

# BANK LOAN ANALYSIS

Name : Hemant Murkute

Mail : [murkutehemant21@gmail.com](mailto:murkutehemant21@gmail.com)

**PROJECT DESCRIPTION** – This project is about a case study related to a Bank Loan. We have to carry out an EDA (Exploratory Data Analysis). Based on our analysis, we will get the solution for required questions.

**APPROACH** – I first analyzed the data. While analyzing, I found out that data had a lot of missing values. So my first task was to get the missing values by performing mean, median and mode functions as required. So, I began by cleaning the data and then finding the outliers so as to make the data standardized.

**TECH STACK USED** – MS Excel 2019

So let's begin with analysis.....

1. Present the overall approach of the analysis. Mention the problem statement and the analysis approach briefly

First we imported the data to excel.

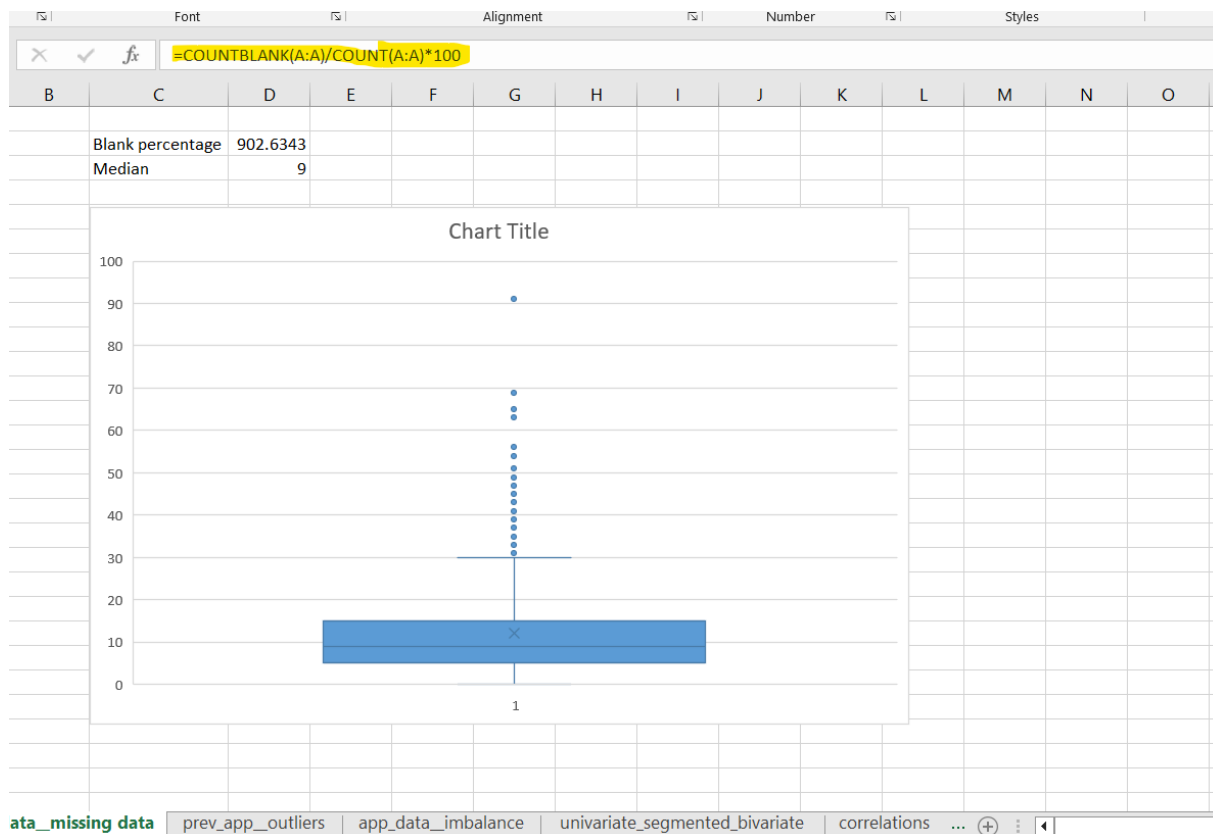
SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	NAME_TYPE_SUITE
100002	1	Cash loans	M	N	Y	0	202500	406597.5	24700.5	351000	Unaccompanied
100003	0	Cash loans	F	N	N	0	270000	1293502.5	35698.5	1129500	Family
100004	0	Revolving loans	M	Y	Y	0	67500	135000	6750	135000	Unaccompanied
100006	0	Cash loans	F	N	Y	0	135000	312682.5	29686.5	297000	Unaccompanied
100007	0	Cash loans	M	N	Y	0	121500	513000	21865.5	513000	Unaccompanied
100008	0	Cash loans	M	N	Y	0	99000	490495.5	27517.5	454500	Spouse, partner
100009	0	Cash loans	F	Y	Y	1	171000	1560726	41301	1395000	Unaccompanied
100010	0	Cash loans	M	Y	Y	0	360000	1530000	42075	1530000	Unaccompanied
100011	0	Cash loans	F	N	Y	0	112500	1019610	33826.5	913500	Children
100012	0	Revolving loans	M	N	Y	0	135000	405000	20250	405000	Unaccompanied
100014	0	Cash loans	F	N	Y	1	112500	652500	21177	652500	Unaccompanied
100015	0	Cash loans	F	N	Y	0	38419.155	148365	10678.5	135000	Children
100016	0	Cash loans	F	N	Y	0	67500	80865	5881.5	67500	Unaccompanied
100017	0	Cash loans	M	Y	N	1	225000	918468	28966.5	697500	Unaccompanied
100018	0	Cash loans	F	N	Y	0	189000	773680.5	32778	679500	Unaccompanied
100019	0	Cash loans	M	Y	Y	0	157500	299772	20160	247500	Family
100020	0	Cash loans	M	N	N	0	108000	509602.5	26149.5	387000	Unaccompanied
100021	0	Revolving loans	F	N	Y	1	81000	270000	13500	270000	Unaccompanied
100022	0	Revolving loans	F	N	Y	0	112500	157500	7875	157500	Other_A
100023	0	Cash loans	F	N	Y	1	90000	544491	17563.5	454500	Unaccompanied
100024	0	Revolving loans	M	Y	Y	0	135000	427500	21375	427500	Unaccompanied
100025	0	Cash loans	F	Y	Y	1	202500	1132573.5	37561.5	927000	Unaccompanied
100026	0	Cash loans	F	N	N	1	450000	497520	32521.5	450000	Unaccompanied
100027	0	Cash loans	F	N	Y	0	83250	239850	23850	225000	Unaccompanied
100029	0	Cash loans	M	Y	N	2	135000	247500	12703.5	247500	Unaccompanied
100030	0	Cash loans	F	N	Y	0	90000	225000	11074.5	225000	Unaccompanied
100031	1	Cash loans	F	N	Y	0	112500	979992	27076.5	702000	Unaccompanied
100032	0	Cash loans	M	N	Y	1	112500	327024	23827.5	270000	Family

Column1	Table	Row	Description	Special
1	application_data	SK_ID_CURR	ID of loan in our sample	
2	application_data	TARGET	Target variable (1 - client with payment difficulties: he/she had late payment more than X days on at	
3	application_data	NAME_CONTRACT_TYPE	Identification if loan is cash or revolving	
4	application_data	CODE_GENDER	Gender of the client	
5	application_data	FLAG_OWN_CAR	Flag if the client owns a car	
6	application_data	FLAG_OWN_REALTY	Flag if client owns a house or flat	
7	application_data	CNT_CHILDREN	Number of children the client has	
8	application_data	AMT_INCOME_TOTAL	Income of the client	
9	application_data	AMT_CREDIT	Credit amount of the loan	
10	application_data	AMT_ANNUITY	Loan annuity	
11	application_data	AMT_GOODS_PRICE	For consumer loans it is the price of the goods for which the loan is given	
12	application_data	NAME_TYPE_SUITE	Who was accompanying client when he was applying for the loan	
13	application_data	NAME_INCOME_TYPE	Clients income type (businessman, working, maternity leave,...)	
14	application_data	NAME_EDUCATION_TYPE	Level of highest education the client achieved	
15	application_data	NAME_FAMILY_STATUS	Family status of the client	
16	application_data	NAME_HOUSING_TYPE	What is the housing situation of the client (renting, living with parents, ...)	
17	application_data	REGION_POPULATION_RELATIVE	Normalized population of region where client lives (higher number means the client lives in more po normalized	
18	application_data	DAYS_BIRTH	Client's age in days at the time of application	time only relative to the application
19	application_data	DAYS_EMPLOYED	How many days before the application the person started current employment	time only relative to the application
20	application_data	DAYS_REGISTRATION	How many days before the application did client change his registration	time only relative to the application
21	application_data	DAYS_ID_PUBLISH	How many days before the application did client change the identity document with which he apply	time only relative to the application
22	application_data	OWN_CAR_AGE	Age of client's car	
23	application_data	FLAG_MOBILE	Did client provide mobile phone (1=YES, 0=NO)	
24	application_data	FLAG_EMP_PHONE	Did client provide work phone (1=YES, 0=NO)	
25	application_data	FLAG_WORK_PHONE	Did client provide home phone (1=YES, 0=NO)	
26	application_data	FLAG_CONT_MOBILE	Was mobile phone reachable (1=YES, 0=NO)	
27	application_data	FLAG_PHONE	Did client provide home phone (1=YES, 0=NO)	
28	application_data	FLAG_EMAIL	Did client provide email (1=YES, 0=NO)	



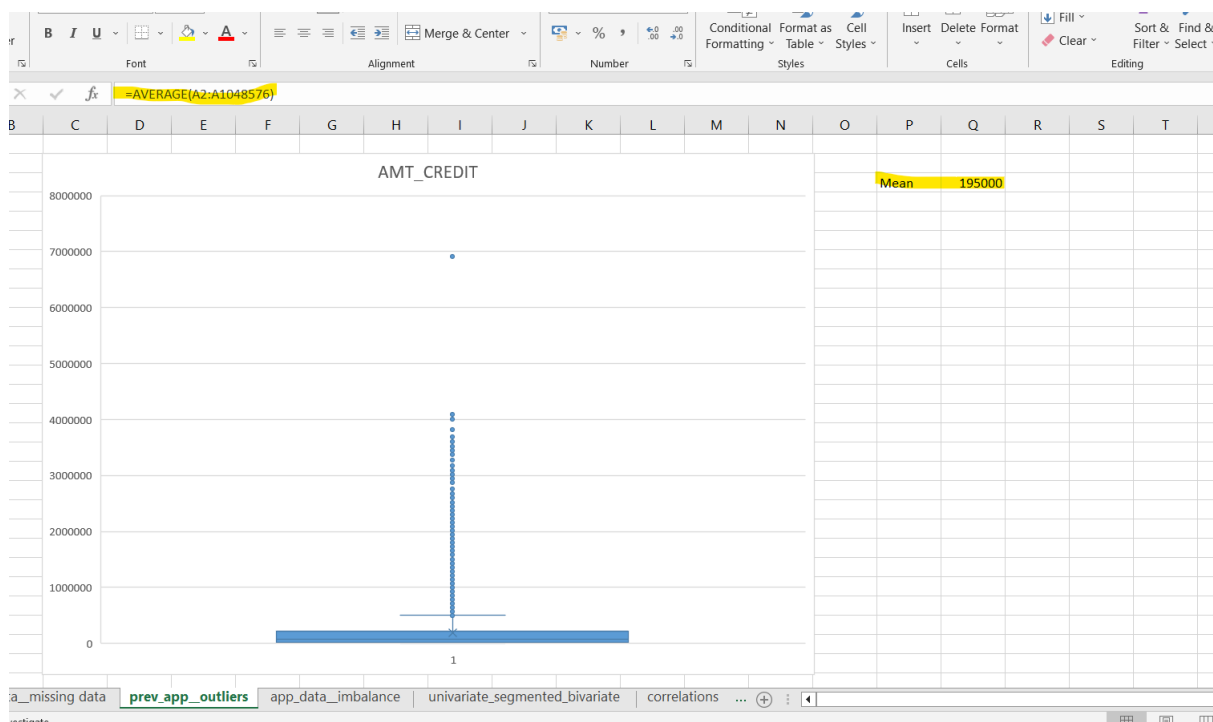
2. **Identify** the missing data and use appropriate method to deal with it. (Remove columns/or replace it with an appropriate value)

I found out the blank percentage and median of the column and filled the empty spaces there. (This is just for one table. Actual cleaning and filling of data is shown in excel file attached for other columns).

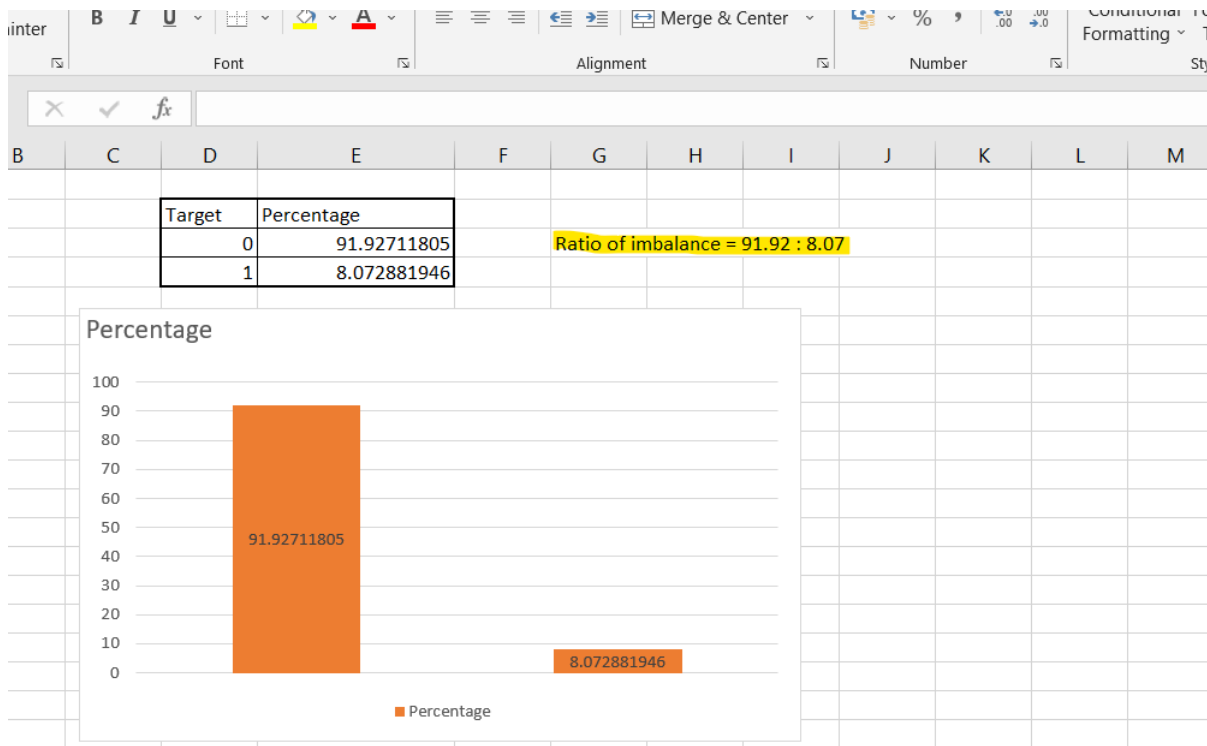


3. Identify if there are **outliers** in the dataset. Also, mention why do you think it is an outlier.

Foe Numerical columns, I found out the outliers and chose the value for the upper whisker as shown below. The credit amount value above 195000 is considered to be an upper whisker.



4. Identify if there is data imbalance in the data. Find the ratio of data imbalance. **The ratio of imbalance for Target Table came out to be 91.92:8.07.**

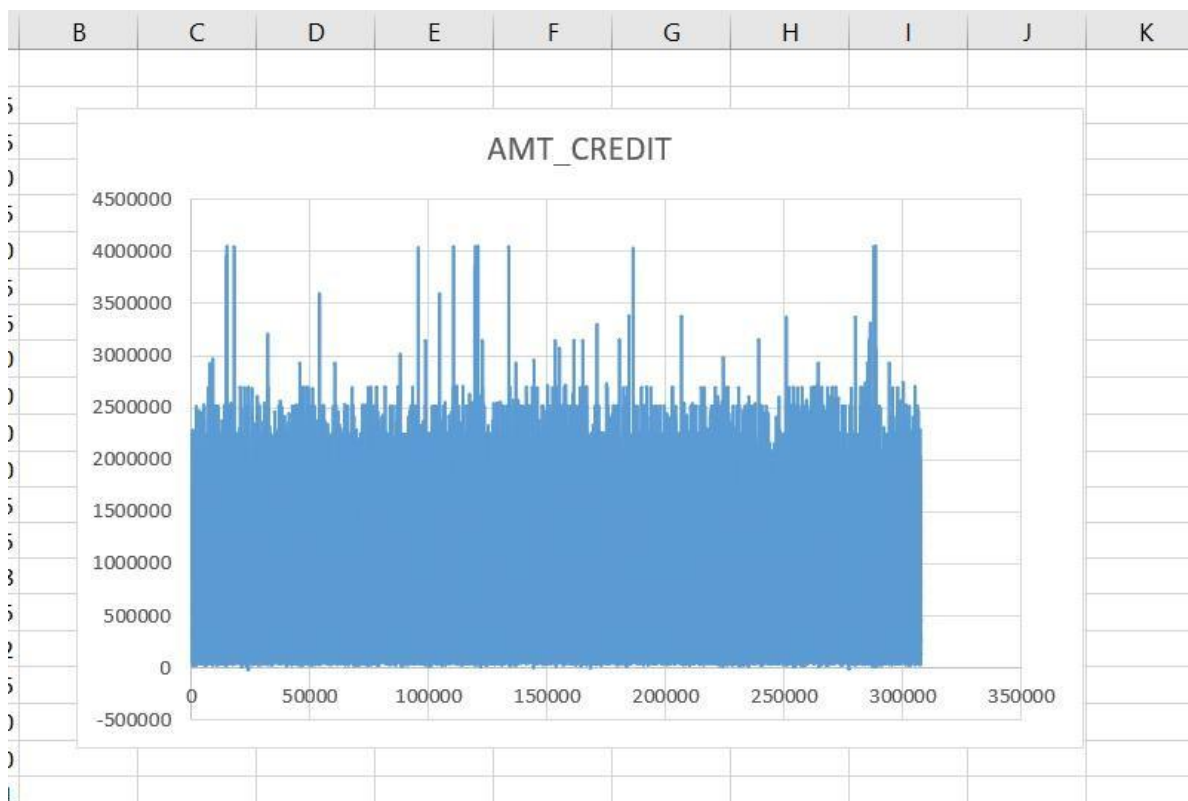


5. Explain the **results of univariate, segmented univariate, bivariate analysis, etc.** in business terms.

The results of **univariate, segmented univariate, bivariate analysis** are as follows –

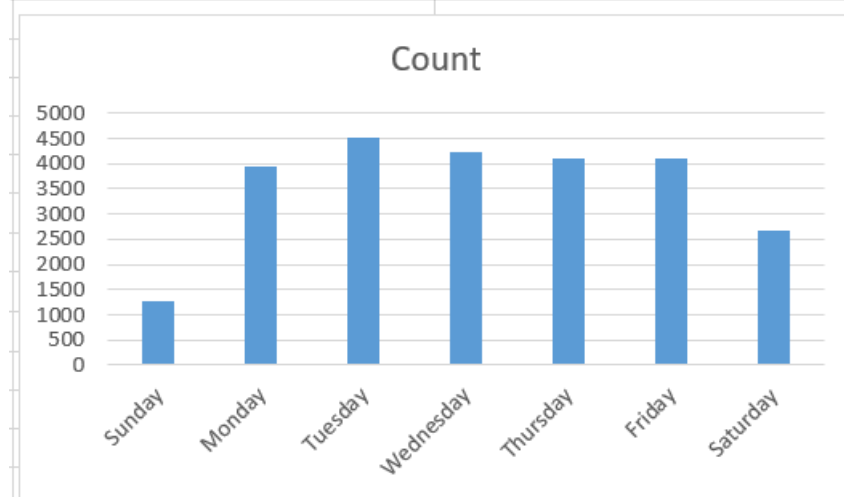
To perform the analysis, I first divided the data into two sets i.e. Target - 0 and Target – 1

## AMT\_CREDIT

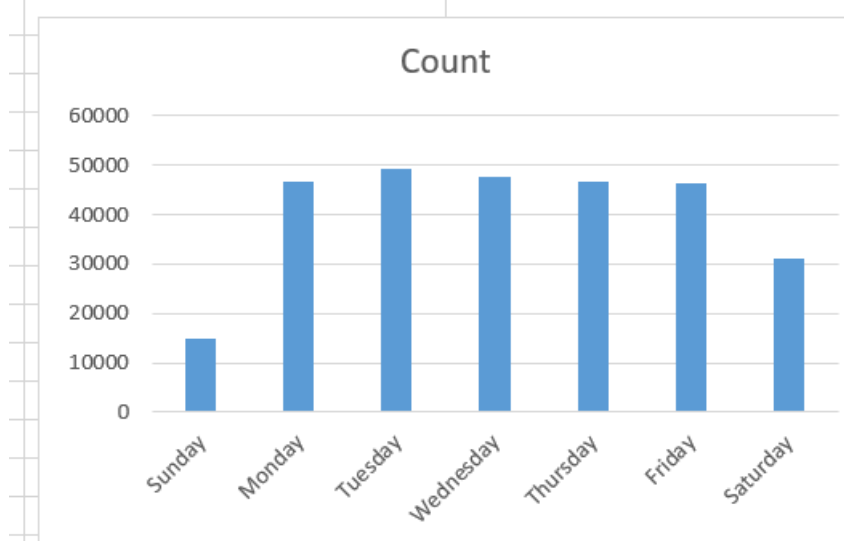


## WEEKDAY\_APPR\_PROCESS\_START

Target - 1	
WEEKDAY_APPR_PROCESS_START	Count
Sunday	1283
Monday	3934
Tuesday	4501
Wednesday	4238
Thursday	4098
Friday	4101
Saturday	2670



Target - 0	
WEEKDAY_APPR_PROCESS_START	Count
Sunday	14898
Monday	46780
Tuesday	49400
Wednesday	47696
Thursday	46493
Friday	46237
Saturday	31182

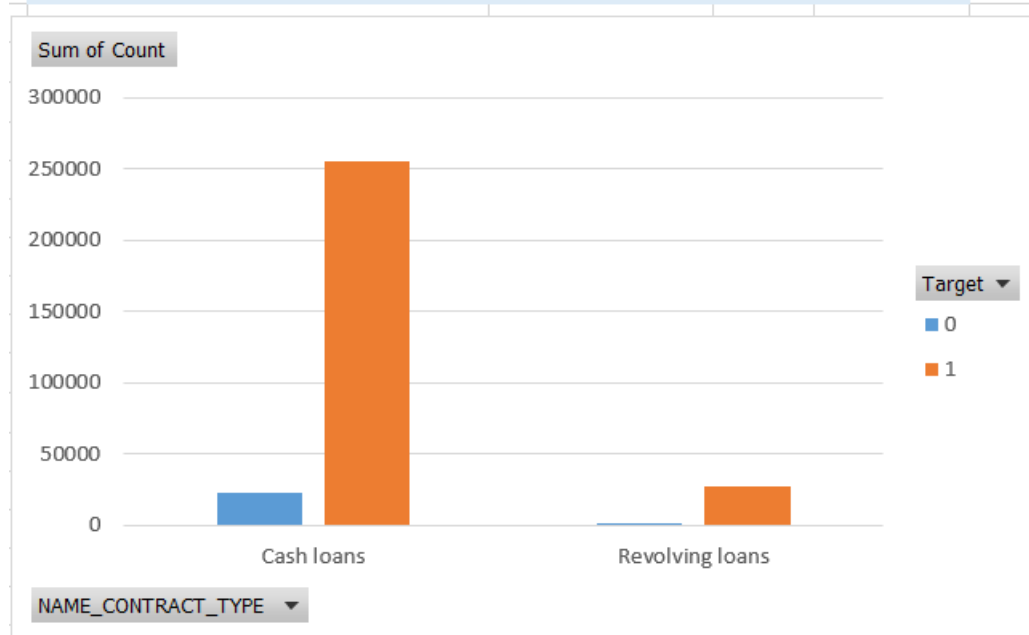


**INSIGHTS – We can conclude that application starting process is less on Saturday and Sunday.**

## NAME\_CONTRACT\_TYPE

Target - 1		
NAME_CONTRACT_TYPE	Count	
Cash loans	23221	
Revolving loans	1604	
Target - 0		
NAME_CONTRACT_TYPE	Count	
Cash loans	255011	
Revolving loans	27675	
NAME_CONTRACT_TYPE	Count	Target
Cash loans	255011	1
Revolving loans	27675	1
Cash loans	23221	0
Revolving loans	1604	0

Sum of Count	Column Labels		
Row Labels	0	1	Grand Total
Cash loans	23221	255011	278232
Revolving loans	1604	27675	29279
<b>Grand Total</b>	<b>24825</b>	<b>282686</b>	<b>307511</b>



**INSIGHTS – We can conclude that people prefer cash type loans more than other. People take more cash loans.**

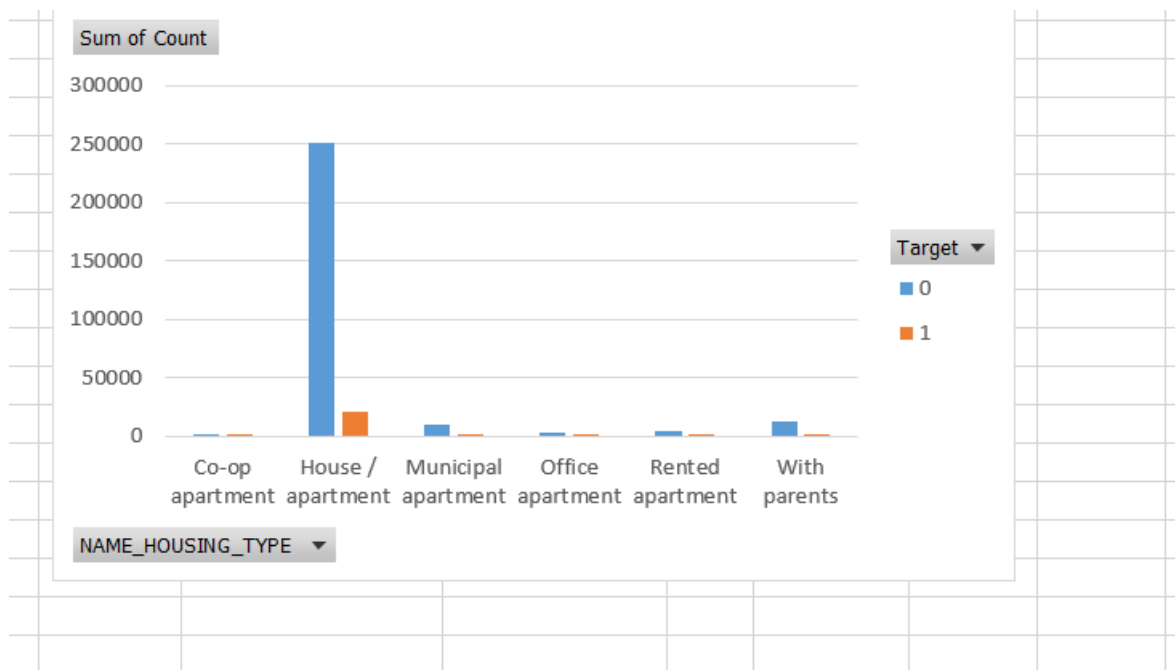


NAME\_HOUSING\_TYPE

Target - 1	
NAME_HOUSING_TYPE	Count
House / apartment	21272
Co-op apartment	89
Municipal apartment	955
Office apartment	172
Rented apartment	601
With parents	1736
Target - 0	
NAME_HOUSING_TYPE	Count
House / apartment	251596
Co-op apartment	1033
Municipal apartment	10228
Office apartment	2445
Rented apartment	4280
With parents	13104

NAME_HOUSING_TYPE	Count	Target
House / apartment	21272	1
Co-op apartment	89	1
Municipal apartment	955	1
Office apartment	172	1
Rented apartment	601	1
With parents	1736	1
House / apartment	251596	0
Co-op apartment	1033	0
Municipal apartment	10228	0
Office apartment	2445	0
Rented apartment	4280	0
With parents	13104	0

Sum of Count		Column Labels		
Row Labels		0	1	Grand Total
Co-op apartment	1033	89		1122
House / apartment	251596	21272		272868
Municipal apartment	10228	955		11183
Office apartment	2445	172		2617
Rented apartment	4280	601		4881
With parents	13104	1736		14840
Grand Total	282686	24825		307511



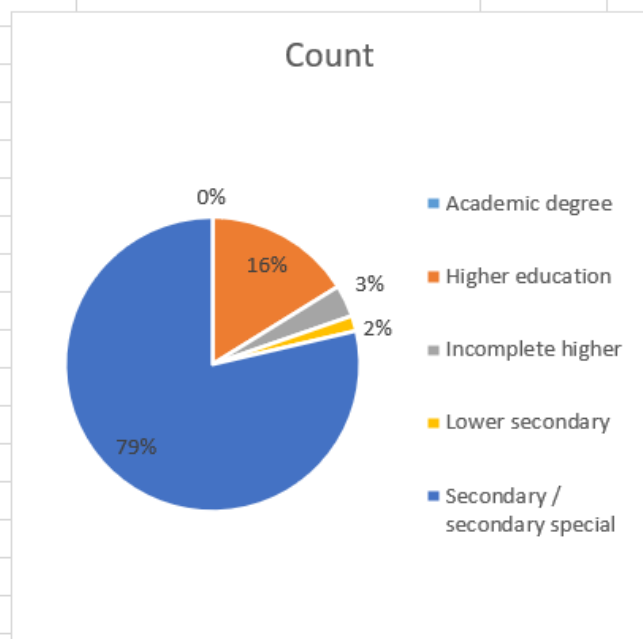
**INSIGHTS – We can conclude that people living in houses fall in both the category of default loans and non-default loans.**

- Find the top 10 **correlation** for the Client with payment difficulties and all other cases (Target variable).

To find the correlation, we again divide the data into two sets based on Targets and consider Target – 1 as defaulters.

### NAME\_EDUCATION\_TYPE

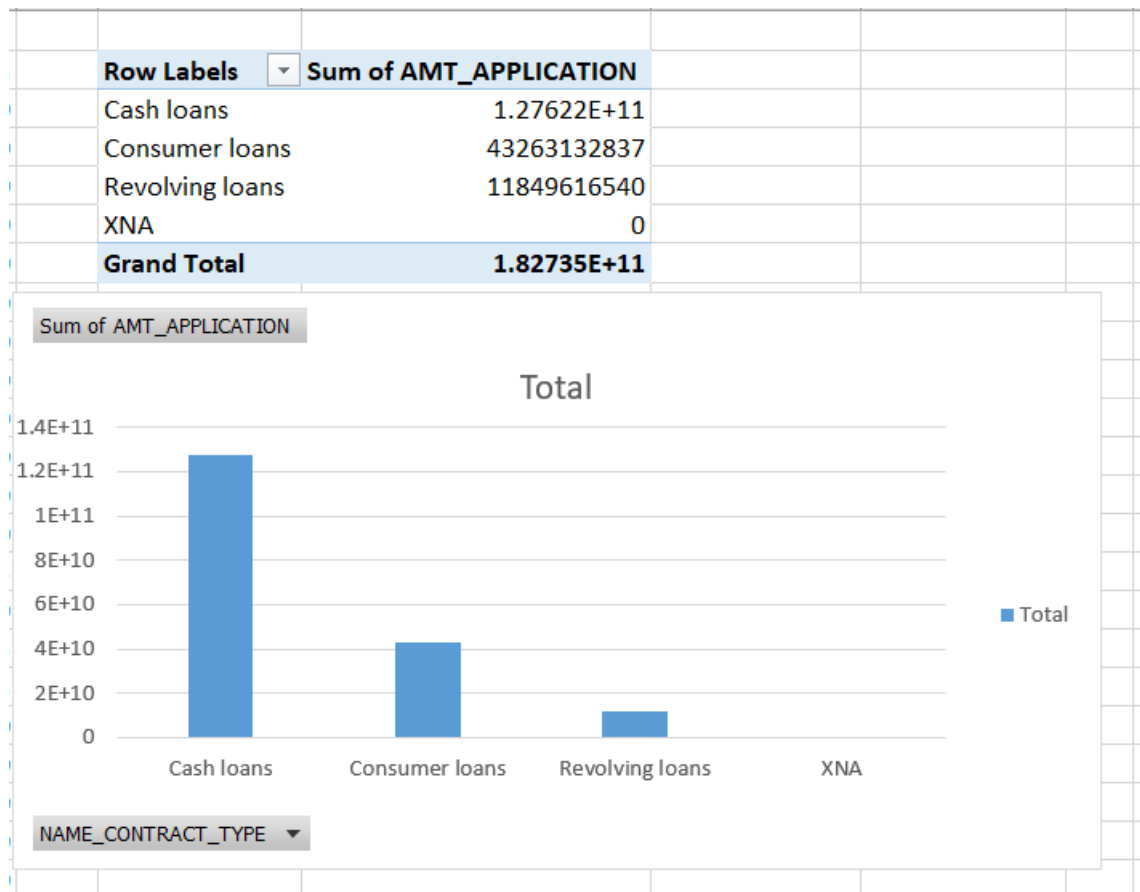
NAME_EDUCATION_TYPE	Count
Academic degree	3
Higher education	4009
Incomplete higher	872
Lower secondary	417
Secondary / secondary special	19524





**INSIGHTS – We can find that people with education type as Secondary/Secondary Special are more likely to default and people with education type Academic degree default the least.**

### NAME\_CONTRACT\_TYPE and AMT\_APPLICATION



**INSIGHTS – If we sum the total amount for loan in applications, we find that that people mostly take cash loans.**

### AMT\_APPLICATION and AMT\_CREDIT

AMT_APPLICATION	AMT_CREDIT		
17145	17145	Correlation Coefficient	
607500	679671	0.975777217	
112500	136444.5		
450000	470790		
337500	404055		
315000	340573.5		
0	0		
0	0		

**INSIGHTS – We find that the correlation coefficient is 0.9758 using excel formula =CORREAL.**

## AMT\_INCOME\_TOTAL and AMT\_ANNUITY

W	X	Y	Z	
AMT_INCOME_TOTAL	AMT_ANNUITY			
202500	24700.5		Correlation Coefficient	
270000	35698.5		0.191657428	
67500	6750			
135000	29686.5			
121500	21865.5			
99000	27517.5			
-----	-----			

**INSIGHTS – We find that the correlation coefficient is 0.19166 using excel formula =CORREAL.**

**CONCLUSION –** From the above analysis, we can find out what kind of people and can repay loan, what kinds of loan people prefer to take, people taking loans come from what background, what is their source of income, for what type of people, the loan applications are refused and based on which conditions.

## RESULTS: -

1. People with academic degree have less defaults.
2. People prefer cash loans more than any other type.
3. People with secondary/secondary special as education type have more chances of defaulting loans.
4. People who have less than 5 years of employment have high default rate.
5. Focused variable for application file – Target.
6. Focused variable for Previous application file – NAME\_CONTRACT\_STATUS.
7. Important fields to consider for loan repayment are –
8. NAME\_EDUCATION\_TYPE
9. AMT\_INCOME\_TOTAL
10. DAYS\_EMPLOYED
11. AMT\_CREDIT
12. People with lower total income are more likely to default.
13. People with high Credit amount are less likely to default

