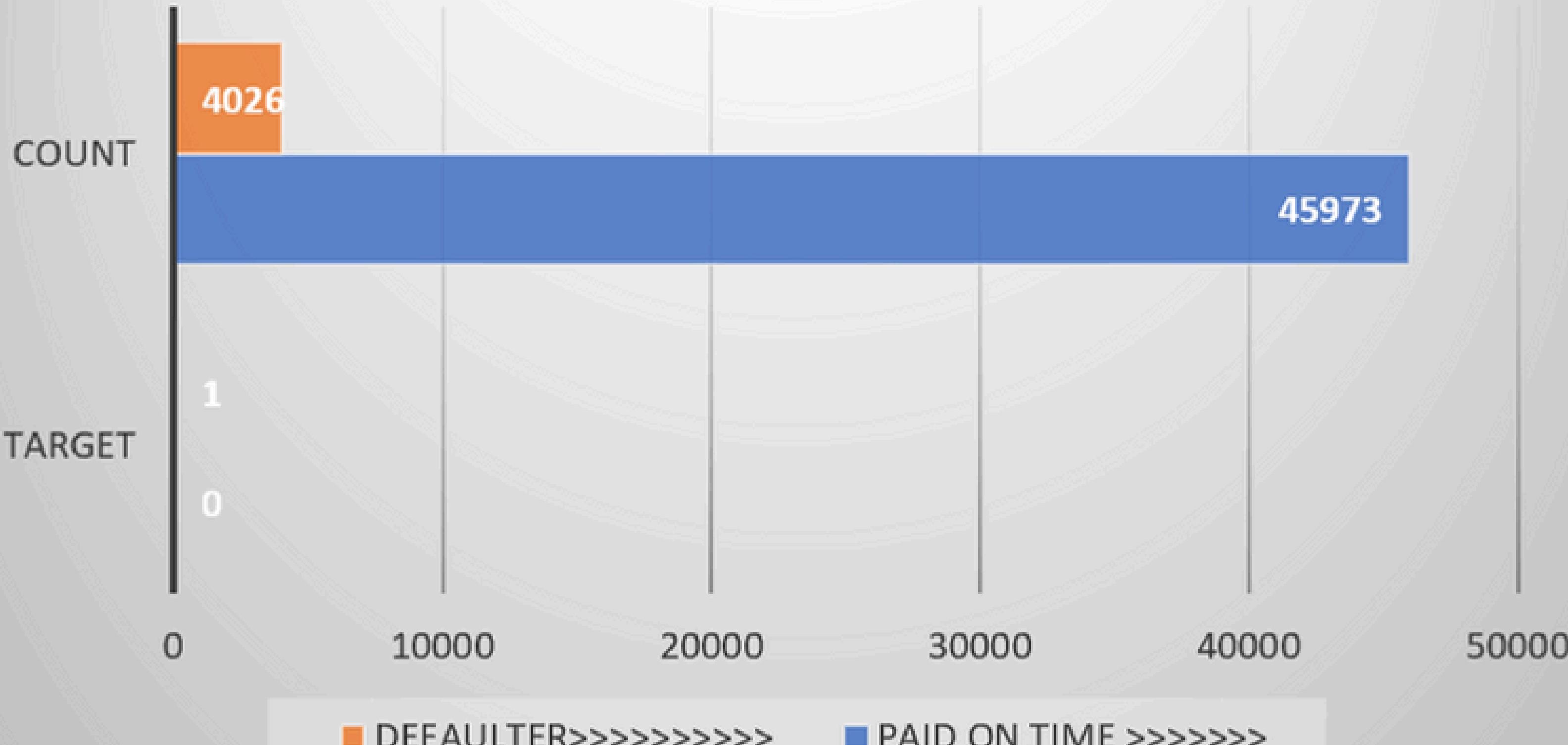
### 3. VISUALIZE THE DISTRIBUTION:

- CREATE A PIE CHART OR BAR CHART:
    - SELECT THE CLASS FREQUENCIES (E.G., COUNT OF CLASS 1 AND CLASS 2).
    - INSERT A PIE CHART OR BAR CHART TO VISUALIZE THE DISTRIBUTION.
  - HIGHLIGHT THE DOMINANT CLASS AND THE SMALLER CLASS FOR BETTER CLARITY



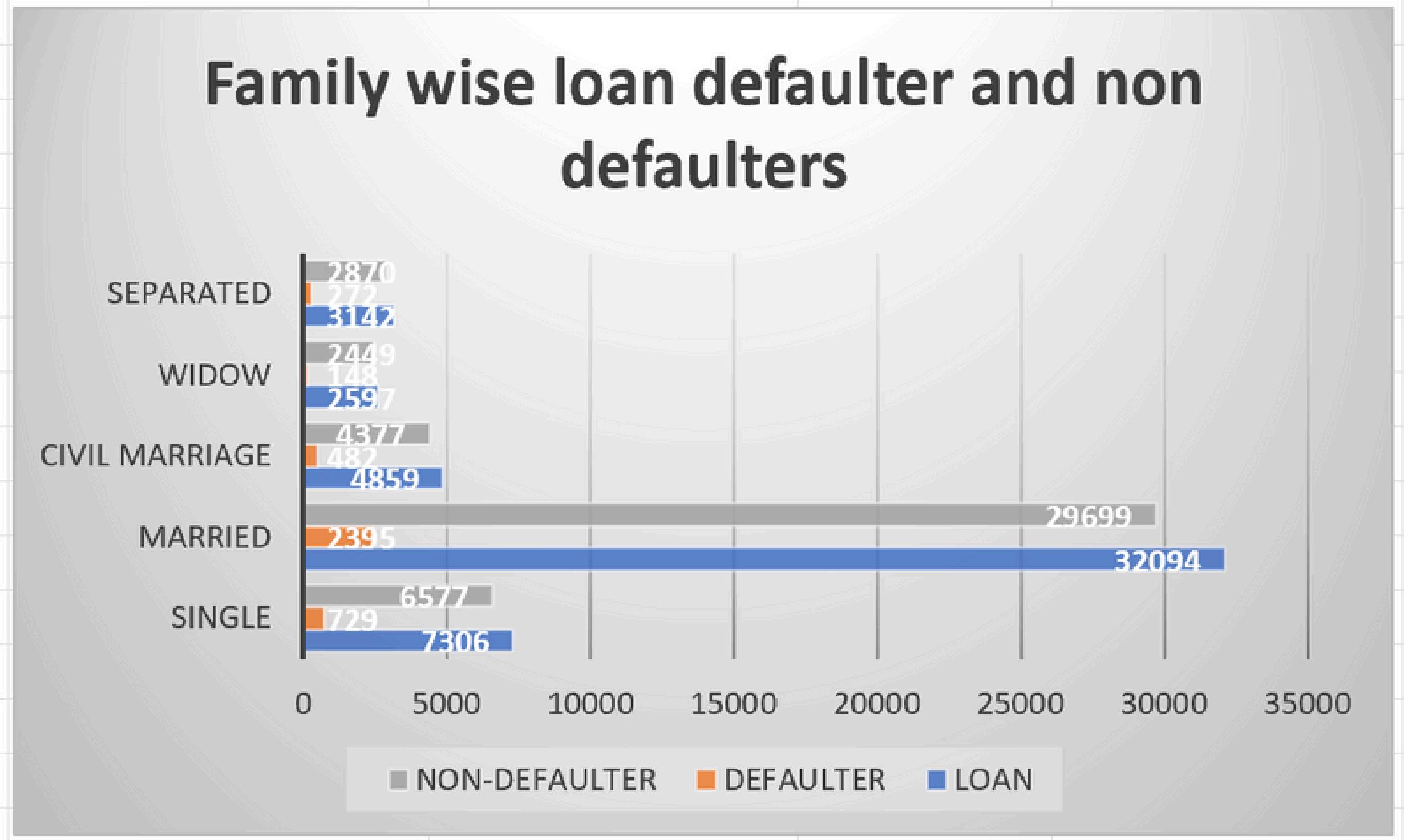


# TARGET COUNT



FAMILY STATUS	LOAN	DEFALUTER	NON-DEFALUTER
Single	7306	729	6577
Married	32094	2395	29699
Civil Marriage	4859	482	4377
Widow	2597	148	2449
Separated	3142	272	2870

## Family wise loan defaulter and non defaulters



# Task 4

## PERFORM UNIVARIATE, SEGMENTED UNIVARIATE, AND BIVARIATE ANALYSIS

To better understand the factors contributing to loan default, it is essential to analyze consumer and loan attributes using various techniques. Here's how you can perform univariate, segmented univariate, and bivariate analyses in Excel:

### 1. UNIVARIATE ANALYSIS:

Understand the distribution of individual variables.

- STEPS:
  - A. USE EXCEL STATISTICAL FUNCTIONS TO CALCULATE KEY METRICS LIKE:
    - MEAN: =AVERAGE(RANGE)
    - MEDIAN: =MEDIAN(RANGE)
    - STANDARD DEVIATION: =STDEV.P(RANGE)
    - COUNTS: =COUNT(RANGE)
  - B. USE HISTOGRAMS OR BAR CHARTS TO VISUALIZE THE DISTRIBUTION:
    - SELECT THE DATA AND GO TO INSERT > CHART > HISTOGRAM/BAR CHART.
  - PURPOSE: IDENTIFY PATTERNS SUCH AS CENTRAL TENDENCY, SPREAD, AND SKEWNESS IN THE DATA.

### 2. SEGMENTED UNIVARIATE ANALYSIS:

Compare variable distributions across different scenarios (e.g., defaulted vs. non-defaulted loans).

- STEPS:
  - A. FILTER OR SORT DATA BY SCENARIOS USING EXCEL FILTERS OR CONDITIONS.

## B. USE PIVOT TABLES FOR GROUPED STATISTICS:

- SELECT DATA AND INSERT A PIVOT TABLE.
- PLACE THE VARIABLE OF INTEREST IN ROWS AND THE SEGMENTATION (E.G., DEFAULTED/Non-DEFAULTED) IN COLUMNS.

- ADD STATISTICS LIKE COUNTS OR AVERAGES TO THE VALUES FIELD.

- CREATE GROUPED OR STACKED BAR CHARTS:

- SELECT PIVOT TABLE RESULTS AND USE INSERT > CHART > STACKED BAR/COLUMN CHART.

- PURPOSE: UNDERSTAND HOW DISTRIBUTIONS DIFFER BETWEEN GROUPS, HIGHLIGHTING KEY CONTRASTS.

## 3. BIVARIATE ANALYSIS:

EXPLORE RELATIONSHIPS BETWEEN VARIABLES AND THE TARGET VARIABLE (E.G., LOAN DEFAULT).

- STEPS:

- CREATE SCATTER PLOTS:

- HIGHLIGHT TWO COLUMNS OF INTEREST AND USE INSERT > SCATTER PLOT.
  - ADD TRENDLINES IF NEEDED TO ASSESS RELATIONSHIPS.

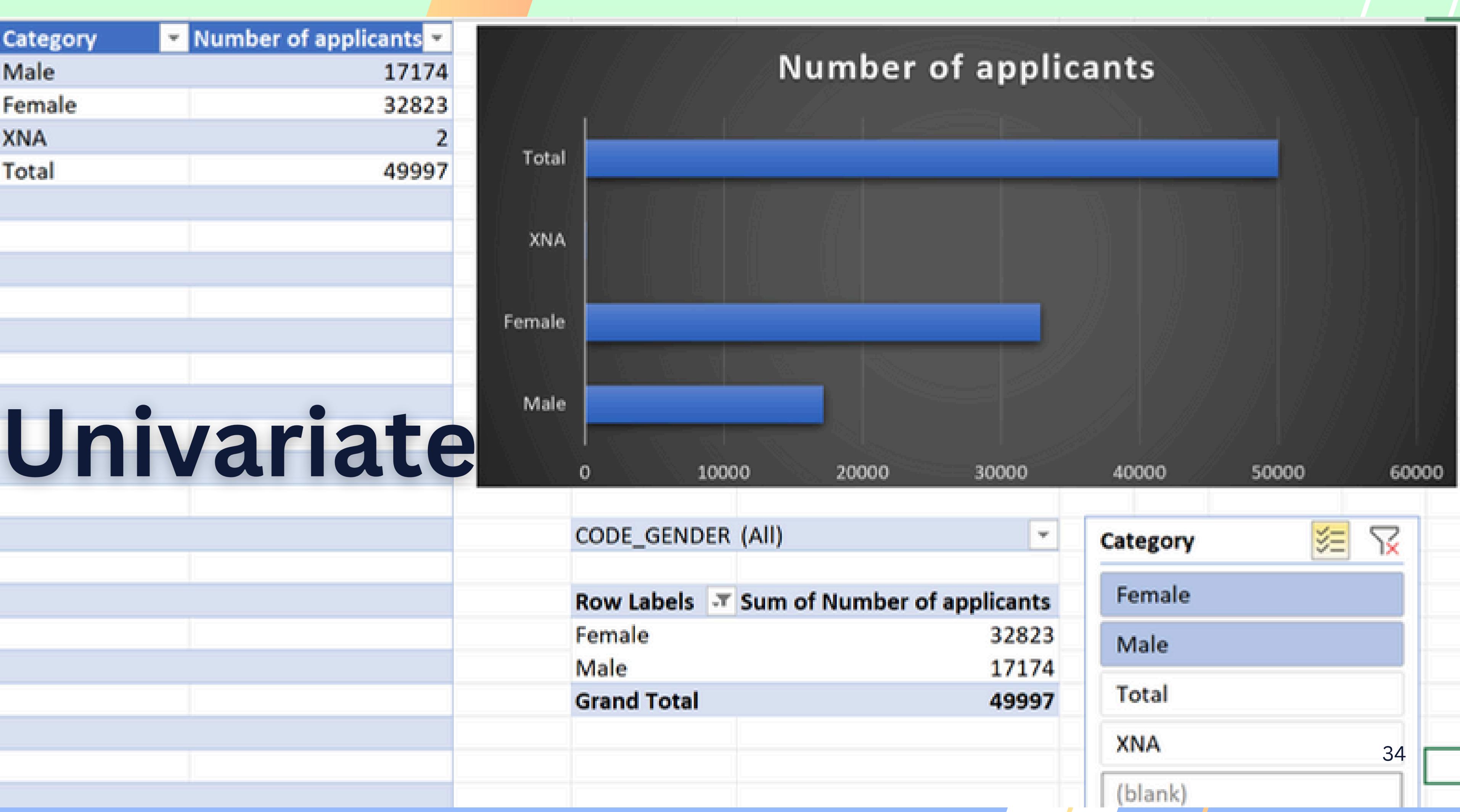
- GENERATE CROSS-TABULATIONS WITH PIVOT TABLES:

- PLACE ONE VARIABLE IN ROWS AND ANOTHER IN COLUMNS, AND USE COUNTS OR PERCENTAGES IN VALUES.

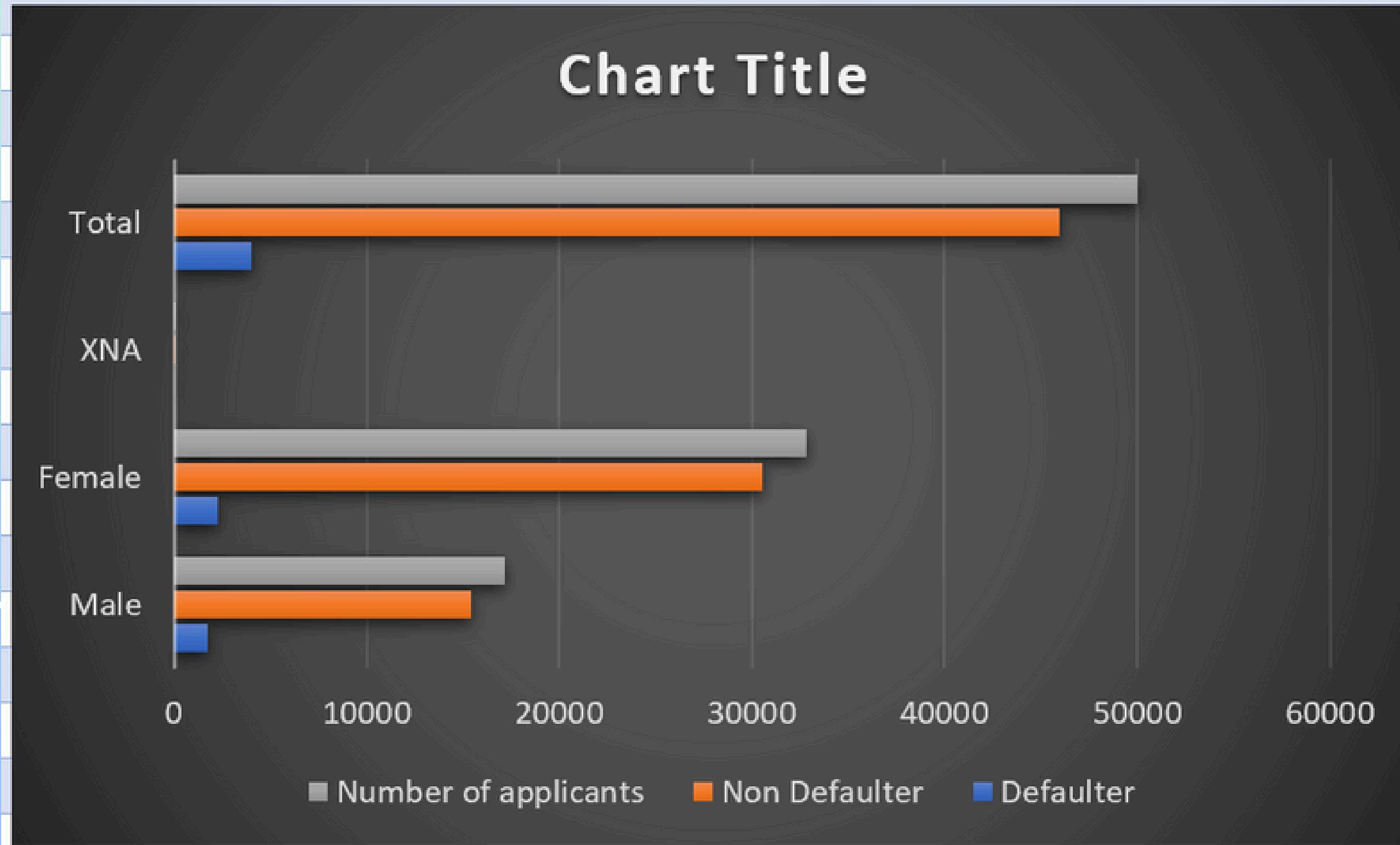
- USE HEATMAPS TO VISUALIZE CORRELATIONS:

- CALCULATE CORRELATIONS USING `=CORREL(RANGE1, RANGE2)` AND COLOR-CODE RESULTS USING CONDITIONAL FORMATTING.

- PURPOSE: DETECT PATTERNS OR TRENDS (E.G., HIGHER LOAN AMOUNTS CORRELATE WITH HIGHER DEFAULT RATES).



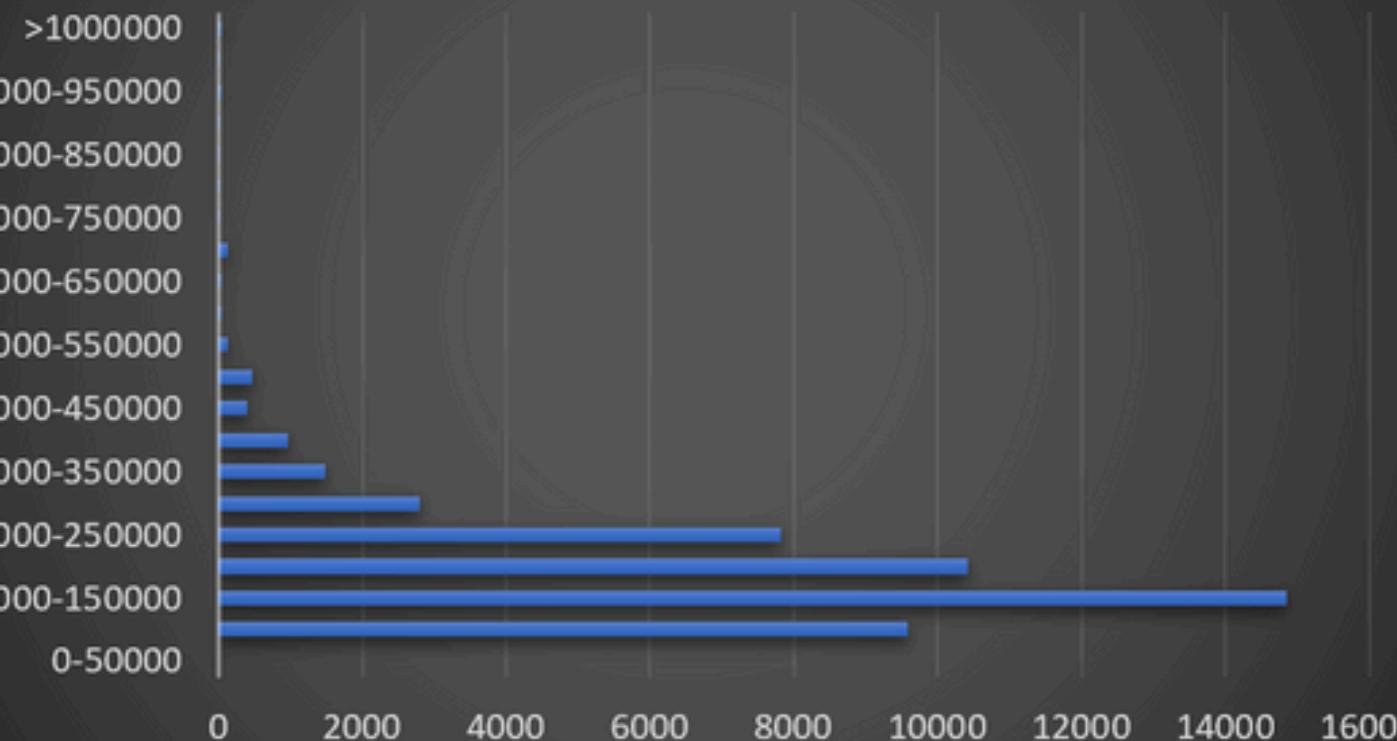
Category	Defaulter	Non Defaulter	Number of applicants
Male	1762	15412	17174
Female	2264	30559	32823
XNA	0	2	2
Total	4026	45973	49997



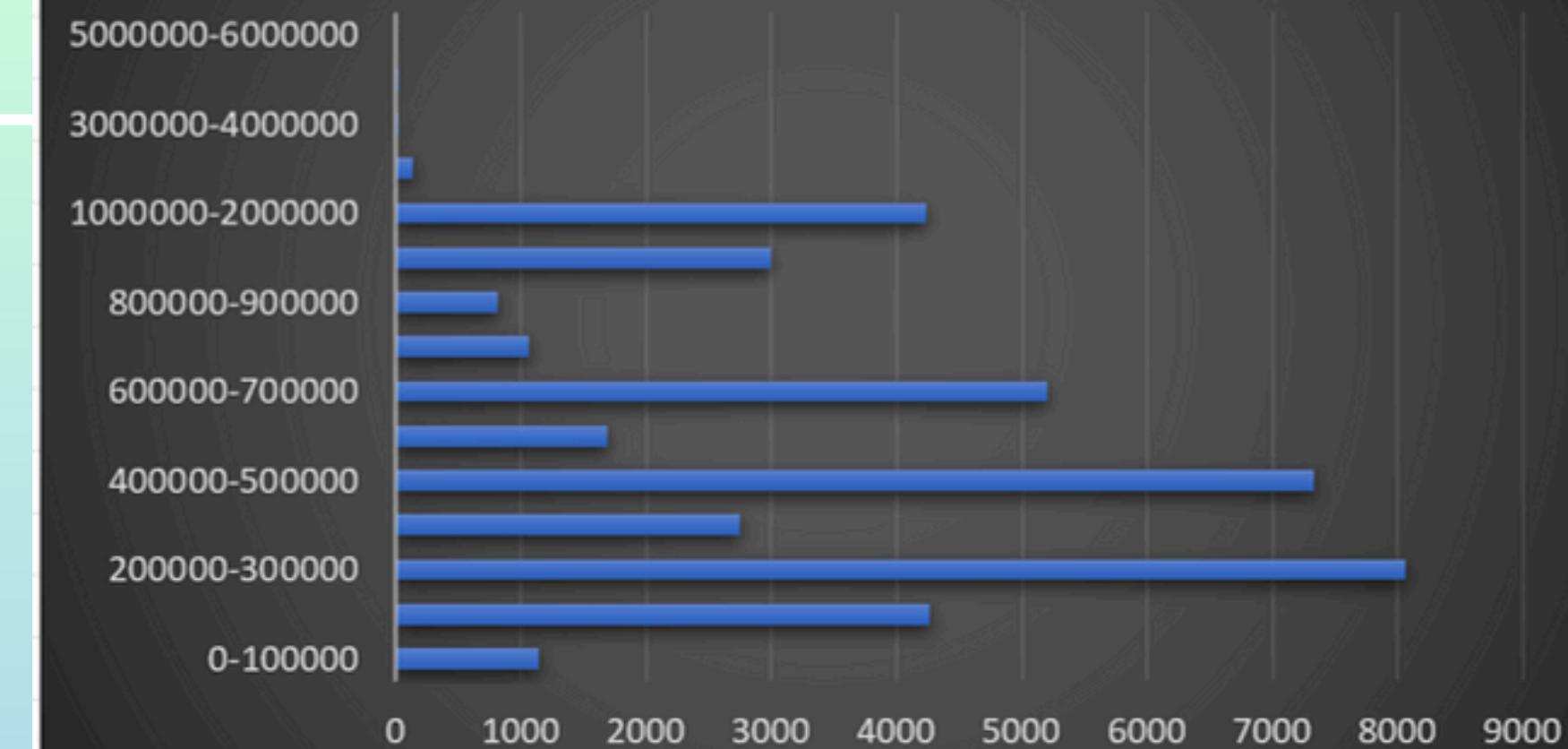
# Segment Univariate

AMT_INCOME_TOTAL	Income_Rang	Total Applicant	Defaulters	Non-Defaulters
202500	0-50000	0 ↓	63 ↓	741
270000	50000-100000	9588 ➔	782 ➔	8806
67500	100000-150000	14852 ↑	1298 ↑	13554
135000	150000-200000	10408 ↑	890 ↑	9518
121500	200000-250000	7818 ➔	576 ➔	7242
99000	250000-300000	2788 ↓	188 ↓	2600
171000	300000-350000	1481 ↓	83 ↓	1398
360000	350000-400000	957 ↓	48 ↓	909
112500	400000-450000	393 ↓	30 ↓	363
135000	450000-500000	456 ↓	37 ↓	419
112500	500000-550000	124 ↓	9 ↓	115
38419.155	550000-600000	43 ↓	5 ↓	38
67500	600000-650000	40 ↓	1 ↓	39
225000	650000-700000	117 ↓	8 ↓	109
189000	700000-750000	22 ↓	1 ↓	21
157500	750000-800000	11 ↓	0 ↓	11
108000	800000-850000	21 ↓	2 ↓	19
81000	850000-900000	5 ↓	0 ↓	5
112500	900000-950000	30 ↓	2 ↓	28
90000	950000-1000000	1 ↓	0 ↓	0
135000	>1000000	40 ↓	3 ↓	0
202500				
450000				36
83250				
135000				

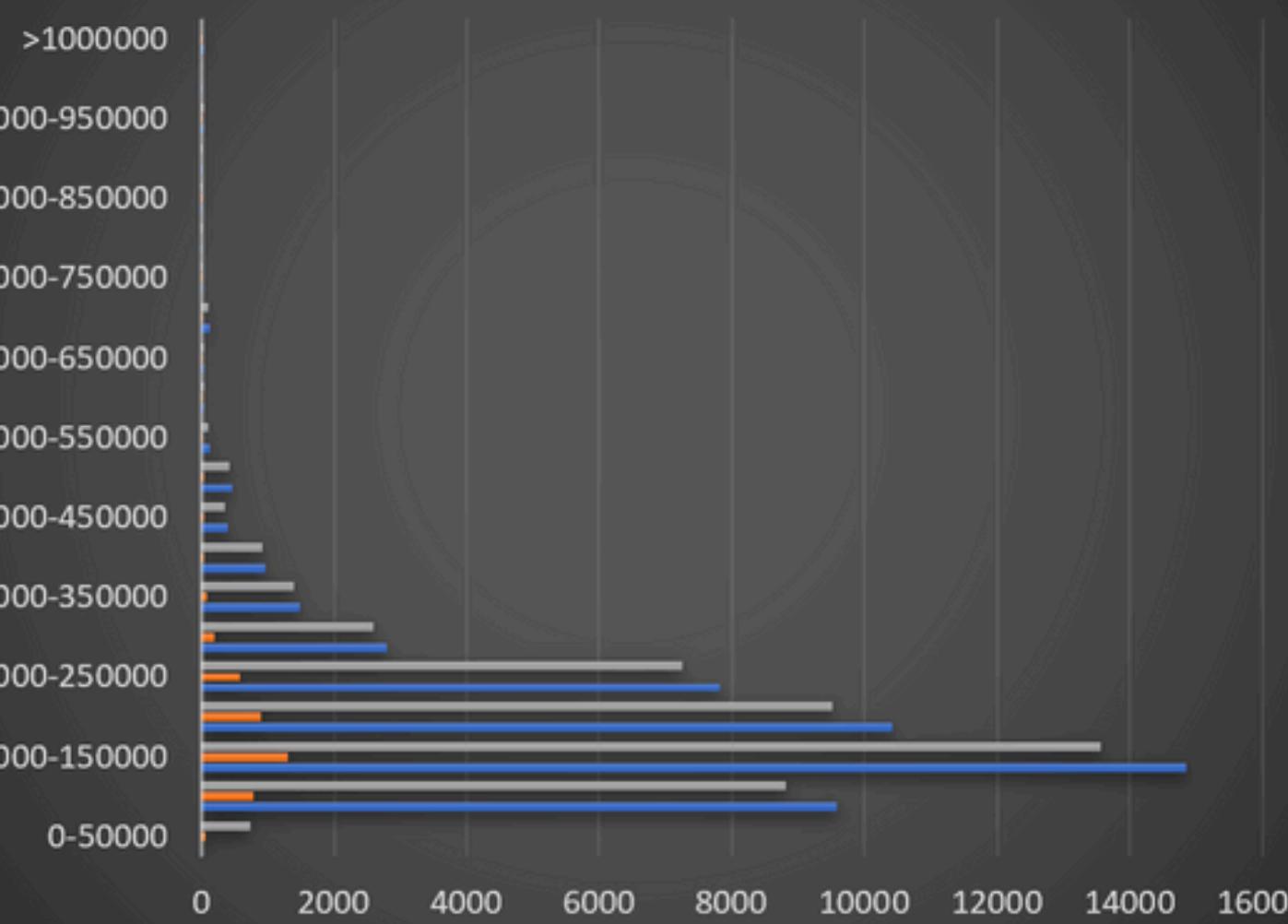
### Income Range vs applicants



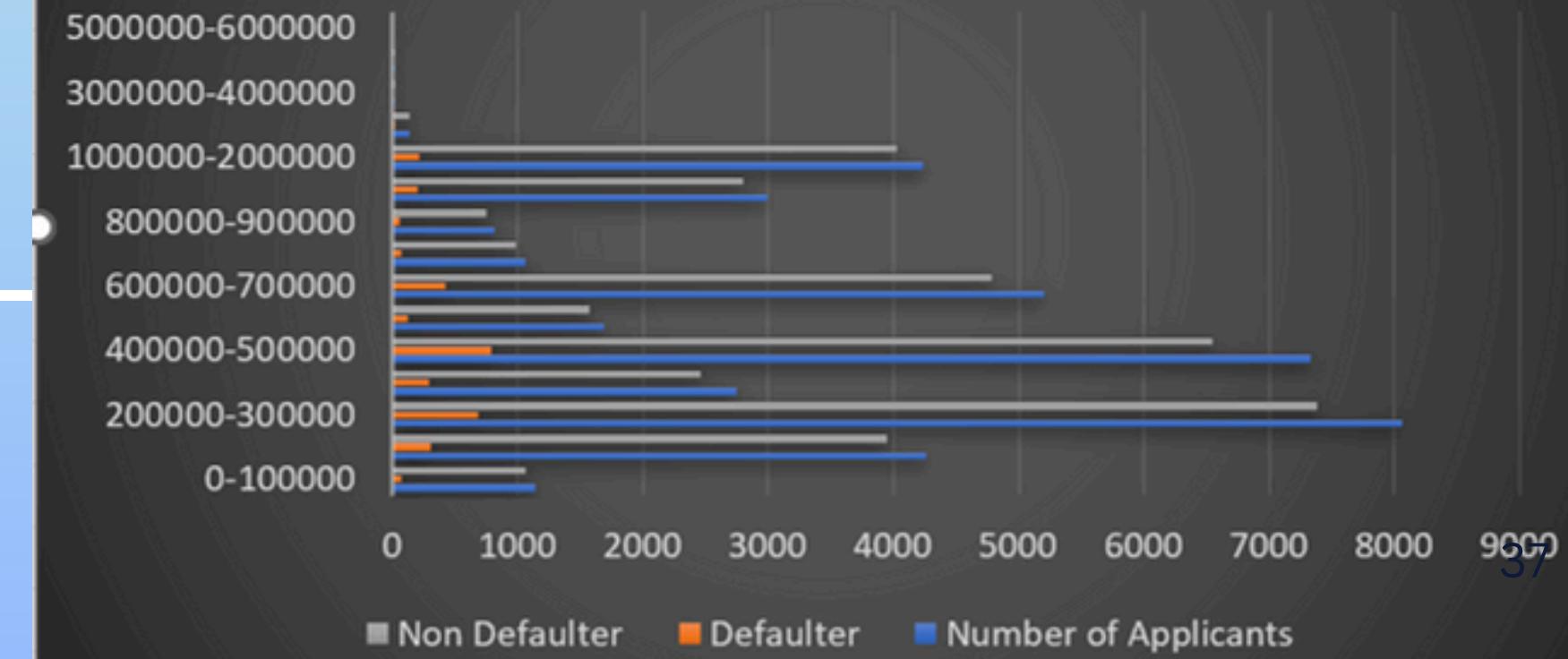
### Good Price vs Number of Applicants



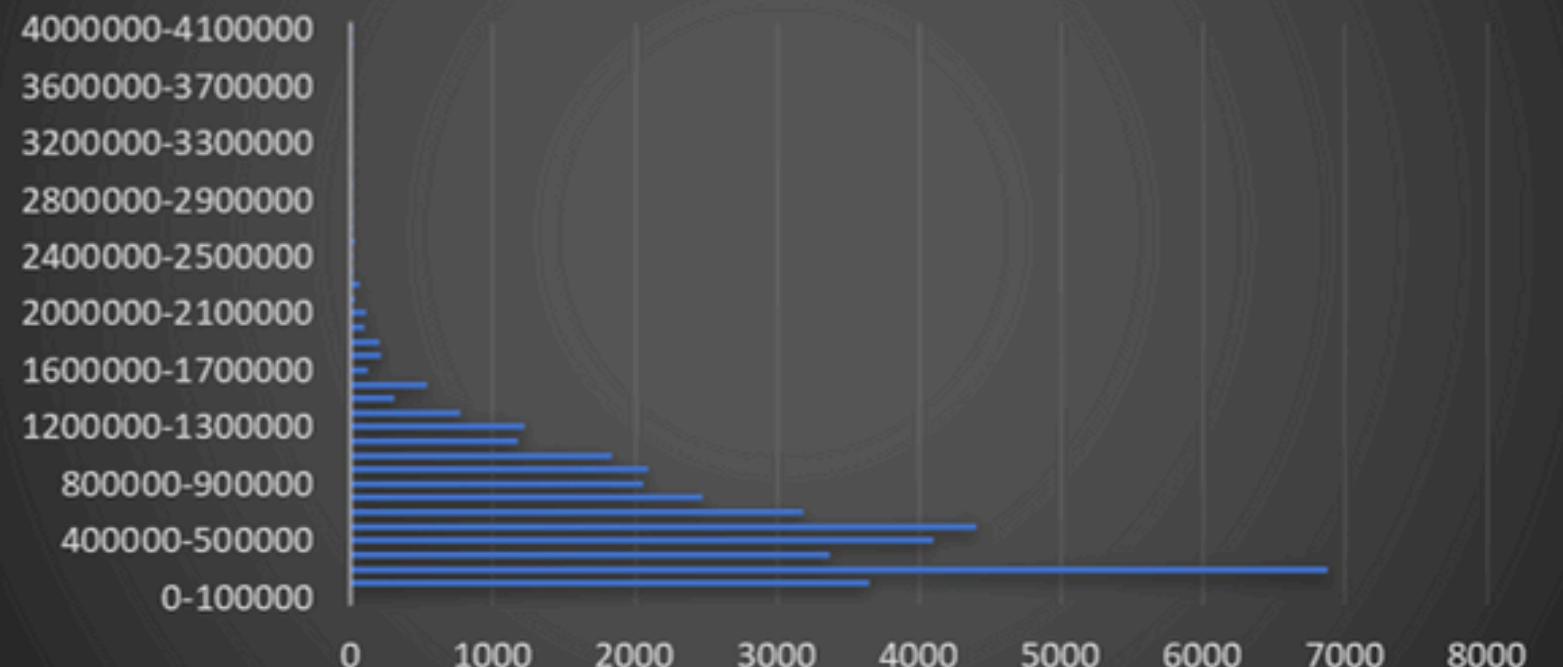
### Income range vs Total applicants vs Defaulters vs Non-Defaulters



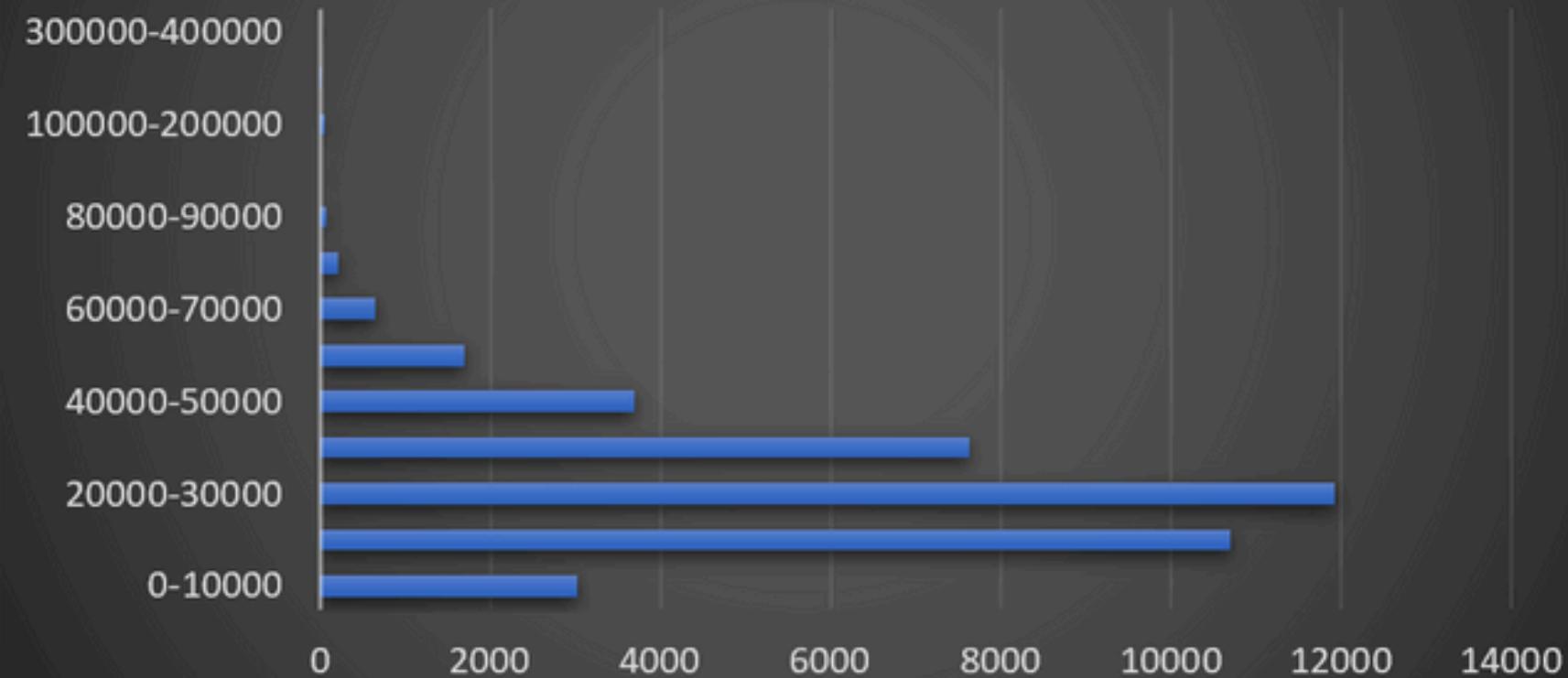
### Good price vs defaulter & non defaulter



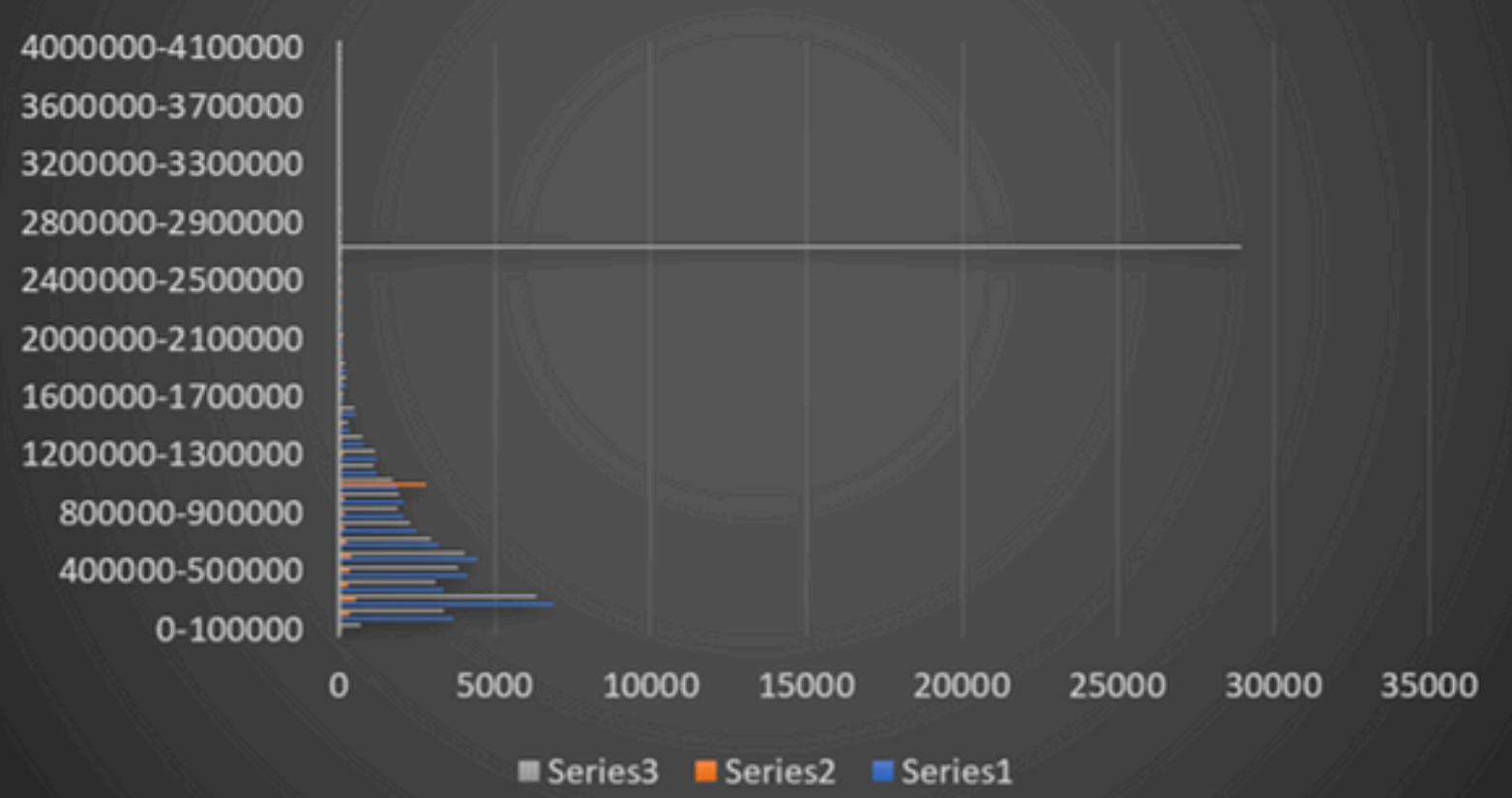
**Income Range vs Number of Applicants**



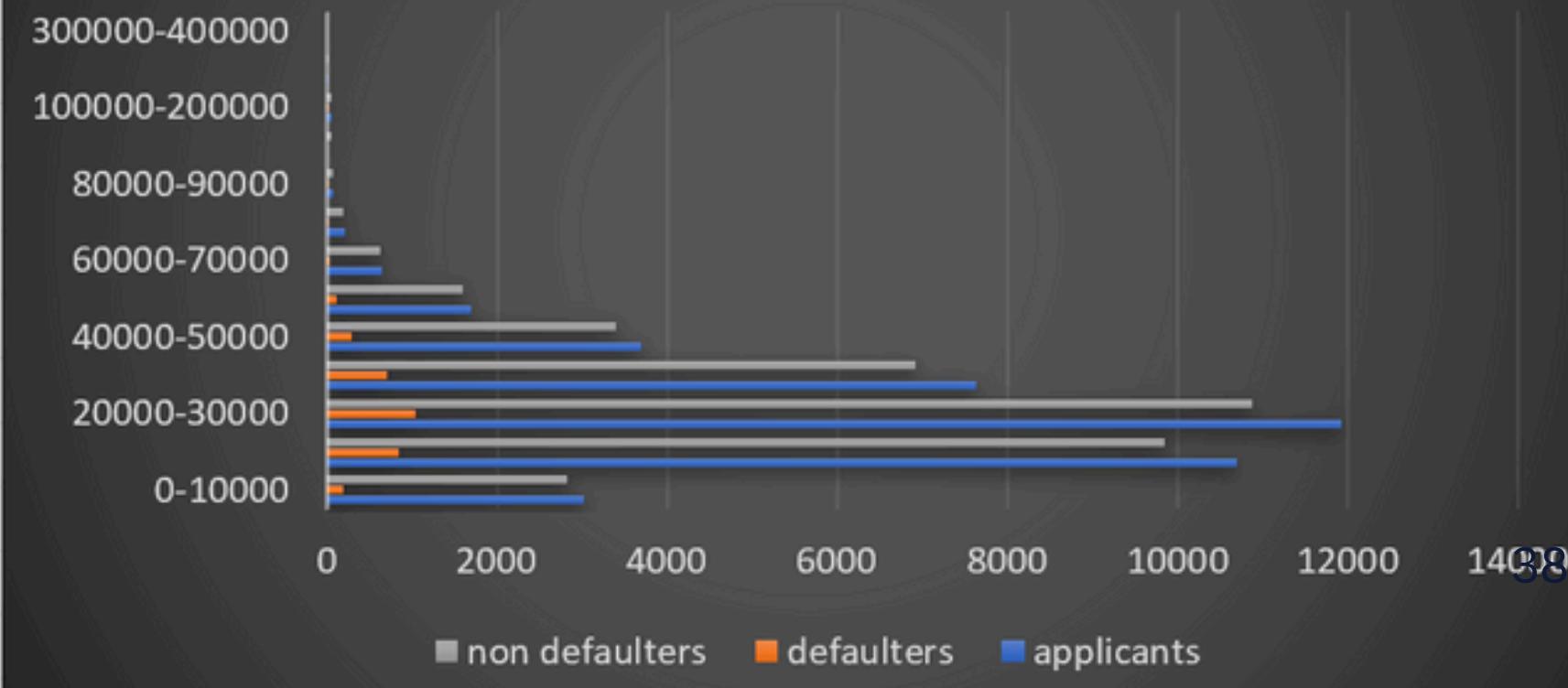
**Income range vs number of applicants**



**Income range vs Defaulter & non defaulter**



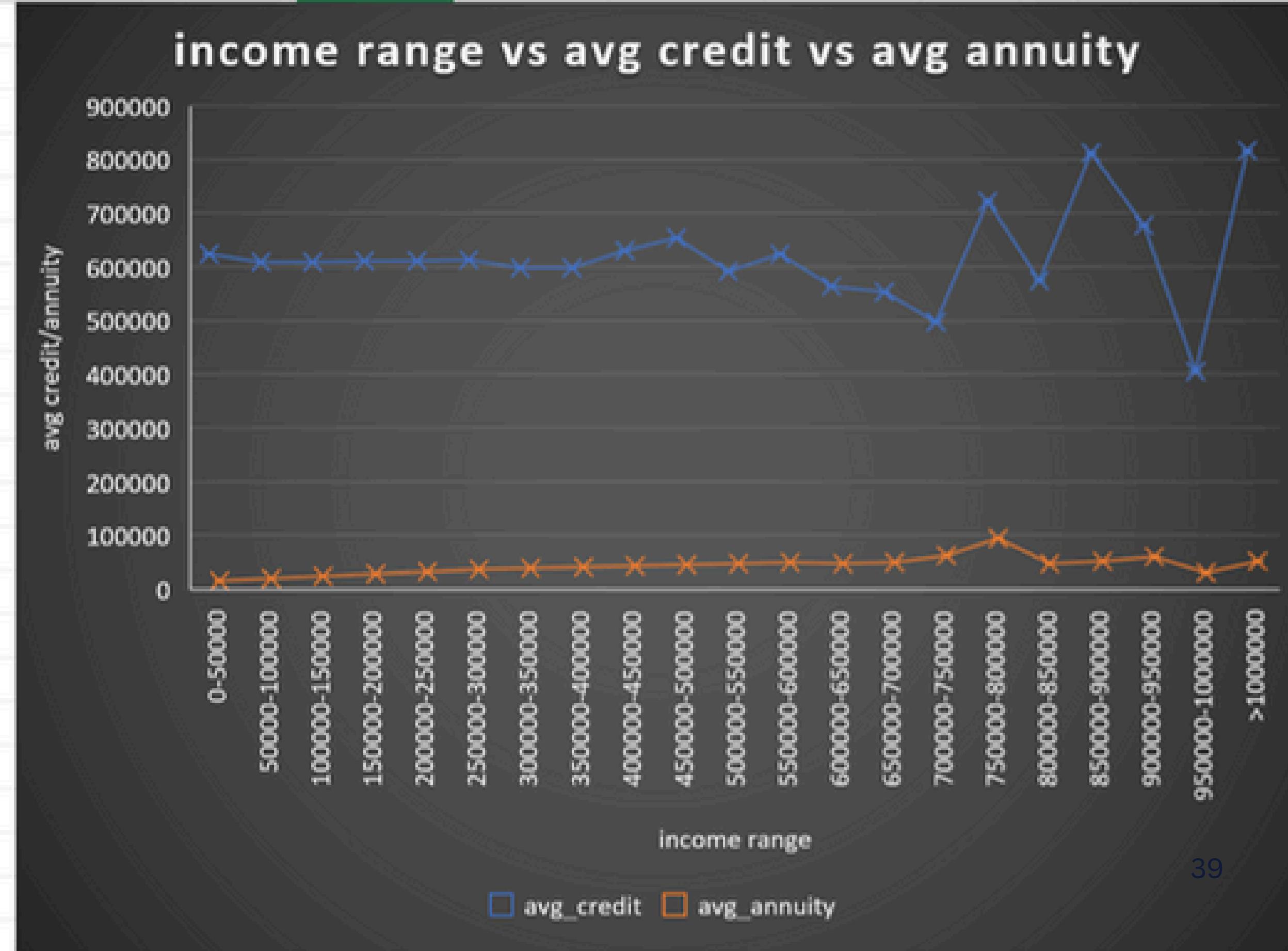
**Income Range vs defaulter & non defaulter**



■ non defaulters   ■ defaulters   ■ applicants

# Bivariate

Income_Rang	avg_credit	avg_annuity
0-50000	623794.649	14151.99685
50000-100000	607827.981	18525.16367
100000-150000	608868.601	23811.0152
150000-200000	610981.252	28601.75194
200000-250000	609620.778	32875.69724
250000-300000	613038.438	36131.96455
300000-350000	598461.257	39045.61974
350000-400000	598587.842	41233.91502
400000-450000	629081.661	43487.16771
450000-500000	653128.563	45862.5625
500000-550000	592347.375	46678.78125
550000-600000	623037.838	48854.38235
600000-650000	564380.016	48108.65625
650000-700000	552699.574	49987.29255
700000-750000	496989	62929.18421
750000-800000	721121.143	95227.71429
800000-850000	574122.214	47424.21429
850000-900000	812246.4	51605.1
900000-950000	677430.173	59482.38462
950000-1000000	405000	30073.5
>1000000	815365.636	51568.22727



# Task 5

## IDENTIFY TOP CORRELATIONS FOR DIFFERENT SCENARIOS

To gain insights into the factors driving loan defaults, you can segment the dataset based on different scenarios (e.g., clients with payment difficulties vs. all other clients) and assess the correlations between variables and the target variable (e.g., loan default). Here's how to conduct this analysis in Excel:

### 1. SEGMENTING THE DATASET

To analyze the data by scenario, first, divide your dataset into relevant segments (e.g., clients with payment difficulties vs. all other clients).

- STEPS:
  - A. USE FILTERS TO SEGMENT THE DATA BASED ON THE RELEVANT SCENARIO (E.G., FILTER OUT CLIENTS WITH PAYMENT DIFFICULTIES OR CREATE A NEW COLUMN INDICATING PAYMENT DIFFICULTY).
  - B. ALTERNATIVELY, USE IF FORMULAS TO CREATE A NEW COLUMN THAT FLAGS CLIENTS IN DIFFERENT SCENARIOS.
    - EXAMPLE: =IF(PAYMENTSTATUS="DIFFICULT", "PAYMENT DIFFICULTY", "No DIFFICULTY").

### 2. CALCULATE CORRELATION COEFFICIENTS WITHIN EACH SEGMENT

Once the dataset is segmented, you can calculate the correlation coefficients between the variables and the target variable (e.g., loan default).

- STEPS:
  - A. SELECT THE SEGMENT YOU WANT TO ANALYZE (E.G., CLIENTS WITH PAYMENT DIFFICULTIES).

## B. USE THE CORREL FUNCTION TO CALCULATE THE CORRELATION COEFFICIENT BETWEEN THE VARIABLES AND THE TARGET VARIABLE:

- EXAMPLE: =CORREL(A2:A100, B2:B100) WHERE A2:A100 IS THE RANGE FOR THE INDEPENDENT VARIABLE (E.G., LOAN AMOUNT) AND B2:B100 IS THE RANGE FOR THE TARGET VARIABLE (E.G., LOAN DEFAULT).
  - REPEAT THE PROCESS FOR ALL RELEVANT VARIABLES IN THE SEGMENT.
  - RANK THE CORRELATION COEFFICIENTS TO IDENTIFY THE STRONGEST CORRELATIONS (USE SORT OR CONDITIONAL FORMATTING TO HIGHLIGHT THE TOP CORRELATIONS).

## 3. RANK AND IDENTIFY TOP INDICATORS OF LOAN DEFAULT

AFTER CALCULATING THE CORRELATION COEFFICIENTS FOR EACH VARIABLE WITHIN THE SEGMENT, RANK THEM TO IDENTIFY WHICH VARIABLES HAVE THE STRONGEST RELATIONSHIP WITH THE TARGET VARIABLE (E.G., LOAN DEFAULT).

- STEPS:
  - CREATE A TABLE THAT LISTS THE VARIABLES AND THEIR RESPECTIVE CORRELATION COEFFICIENTS FOR EACH SEGMENT.
  - SORT THE CORRELATION COEFFICIENTS IN DESCENDING ORDER (HIGHEST TO LOWEST) TO IDENTIFY THE TOP INDICATORS.
  - HIGHLIGHT THE TOP CORRELATIONS USING CONDITIONAL FORMATTING OR COLOR CODING.

## 4. VISUALIZE CORRELATIONS USING A HEATMAP OR CORRELATION MATRIX

A VISUAL REPRESENTATION OF THE CORRELATIONS WILL HELP QUICKLY IDENTIFY WHICH VARIABLES ARE MOST STRONGLY CORRELATED WITH THE TARGET VARIABLE.

- STEPS:
  - USE CONDITIONAL FORMATTING TO COLOR-CODE THE CORRELATION MATRIX:
    - SELECT THE CORRELATION COEFFICIENTS IN A TABLE.

- Go to Home > Conditional Formatting > Color Scales and choose a color scale.
- Alternatively, create a correlation matrix by listing the variables in both rows and columns and filling the cells with the correlation coefficients.
- Use a heatmap to visually highlight the strongest correlations. Excel offers options to format cells with color gradients based on correlation values.

# Defaulters

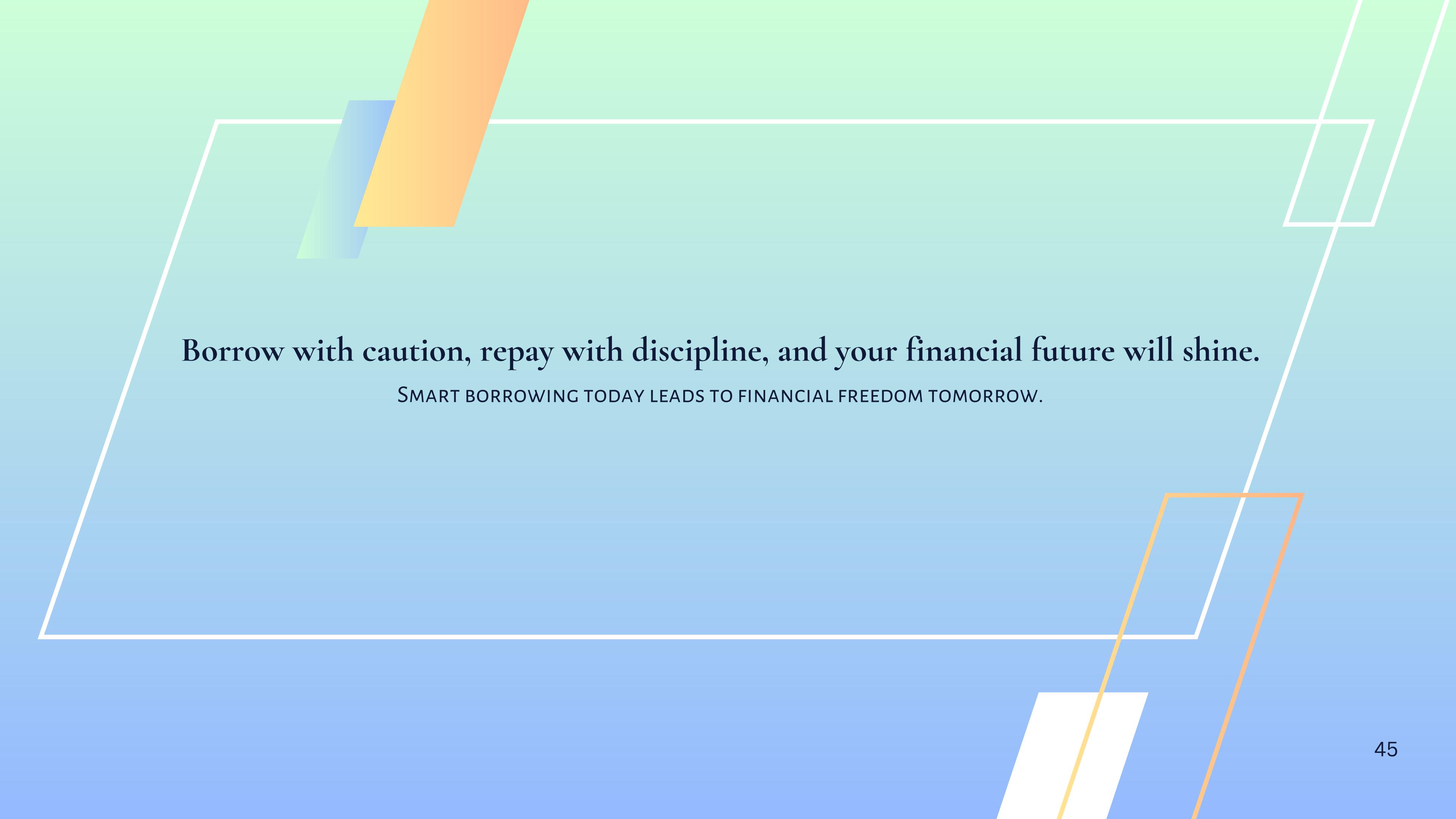
	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE	Year_BIRTH	Year_EMPLOYED	Year_REGISTERED	Days_ID	
CNT_CHILDREN	100.00%	0.92%	0.54%	-0.50%	-0.14%		-2.32%	-34.28%	-24.32%	-18.14%	3.84%
AMT_INCOME_TOTAL	0.92%	100.00%	6.09%	0.00%	0.28%		2.49%	-1.64%	-2.88%	-0.85%	-0.44%
AMT_CREDIT	0.54%	6.09%	100.00%	0.33%	0.58%		9.55%	5.02%	-7.27%	-0.47%	0.10%
AMT_ANNUITY	-0.50%	0.00%	0.33%	100.00%	77.45%		0.42%	0.61%	0.45%	1.25%	0.45%
AMT_GOODS_PRICE	-0.14%	0.28%	0.58%	77.45%	100.00%		0.61%	0.40%	0.56%	0.93%	0.38%
REGION_POPULATION_RELATIVE	-2.32%	2.49%	9.55%	0.42%	0.61%	100.00%	3.18%	-0.01%	5.84%	0.19%	
Year_BIRTH	-34.28%	-1.64%	5.02%	0.61%	0.40%		3.18%	100.00%	61.93%	33.51%	24.53%
Year_EMPLOYED	-24.32%	-2.88%	-7.27%	0.45%	0.56%		-0.01%	61.93%	100.00%	20.55%	26.28%
Year_REGISTERED	-18.14%	-0.85%	-0.47%	1.25%	0.93%		5.84%	33.51%	20.55%	100.00%	8.78%
Days_ID	3.84%	-0.44%	0.10%	0.45%	0.38%		0.19%	24.53%	26.28%	8.78%	100.00%

# Non Defaulters

	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE	Year_BIRTH	Year_EMPLOYED	Year_REGISTERED	Days_ID	
CNT_CHILDREN	100.00%	3.57%	0.60%	-0.15%	0.26%		-2.15%	-34.92%	-24.69%	-18.35%	3.76%
AMT_INCOME_TOTAL	3.57%	100.00%	37.28%	0.28%	0.67%		17.64%	-8.05%	-16.29%	-6.95%	-4.20%
AMT_CREDIT	0.60%	37.28%	100.00%	0.46%	0.72%		9.56%	4.28%	-7.84%	-0.82%	-0.20%
AMT_ANNUITY	-0.15%	0.28%	0.46%	100.00%	77.42%		0.31%	0.68%	0.44%	1.06%	0.63%
AMT_GOODS_PRICE	0.26%	0.67%	0.72%	77.42%	100.00%		0.49%	0.40%	0.49%	0.90%	0.51%
REGION_POPULATION_RELATIVE	-2.15%	17.64%	9.56%	0.31%	0.49%	100.00%	-0.26%	-0.26%	-0.26%	5.73%	-0.01%
Year_BIRTH	-34.92%	-8.05%	4.28%	0.68%	0.40%		2.93%	100.00%	62.08%	33.69%	24.39%
Year_EMPLOYED	-24.69%	-16.29%	-7.84%	0.44%	0.49%		-0.26%	62.08%	100.00%	20.56%	26.50%
Year_REGISTERED	-18.35%	-6.95%	-0.82%	1.06%	0.90%		5.73%	33.69%	20.56%	100.00%	8.87%
Days_ID	3.76%	-4.20%	-0.20%	0.63%	0.51%		-0.01%	24.39%	26.50%	8.87%	100.00%

# Top Correlations

TOP CORRELATION FOR DEFAULTERS			TOP CORRELATIONS FOR NON-DEFAULTERS		
VARIABLE 1	VARIABLE 2	CORRELATION	VARIABLE 1	VARIABLE 2	CORRELATION
AMT_ANNUITY	AMT_GOODS_PRICE	77.45%	AMT_ANNUITY	AMT_GOODS_PRICE	77.42%
AMT_GOODS_PRICE	AMT_ANNUITY	77.45%	AMT_GOODS_PRICE	AMT_ANNUITY	77.42%
Year BIRTH	Year EMPLOYED	61.93%	Year BIRTH	Year EMPLOYED	62.08%
Year EMPLOYED	Year BIRTH	61.93%	Year EMPLOYED	Year BIRTH	62.08%
Year REGISTERED	Year BIRTH	33.51%	AMT_INCOME_TOTAL	AMT_CREDIT	37.28%



**Borrow with caution, repay with discipline, and your financial future will shine.**

SMART BORROWING TODAY LEADS TO FINANCIAL FREEDOM TOMORROW.

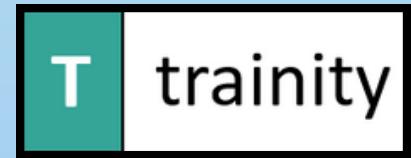


Link to the project  
**CLICK: LINK TO PROJECT DRIVE**



# Ankita Taneja

DATA ANALYST



# Your Reviews

WHAT CAN YOU SAY ABOUT MY PROJECTS? SHARE IT HERE!



# Contact Us

-  GURUGRAM
-  0-7011-334-048
-  VENUSGIRLATWORK@GMAIL.COM
-  @IG\_SHE\_HAS\_NO\_IDEA
-  [WWW.LINKEDIN.COM/IN/ANKITA-TANEJA/](https://www.linkedin.com/in/ankita-taneja/)





The background features several abstract geometric shapes. In the top left, there's a white right-angled triangle pointing downwards. Next to it is a trapezoid with orange sides and a light green interior. Above the trapezoid is a blue parallelogram. In the top right, there's a yellow parallelogram and a smaller orange parallelogram. At the bottom right, there's a large orange parallelogram containing a white rectangle, and below it is a white right-angled triangle pointing upwards. The overall background is a blue gradient.

# Thank You!