NAME-AKSHAY ANAND

PROJECT DESCRIPTION – This project is about a case study relayed to a Bank Loan. We have to carry out an EDA (Exploratory Data Analysis).
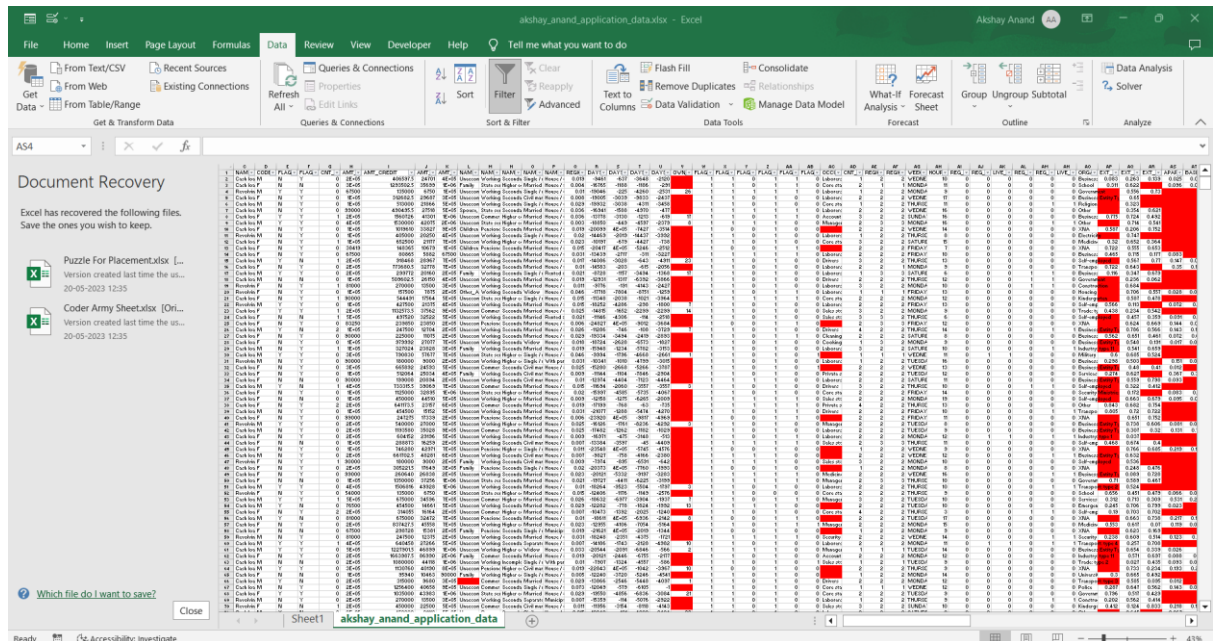
- Present the overall approach of the **analysis**. Mention the problem statement and the analysis approach briefly

SINCE THE DATASET HAD A LOT OF MISSING VALUES WE FILLED THE BLANK VALUES WITH MEAN OR MEDIAN DEPENDING ON OUTLIERS AND THEN USED PIVOT TABLES,FILTERS AND OTHER FEATURES OF EXCEL TO PERFORM FURTHER ANALYSIS.
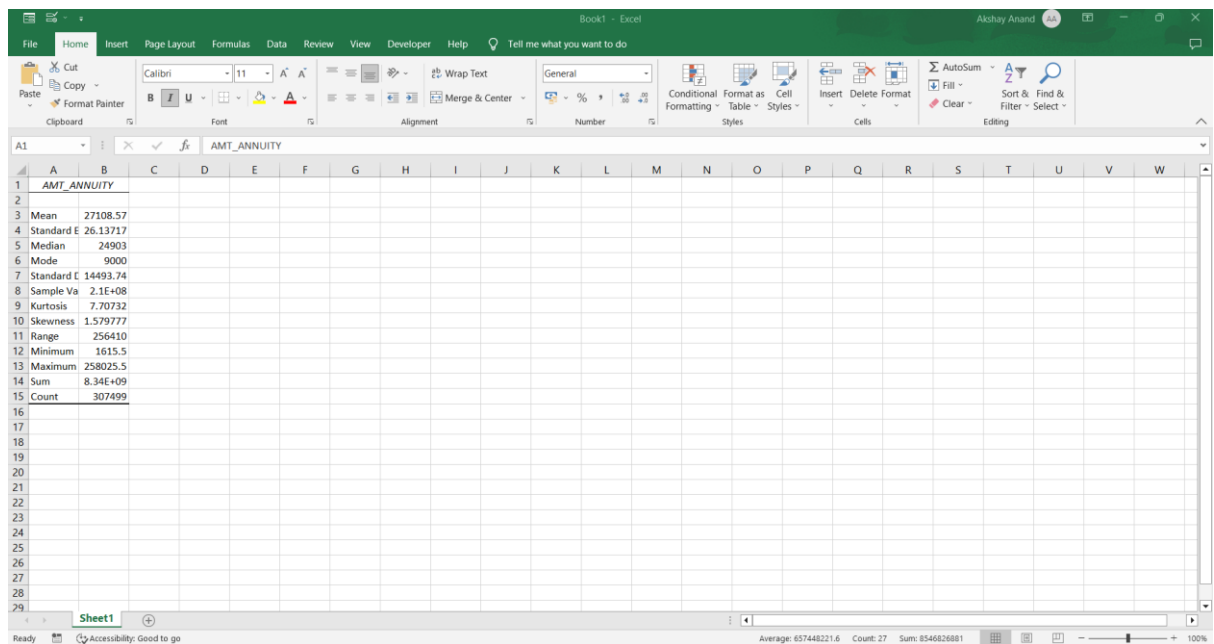
TECH STACK USED: MS EXCEL

- **Indentify** the missing data and use appropriate method to deal with it. (Remove columns/or replace it with an appropriate value)
  *Hint: Note that in EDA, since it is not necessary to replace the missing value, but if you have to replace the missing value, what should be the approach. Clearly mention the approach.*

TO FIND OUT BLANK OR MISSING VALUES WE CAN USE FILTERS AND SELECT BLANK OR USING FIND & REPLACE WE REPLACE ALL BLANK CELLS BY ANY COLOUR.
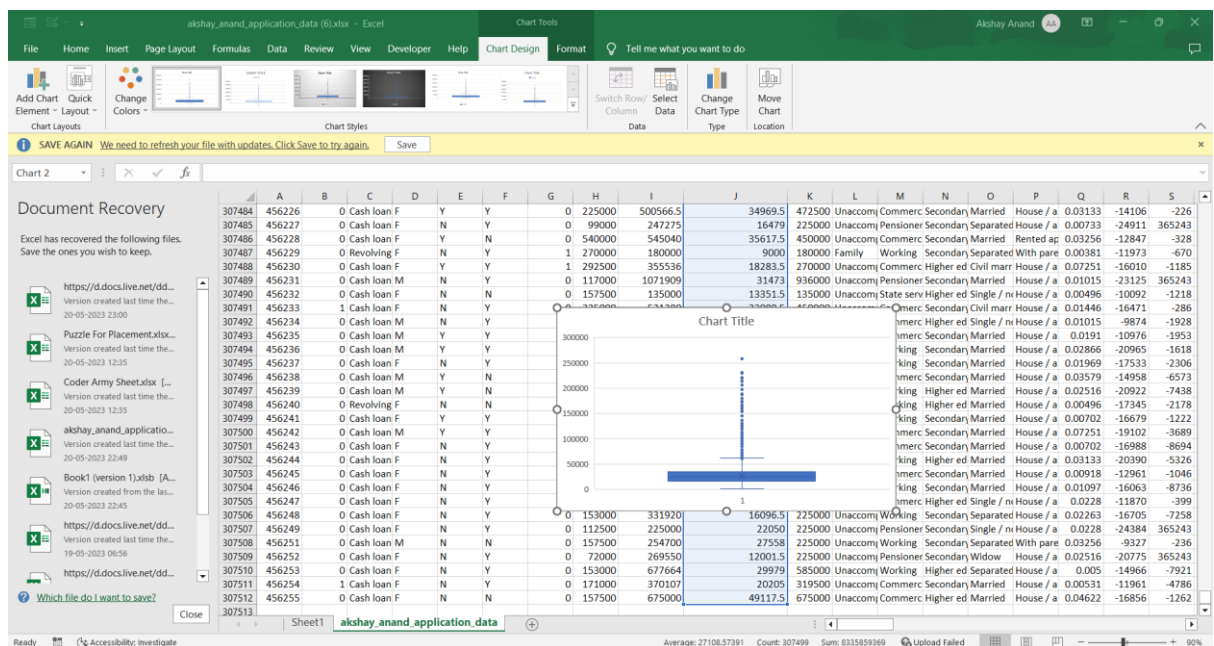


WE CAN USE DESCRIPTIVE STATISTICS BY SELECTING DATA AND THEN DATA ANALYSIS.

SINCE IT IS POSSIBLE THAT DATA HAS OUTLIERS SO WE REPLACE ALL MISSING VALUES BY MEDIAN.

SINCE THE MEDIAN FOR AMT_ANNUITY COLUMN IS 24903 WE REPLACE ALL BLANK VALUES IN THAT COLUMN WITH THIS MEDIAN VALUE.
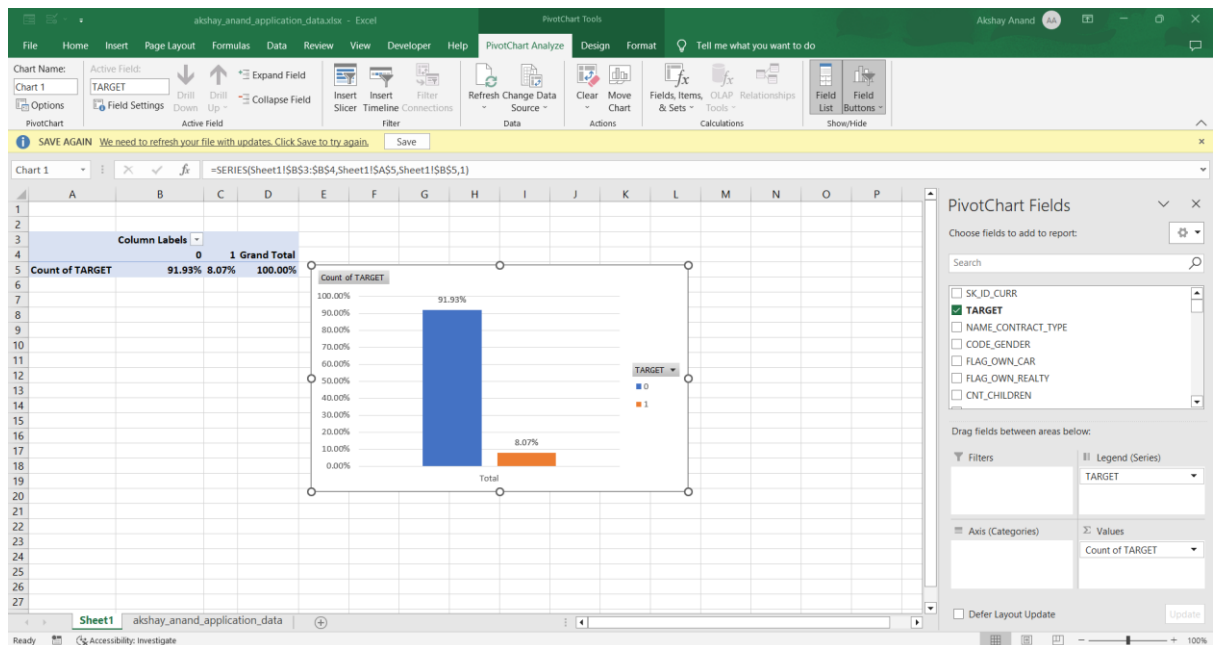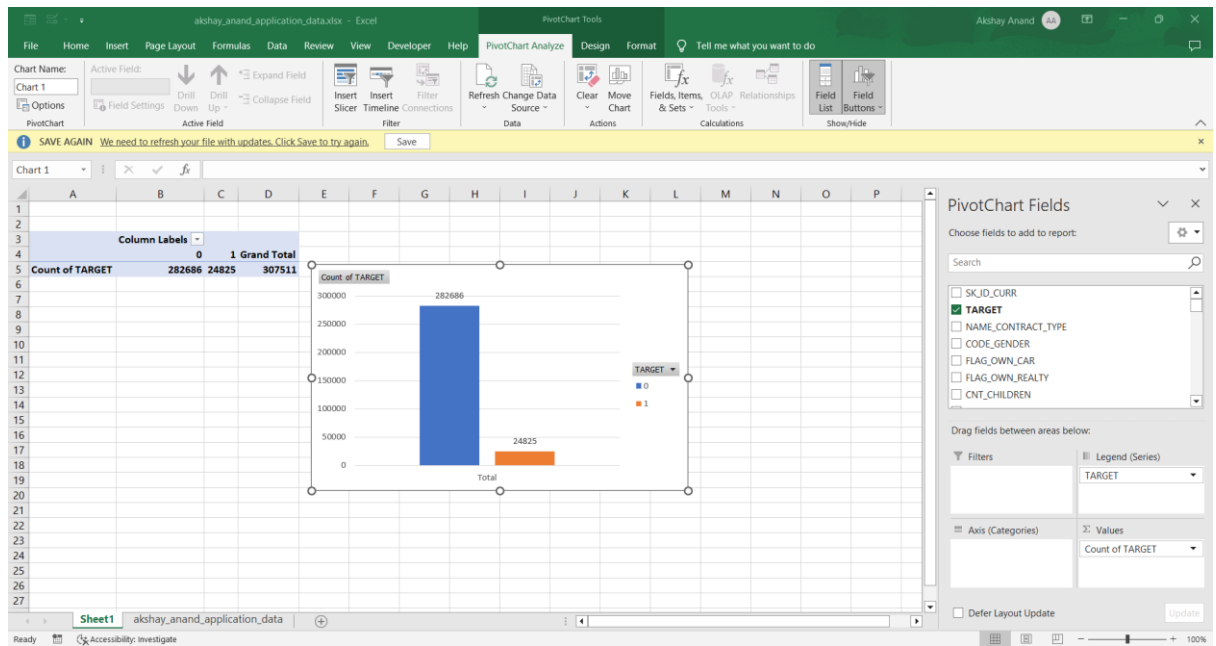


- Identify if there are **outliers** in the dataset. Also, mention why do you think it is an outlier. Again, remember that for this exercise, it is not necessary to remove any data points.
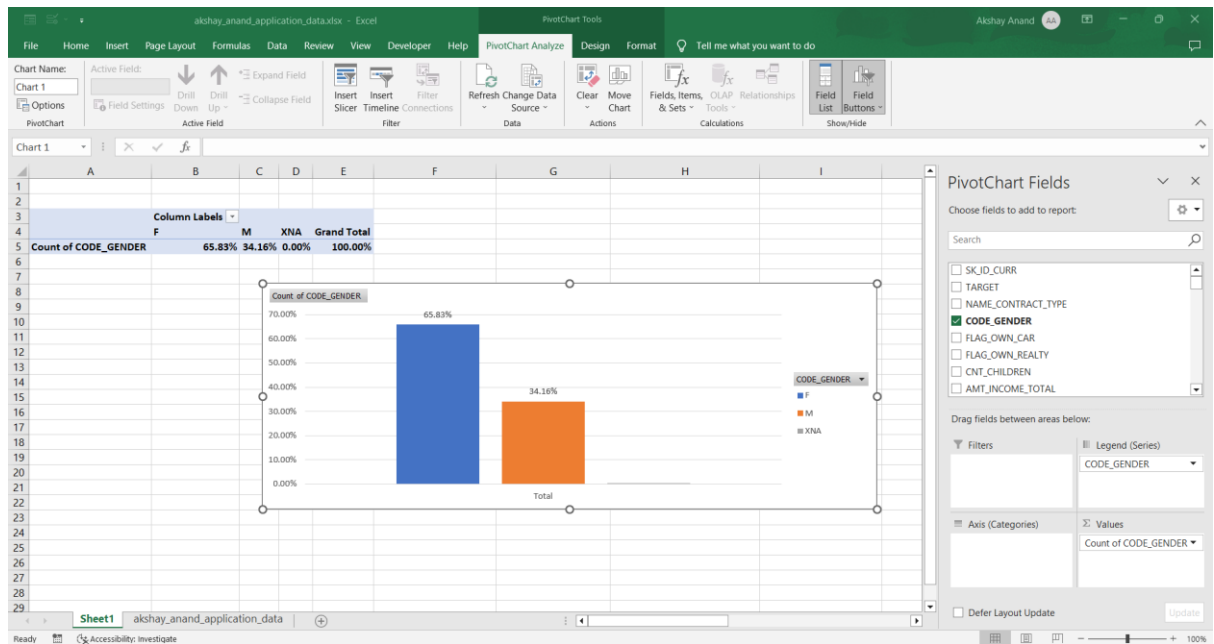
THE POINTS WHICH ARE LOCATED OUTSIDE THE WHISKERS OF THE BOX PLOT ARE OUTLIERS.

ALSO, THE MEAN IS GREATLY AFFECTED BY OUTLIERS BUT MODE AND MEDIAN ARE NOT MUCH AFFECTED BY OUTLIERS SO IF ANY DATA HAS A SIGNIFICANT DIFFERENCE BETWEEN MEAN AND MEDIAN THEN WE CAN SAY THAT THERE ARE OUTLIERS IN THE DATA.

FOR AMT_CREDIT COLUMN

FOR AMT_ANNUITY COLUMN



SO WE SEE THE POINTS WHICH ARE VERY FAR AWAY FROM THE WHISKER PLOT THOSE POINTS ARE OUTLIERS.

- Identify if there is data imbalance in the data. Find the ratio of data imbalance. *Hint: Since there are a lot of columns, you can run your analysis in loops for the appropriate columns and find the insights.*
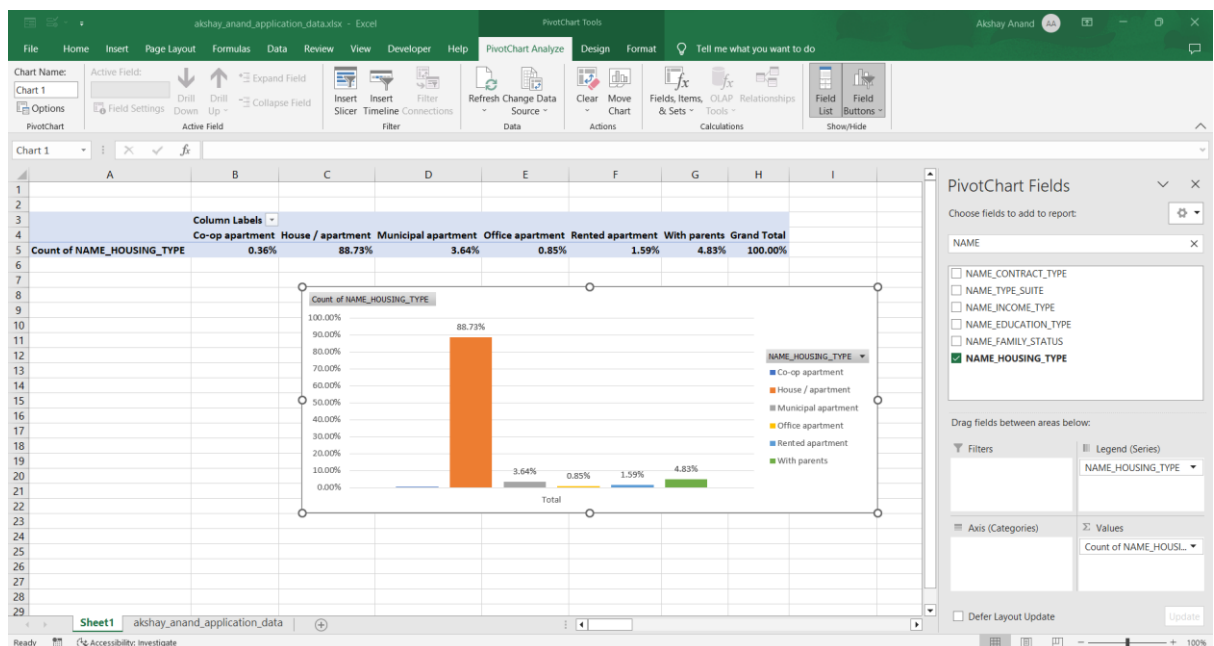
ALMOST 92 % OF VALUES ARE 0 AND 8% ARE 1. HENCE DATA IMBALANCE IS THERE.

FEMALE IS 65%, MALE IS 35%, XNA(BECAUSE OF VERY LESS COUNT) IS 0%. HENCE THERE IS DATA IMBALANCE.
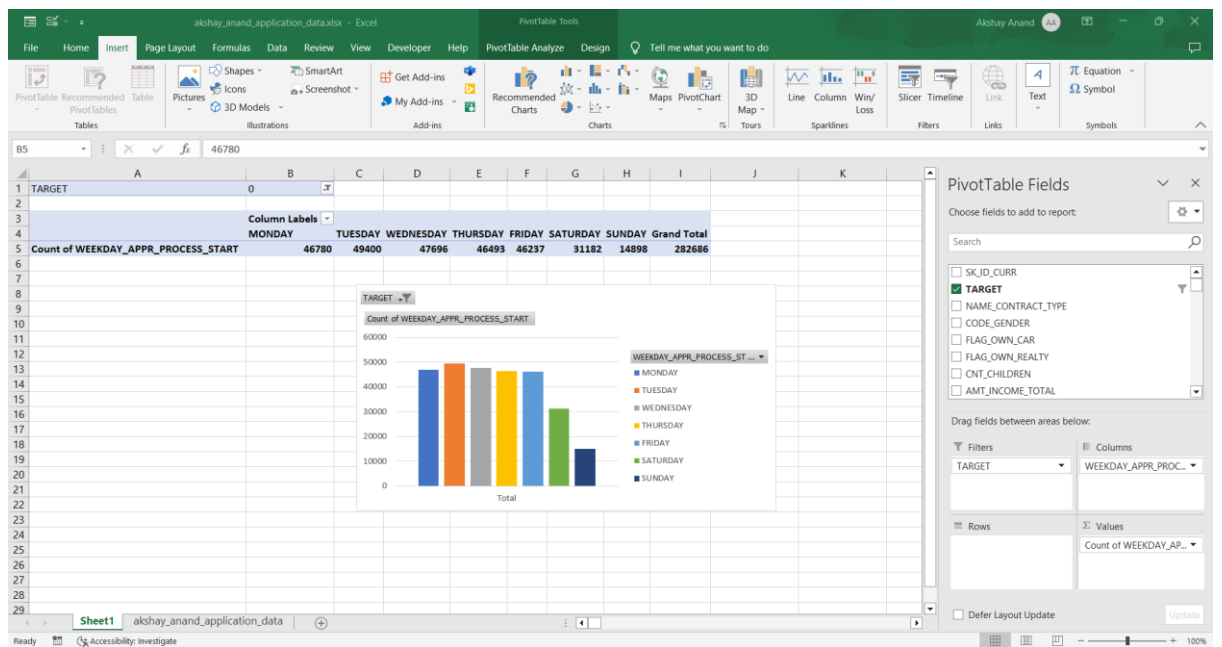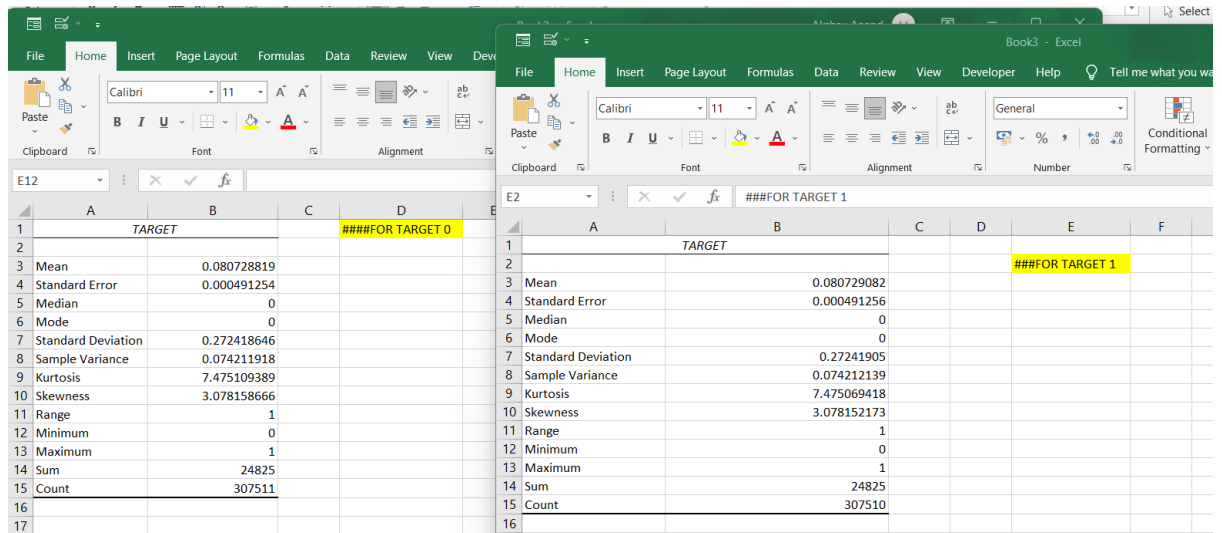


HOUSE/APARTMENT'S SHARE IS 89% AS COMPARED TO OTHER HOUSING TYPES. HENCE THERE IS A DATA IMBALANCE.

- Explain the **results of univariate, segmented univariate, bivariate analysis, etc.** in business terms.
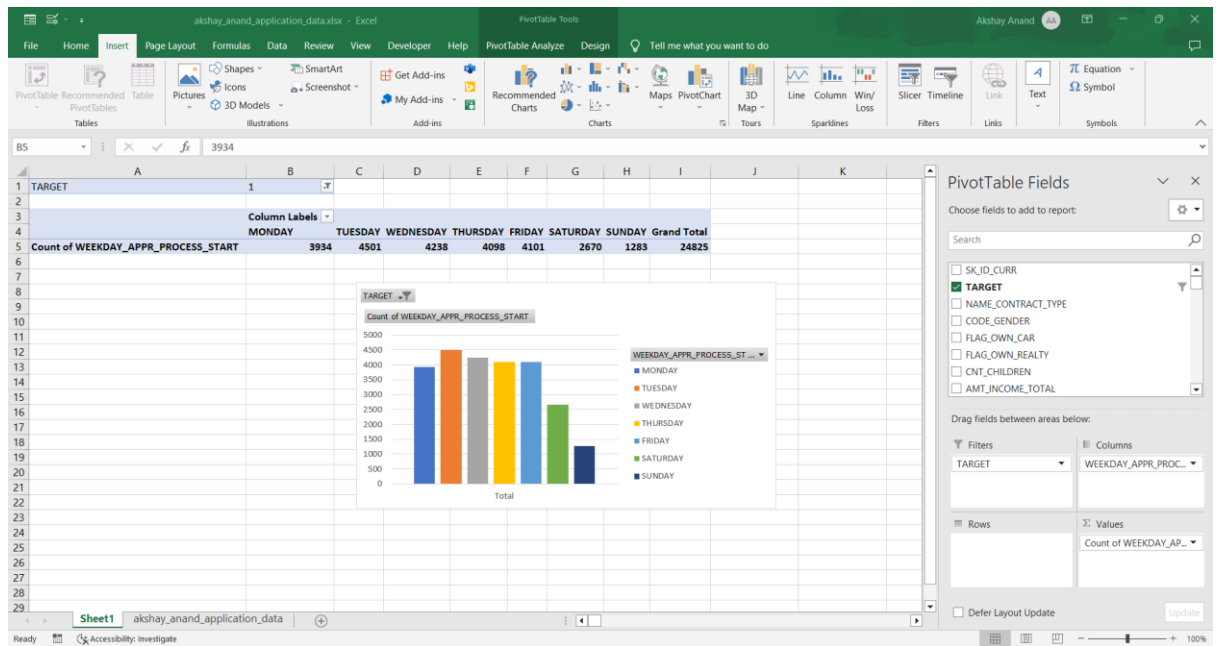
WE CAN DO DESCRIPTIVE UNIVARIATE ANALYSIS BY USING THE DATA ANALYSIS TOOL IN THE DATA TAB.

OR WE CAN USE PIVOT TABLES AND GRAPHS FOR UNIVARIATE AND OTHER TYPES OF ANALYSIS.

First Excel window (TARGET 0):

| | A | B | C | D |
|---|---|---|---|---|
| 1 | | TARGET | | ####FOR TARGET 0 |
| 2 | | | | |
| 3 | Mean | 0.080728819 | | |
| 4 | Standard Error | 0.000491254 | | |
| 5 | Median | 0 | | |
| 6 | Mode | 0 | | |
| 7 | Standard Deviation | 0.272418646 | | |
| 8 | Sample Variance | 0.074211918 | | |
| 9 | Kurtosis | 7.475109389 | | |
| 10 | Skewness | 3.078158666 | | |
| 11 | Range | 1 | | |
| 12 | Minimum | 0 | | |
| 13 | Maximum | 1 | | |
| 14 | Sum | 24825 | | |
| 15 | Count | 307511 | | |
| 16 | | | | |
| 17 | | | | |

Second Excel window (TARGET 1):

| | A | B | C | D | E |
|---|---|---|---|---|---|
| 1 | | TARGET | | | |
| 2 | | | | | ###FOR TARGET 1 |
| 3 | Mean | 0.080729082 | | | |
| 4 | Standard Error | 0.000491256 | | | |
| 5 | Median | 0 | | | |
| 6 | Mode | 0 | | | |
| 7 | Standard Deviation | 0.27241905 | | | |
| 8 | Sample Variance | 0.074212139 | | | |
| 9 | Kurtosis | 7.475069418 | | | |
| 10 | Skewness | 3.078152173 | | | |
| 11 | Range | 1 | | | |
| 12 | Minimum | 0 | | | |
| 13 | Maximum | 1 | | | |
| 14 | Sum | 24825 | | | |
| 15 | Count | 307510 | | | |
| 16 | | | | | |

Pivot table screenshot:

| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| 1 | TARGET | 0 | | | | | | | |
| 3 | | Column Labels | | | | | | | |
| 4 | | MONDAY | TUESDAY | WEDNESDAY | THURSDAY | FRIDAY | SATURDAY | SUNDAY | Grand Total |
| 5 | Count of WEEKDAY_APPR_PROCESS_START | 46780 | 49400 | 47696 | 46493 | 46237 | 31182 | 14898 | 282686 |



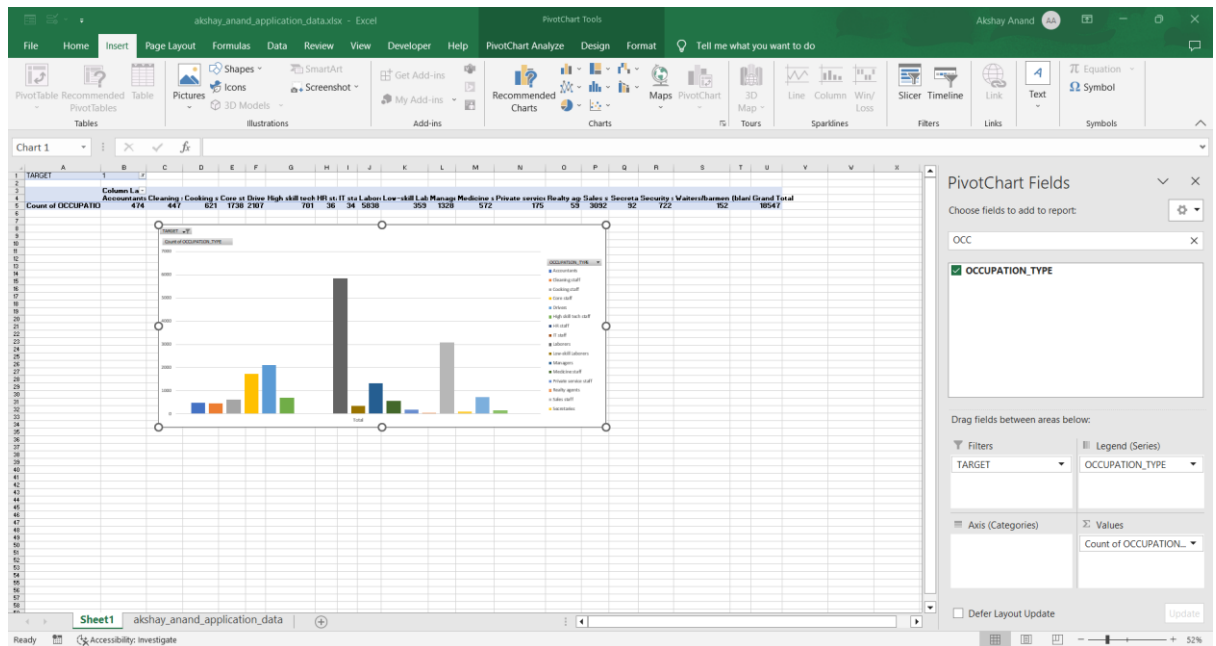ABOVE IS A UNIVARIATE ANALYSIS FOR WEEKDAY_APPR_PROCESS_START FOR TARGET=0.
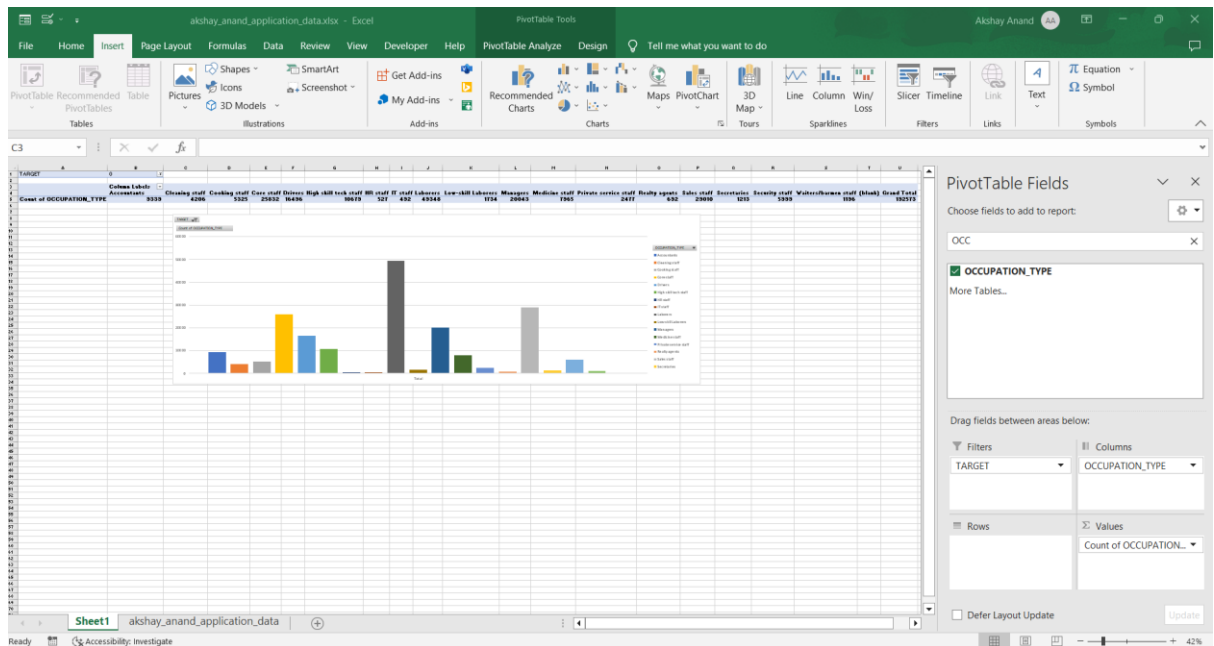
SIMILARLY FOR TARGET=1.

BY THE ABOVE WE CONCLUDE THAT THE APPLICATION STARTING PROCESS IS LESS ON SATURDAY AND SUNDAY IN BOTH CASES OF TARGET = 0 AND 1.

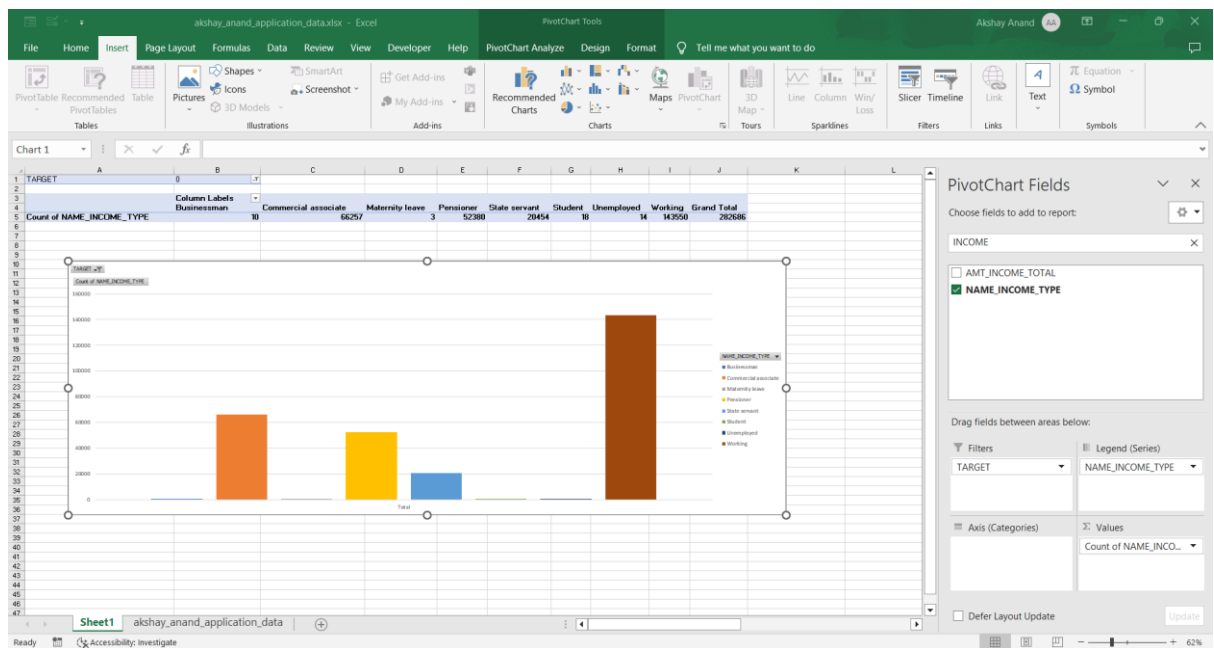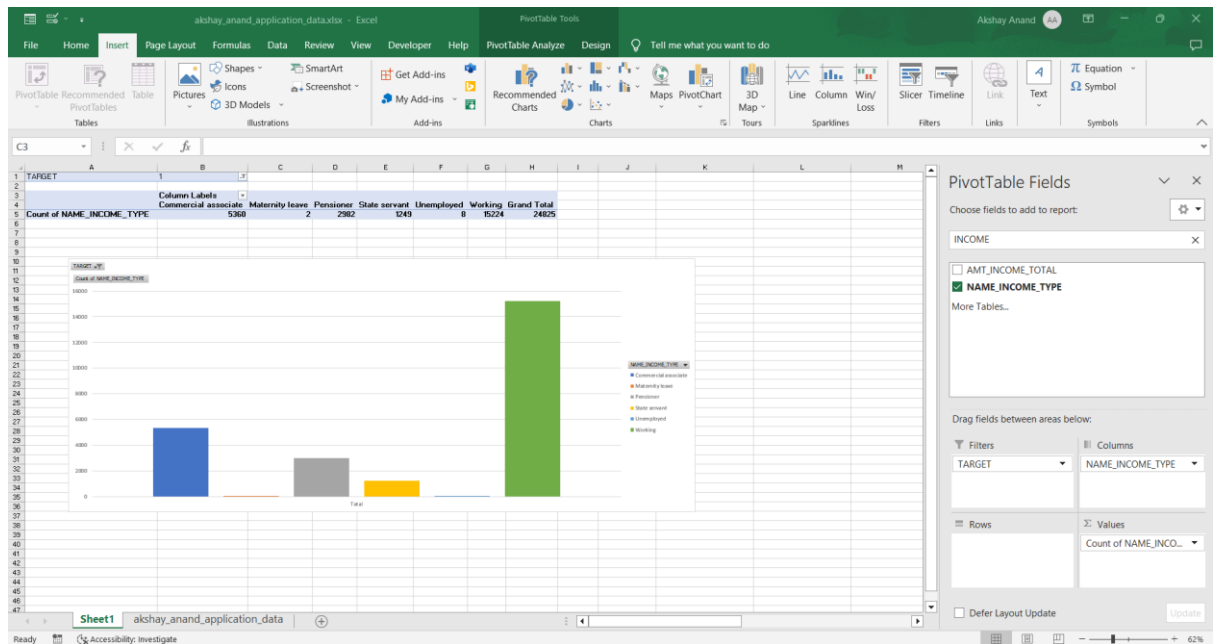NOW FOR OCCUPATION_TYPE

TARGET=1



TARGET=0

IT IS CLEAR FROM THE BAR PLOT THAT CORRESPONDING TO LABOURERS THERE IS A LARGE NO. OF PAYMENT ISSUES IN BOTH CASES OF TARGET= 0 AND 1 WHERE TARGET=0 FOR CLIENTS WITH NO PAYMENT ISSUES AND 1 IS FOR CLIENTS WITH PAYMENT ISSUES.

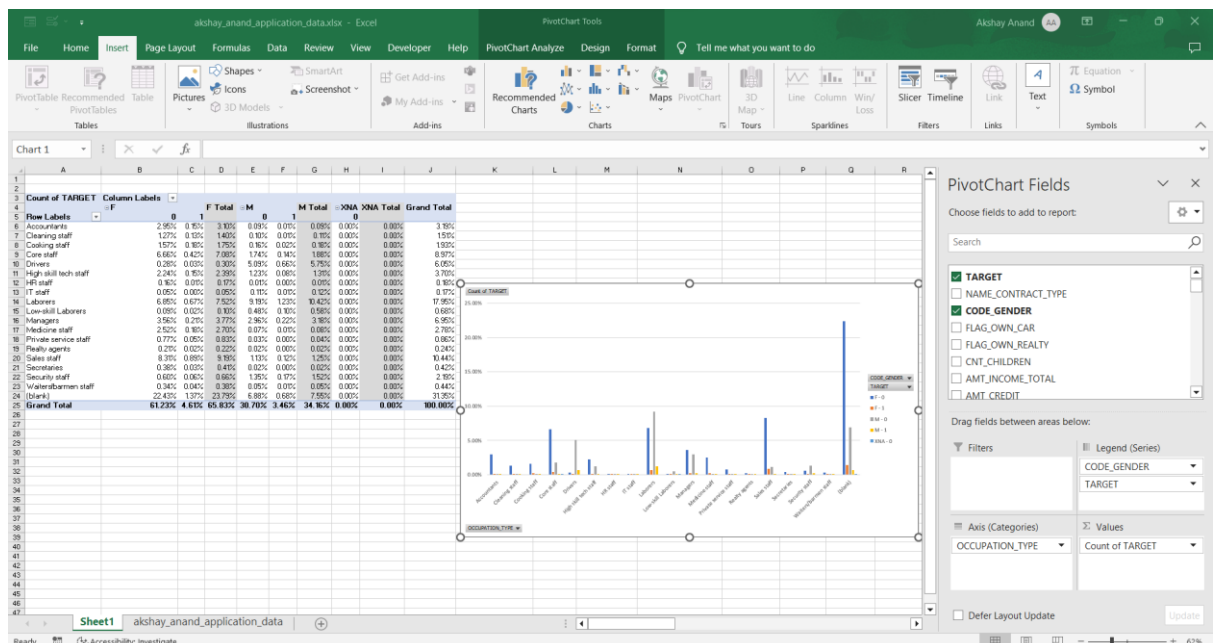FOR NAME_INCOME_TYPE

TARGET=0



TARGET=1

FROM THE ABOVE BAR PLOT WE CAN INFER THAT PEOPLE WITH WORKING TYPE AS WORKING HAVE THE HIGHEST COUNT FOR TARGET =0 AND 1 I.E. FOR CLIENTS WITH AND WITHOUT PAYMENT ISSUES.

SIMILARLY, WE CAN DO UNIVARIATE ANALYSIS FOR OTHER COLUMNS.

UNIVARIATE ANALYSIS IS THE ANALYSIS OF A SINGLE VARIABLE W.R. TO ANOTHER VARIABLE.

SINCE THERE ARE MORE VARIABLES BIVARIATE ANALYSIS IS DONE TO ANALYZE CONCURRENT RELATIONS BETWEEN 2 VARIABLES OR ATTRIBUTES.



ABOVE IS A BIVARIATE ANALYSIS OF OCCUPATION_TYPE VS GENDER W.R. TO TARGET VARIABLE.

- Find the top 10 **correlation** for the Client with payment difficulties and all other cases (Target variable). Note that you have to find the top correlation by
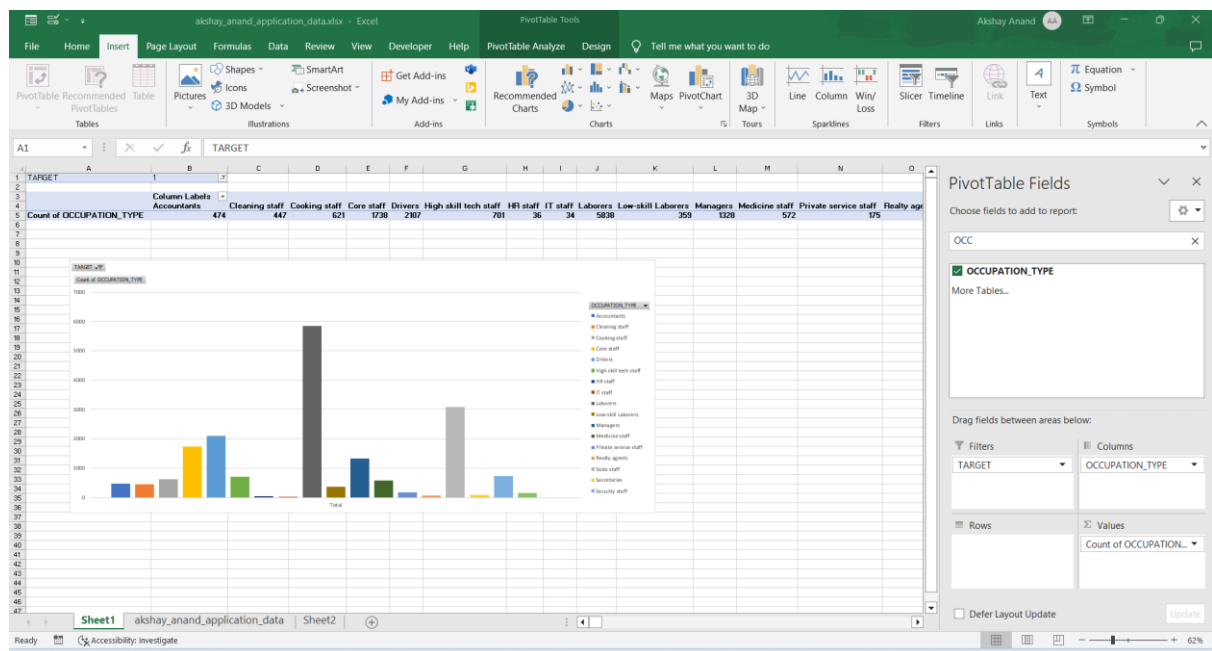
segmenting the data frame w.r.t to the target variable and then find the top correlation for each of the segmented data and find if any insight is there. Say, there are 5+1(target) variables in a dataset: Var1, Var2, Var3, Var4, Var5, Target. And if you have to find top 3 correlation, it can be: Var1 & Var2, Var2 & Var3, Var1 & Var3. Target variable will not feature in this correlation as it is a categorical variable and not a continuous variable which is increasing or decreasing.

FOR TARGET=1 WE FIND A CORRELATION USING THE CORREL FUNCTION BETWEEN AMT_CREDIT AND AMT_CREDIT_TOTAL AND FOUND A CORRELATION COEFFICIENT OF 1.

USED: =CORREL(Sheet1!I2:I307511,Sheet1!H2:H307511)

SIMILARLY WE FINAL CORRELATION COEFFICIENT 1 IN THE CASE OF COUNT_CHILDEREN VS AMT_TOTAL_INCOME. USED: =CORREL(Sheet1!G2:G307511,Sheet1!H2:H307511)

- **Include visualizations** and **summarize** the most important results in the presentation. You are free to choose the graphs which explain the numerical/categorical variables. Insights should explain why the variable is important for differentiating the clients with payment difficulties with all other cases.



LABOURERS' OCCUPATION TYPE HAS THE MOST PAYMENT ISSUES.

OTHER IMPORTANT RESULTS:

1)PEOPLE WITH ACADEMIC DEGREES DO FEWER DEFAULTS.

2)PEOPLE WITH LOW TOTAL INCOME ARE MORE LIKELY TO DEFAULT.

3) PEOPLE WITH HIGH CREDIT AMOUNTS ARE LESS LIKELY TO DEFAULT.

4)THE PROPORTION OF DEFAULTERS I.E. TARGET=1 IS 8% WHILE THAT OF NON-DEFAULTERS I.E. TARGET= 0 IS 92%.

5) PEOPLE BELONGING TO THE WORKING CLASS TEND TO PAY OFF THEIR LOANS ON TIME.

SIMILARLY, BY PERFORMING UNIVARIATE AND BIVARIATE ANALYSES AS DONE ABOVE, WE CAN GET MORE RESULTS.

RESULTS: THROUGH THIS EXPERIMENT, I UNDERSTOOD OUTLIER TREATMENT, UNIVARIATE ANALYSIS, PIVOT TABLES AND PLOTS. I ALSO USED AND UNDERSTOOD THE SIGNIFICANCE OF FILTERING IN EXCEL.