

- First we loaded the application\_data.csv file and checked its structure.
- Then we found the percentage of missing values of all the columns and removed columns with higher percentage (more than 50%) of null values .

```
df = df.loc[:,round(100*(df.isnull().sum()/Len(df.index)),2)<50]
```

- Selected columns in variable which have low percentage of missing values, then checked its description using describe. After that imputed the null values in these columns using best metric (Mean/ Median).

```
Input_cols = ['AMT_GOODS_PRICE', 'EXT_SOURCE_2', 'EXT_SOURCE_3',  
              'OBS_30_CNT_SOCIAL_CIRCLE', 'DEF_60_CNT_SOCIAL_CIRCLE', 'CNT_FAM_MEMBERS']
```

AMT\_GOODS\_PRICE = will impute this columns null value with mean as it don't have any outliers.

EXT\_SOURCE\_2 = will impute this columns null value with mean as it don't have any outliers.

EXT\_SOURCE\_3 = will impute this columns null value with mean as it don't have any outliers.

OBS\_30\_CNT\_SOCIAL\_CIRCLE = will impute this columns null value with median as this column has outliers which affect the mean.

DEF\_60\_CNT\_SOCIAL\_CIRCLE = will impute this columns null value with median as this column has outliers which affect the mean.

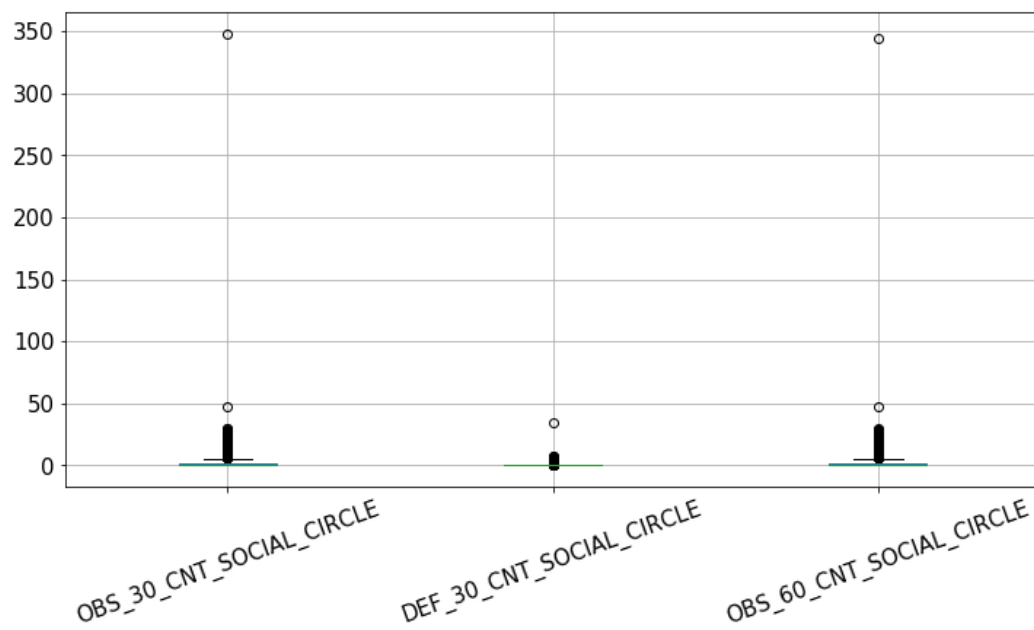
CNT\_FAM\_MEMBERS = will impute this columns null value with median as this column has outliers which affect the mean.

- Imbalance Percentage for Variable Target is below:

Target = 1 then = 8.07

Target = 0 then 92.93

Imbalance % = 8.07/92.93



As we can see in Boxplot of 3 variables that there are some values present in this which are out of the box so they are clearly outliers. So the Variables 'OBS\_30\_CNT\_SOCIAL\_CIRCLE', 'DEF\_30\_CNT\_SOCIAL\_CIRCLE', 'OBS\_60\_CNT\_SOCIAL\_CIRCLE' have outliers in their values.

So to handle Outliers, we can remove them from the datasets or in some cases we can replace them also. In this we have replaced them by median.

	REGION_POPULATION_RELATIVE	Poppulation_level	CNT_FAM_MEMBERS	Family Size
0	0.018801	Medium Poppulated	1.0	Small Family
1	0.003541	Less Poppulated	2.0	Small Family
2	0.010032	Medium Poppulated	1.0	Small Family
3	0.008019	Less Poppulated	2.0	Small Family
4	0.028663	Medium Poppulated	1.0	Small Family
...	...	...	...	...
307506	0.032561	Medium Poppulated	1.0	Small Family
307507	0.025164	Medium Poppulated	1.0	Small Family
307508	0.005002	Less Poppulated	1.0	Small Family
307509	0.005313	Less Poppulated	2.0	Small Family
307510	0.046220	Medium Poppulated	2.0	Small Family

307511 rows × 4 columns

For Binning, we have selected 2 variables (REGION\_POPULATION\_RELATIVE, CNT\_FAM\_MEMBERS ).

We have assigned the binning Values as Below:  
'REGION\_POPULATION\_RELATIVE:

REGION_POPULATION_RANGE	Bin
0-0.01	Less Poppulated
0.01-0.1	Medium Poppulated
0.1 - 1.0	Highly Poppulated

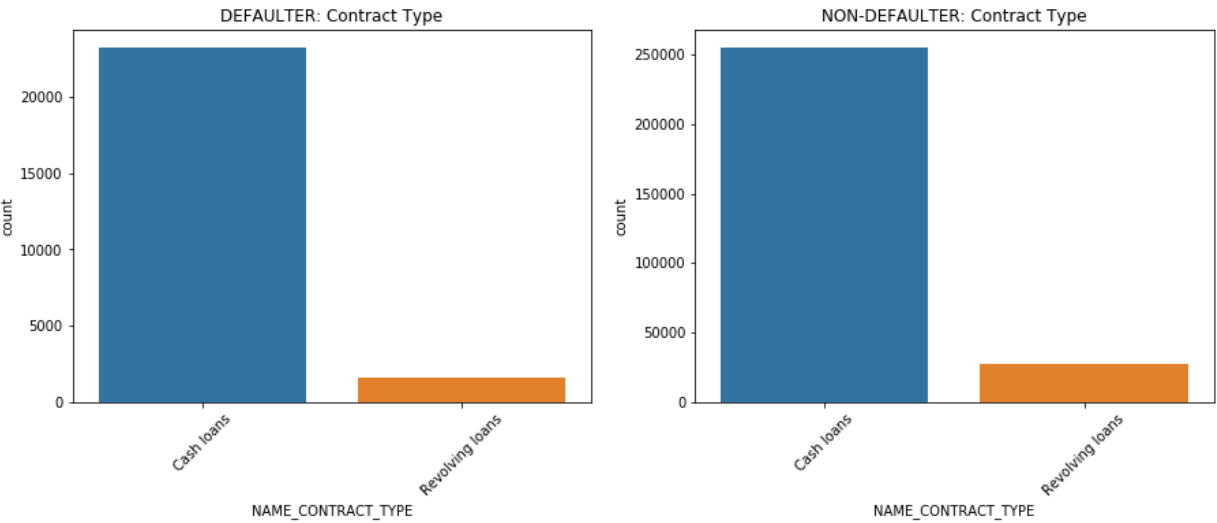
CNT\_FAM\_MEMBERS

CNT_FAM_MEMBERS_RANGE	Bin
0-4	Small Family
4-10	Medium Family
10-20	Large Family

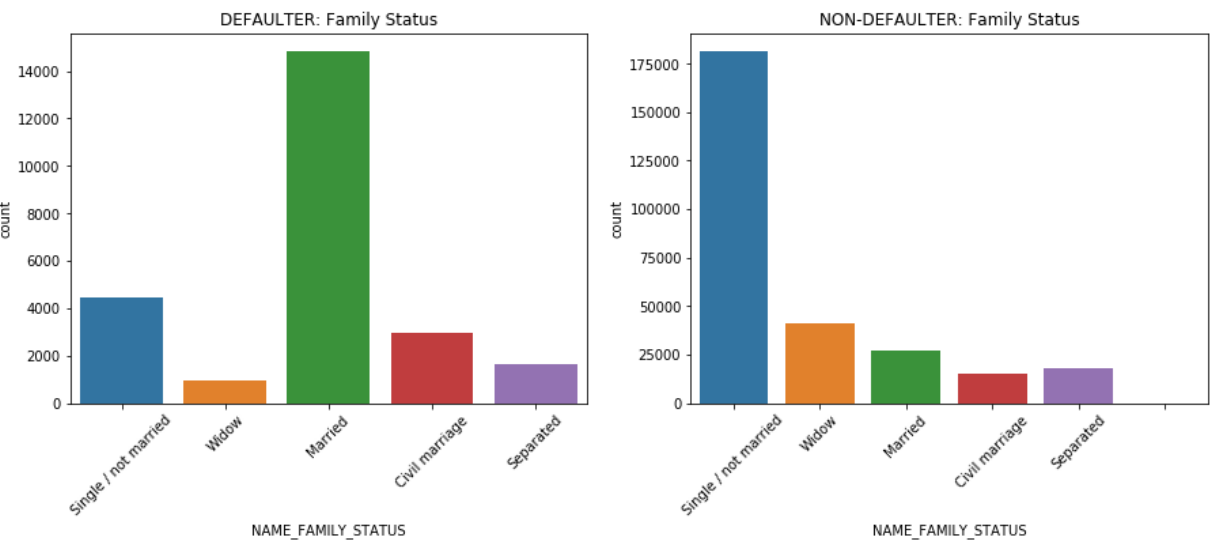
# Univariate analysis for categorical variables for both Target 0 and 1.

Below are the Comparison of the target variable across categories of categorical variables. Left plot is for defaulter in all plots.

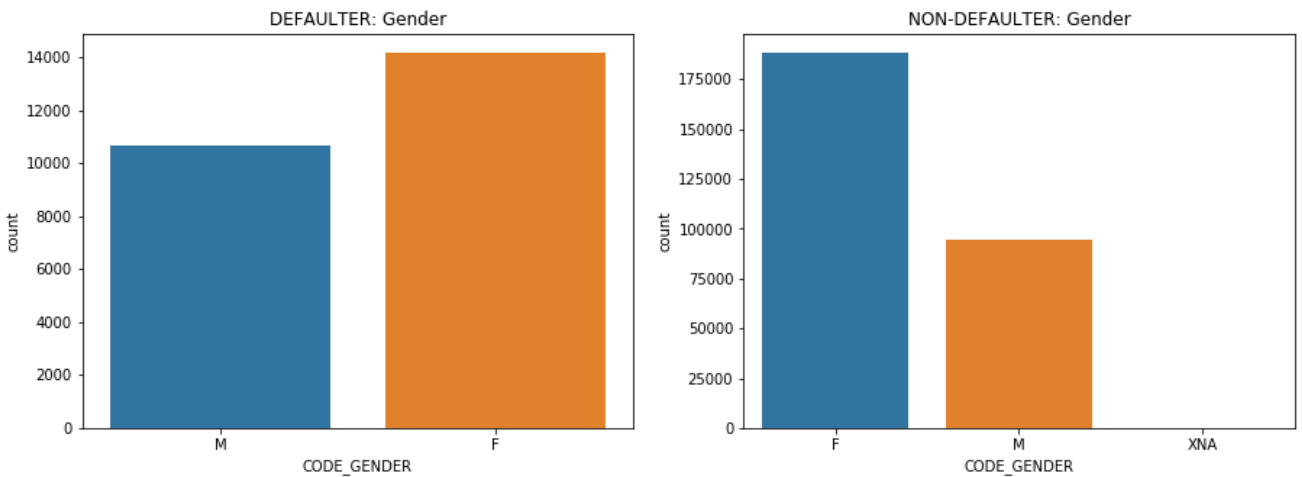
- NAME\_CONTRACT\_TYPE: Loan type for both Defaulter and Non Defaulter have as cash loans.



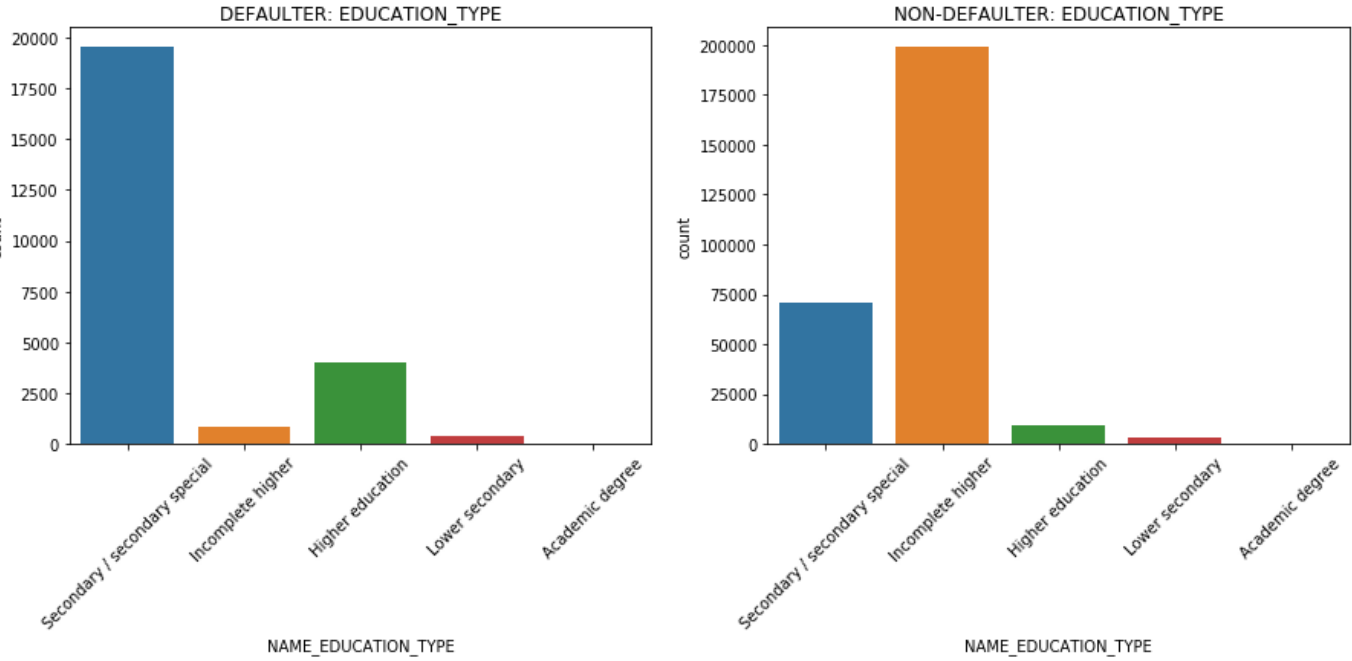
- NAME\_FAMILY\_STATUS: Married people are more likely to be Defaulter and single/not married are more likely to be Non Defaulter.



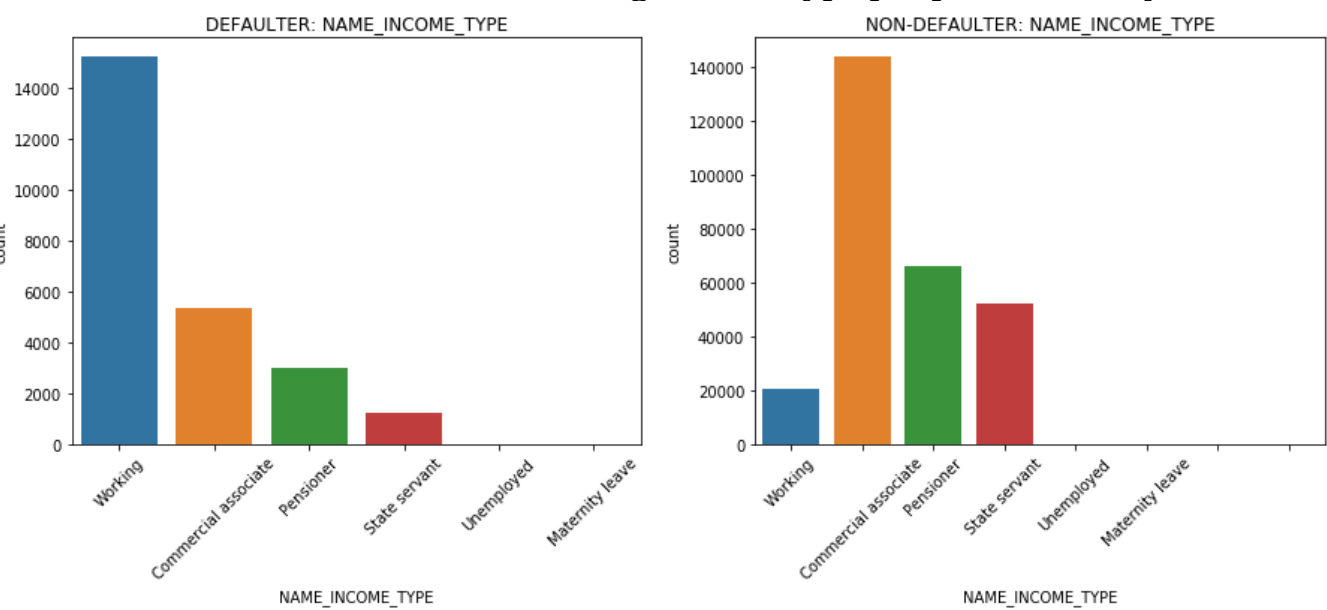
- CODE\_GENDER: Most Females are Defaulter and most males are Non Defaulter.



- NAME EDUCATION TYPE: Secondary specials educated are mostly Defaulter and Incomplete Higher Education are mostly Non Defaulter.



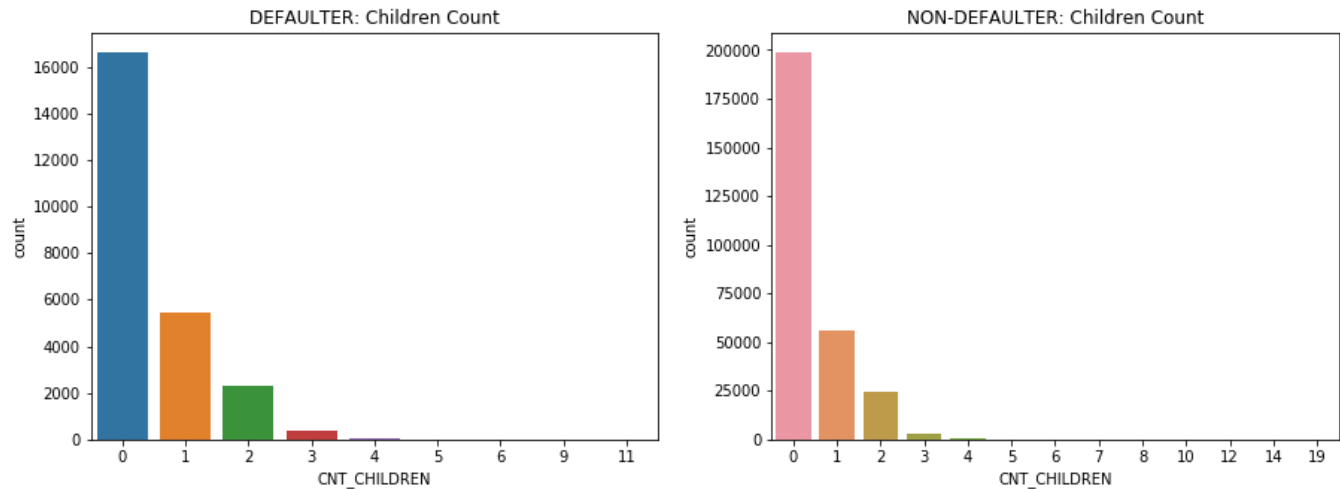
- NAME\_INCOME\_TYPE: Working Income type people are mostly defaulter and Commercial associate are mostly Non Defaulter.



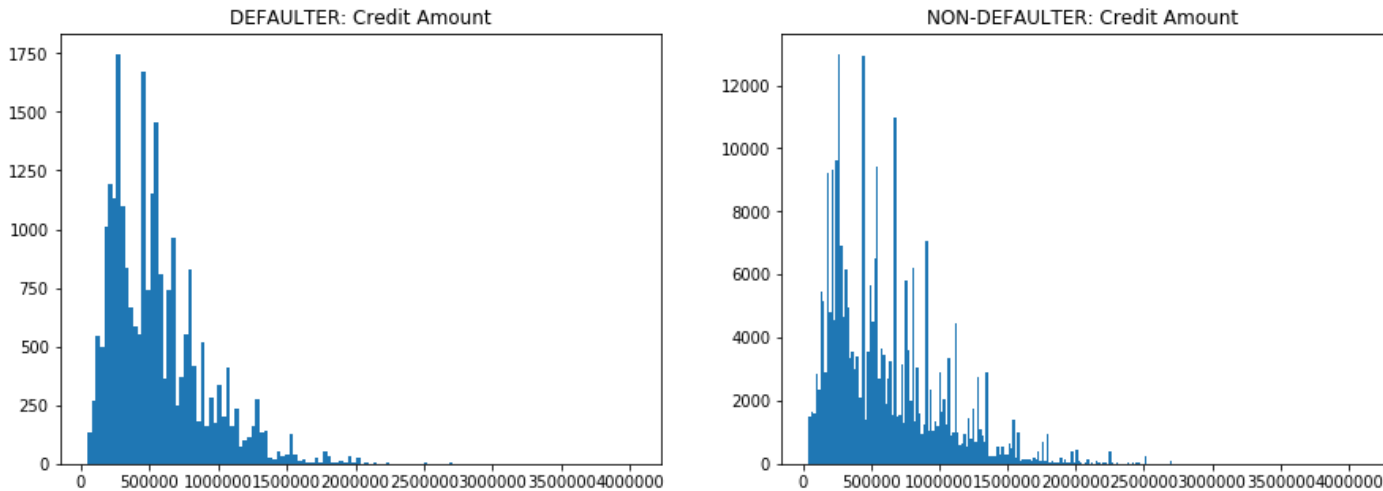
# Univariate analysis for numerical variables for both Target 0 and 1.

Below are the Comparison of the target variable across categories of Continues variables. Left plot is for defaulter in all plots.

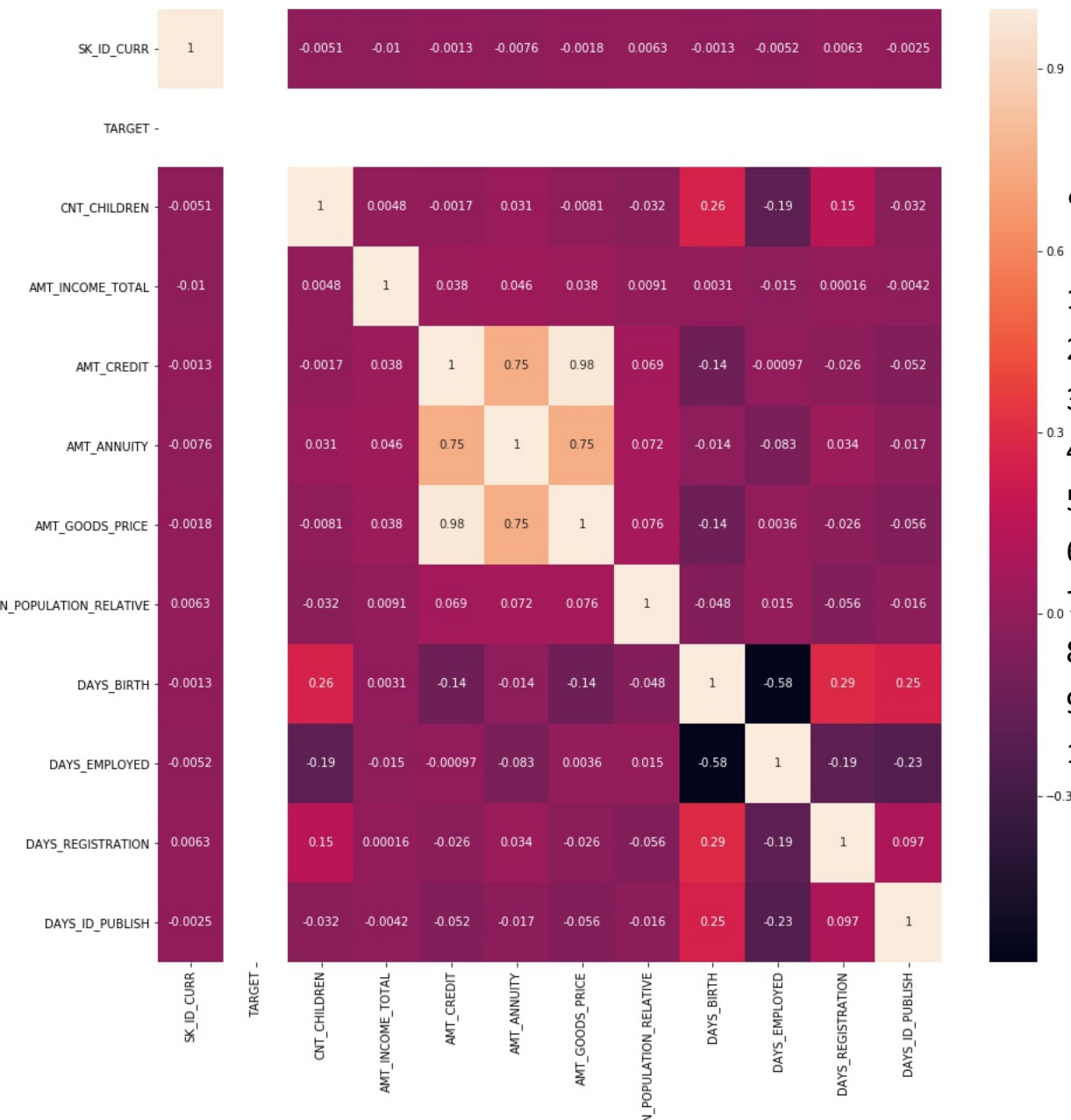
- CNT\_CHILDREN: Mostly Defaulter and Non-Defaulter have same children's count i.e. 0.



- AMT\_CREDIT: Credit amount is mostly distributed between 0-500000 for defaulter and for Non-Defaulter it is distributed equally with high spike between.



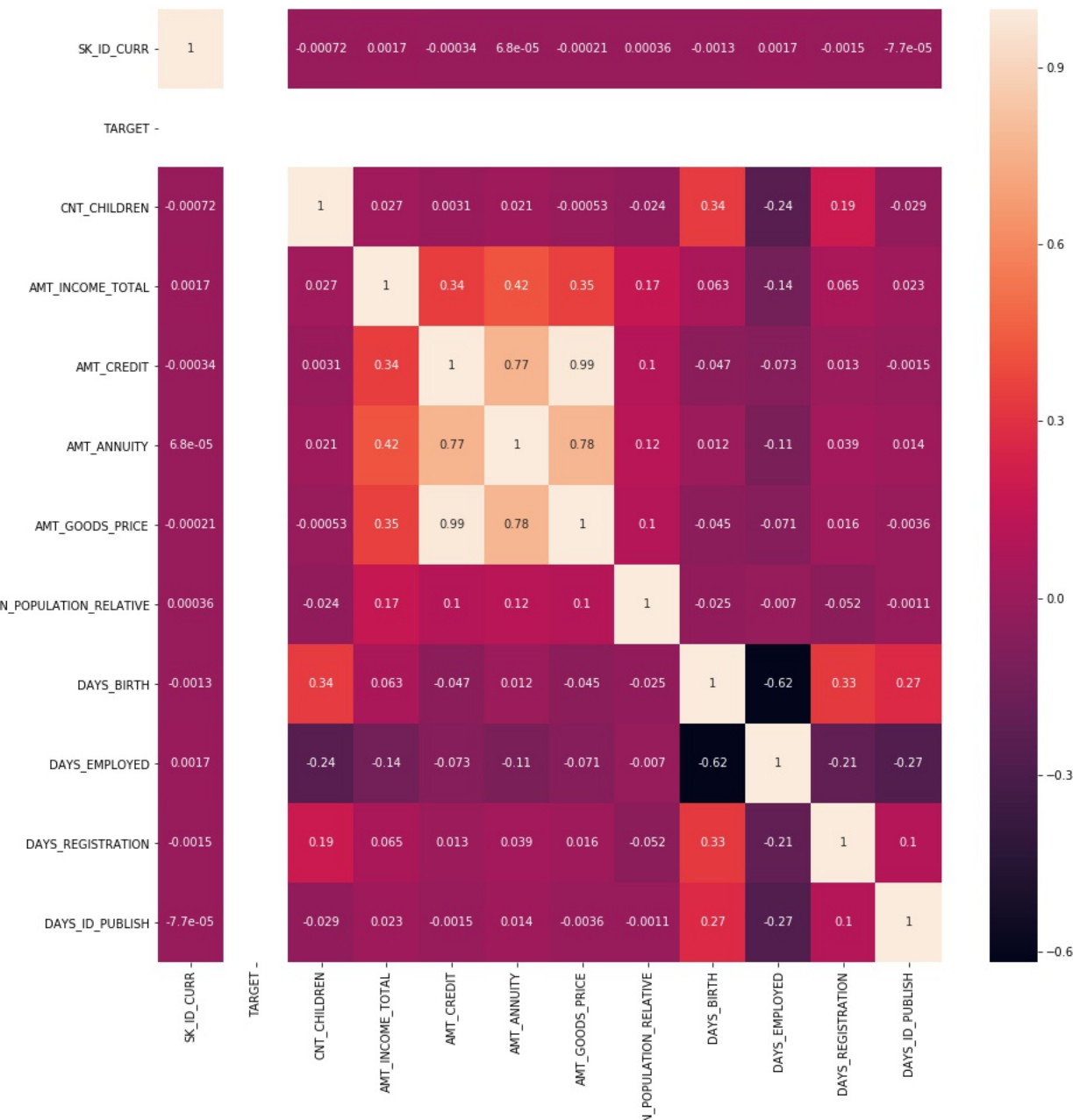
## ● HEATMAP FOR DEFAULTER CATEGORY:



● As per Defaulter Category Heatmap the top 10 correlation for the **Client with payment difficulties** are below:

1. AMT\_GOODS\_PRICE and AMT\_CREDIT
2. AMT\_GOODS\_PRICE and AMT\_AMMUNITY
3. AMT\_AMMUNITY and AMT\_CREDIT
4. DAYS\_BIRTH and DAYS\_REGISTRATION
5. CNT\_CHILDREN and DAYS\_BIRTH
6. DAYS\_BIRTH and DAYS\_ID\_PUBLISH
7. CNT\_CHILDREN and DAYS\_REGISTRATION
8. DAYS\_REGISTRATION and DAYS\_ID\_PUBLISH
9. REGION\_POPULATION\_RELATIVE and AMT\_GOODS\_PRICE
10. REGION\_POPULATION\_RELATIVE and AMT\_AMMUNITY

## ● HEATMAP FOR NON-DEFAULTER CATEGORY:



● As per Non-Defaulter Category Heatmap, the top 10 correlated variables for the **Client without any payment difficulties** are below:

1. AMT\_GOODS\_PRICE and AMT\_CREDIT
2. AMT\_GOODS\_PRICE and AMT\_AMMUNITY
3. AMT\_AMMUNITY and AMT\_CREDIT
4. AMT\_AMMUNITY and AMT\_INCOME\_TOTAL
5. AMT\_GOODS\_PRICE and AMT\_INCOME\_TOTAL
6. AMT\_INCOME\_TOTAL and AMT\_CREDIT
7. CNT\_CHILDREN and DAYS\_BIRTH
8. DAYS\_BIRTH and DAYS\_REGISTRATION
9. DAYS\_BIRTH and DAYS\_ID\_PUBLISH
10. CNT\_CHILDREN and DAYS\_REGISTRATION

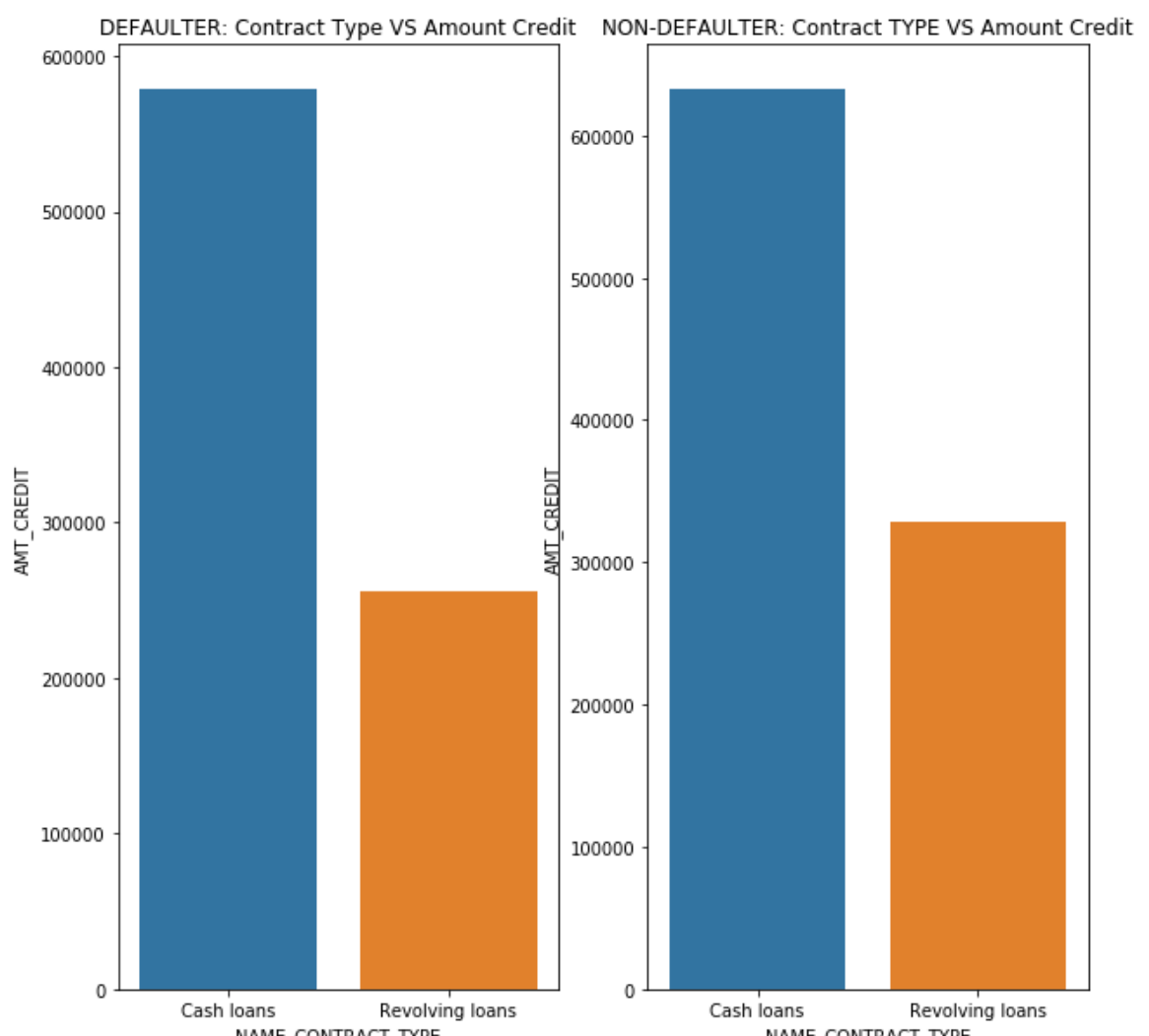
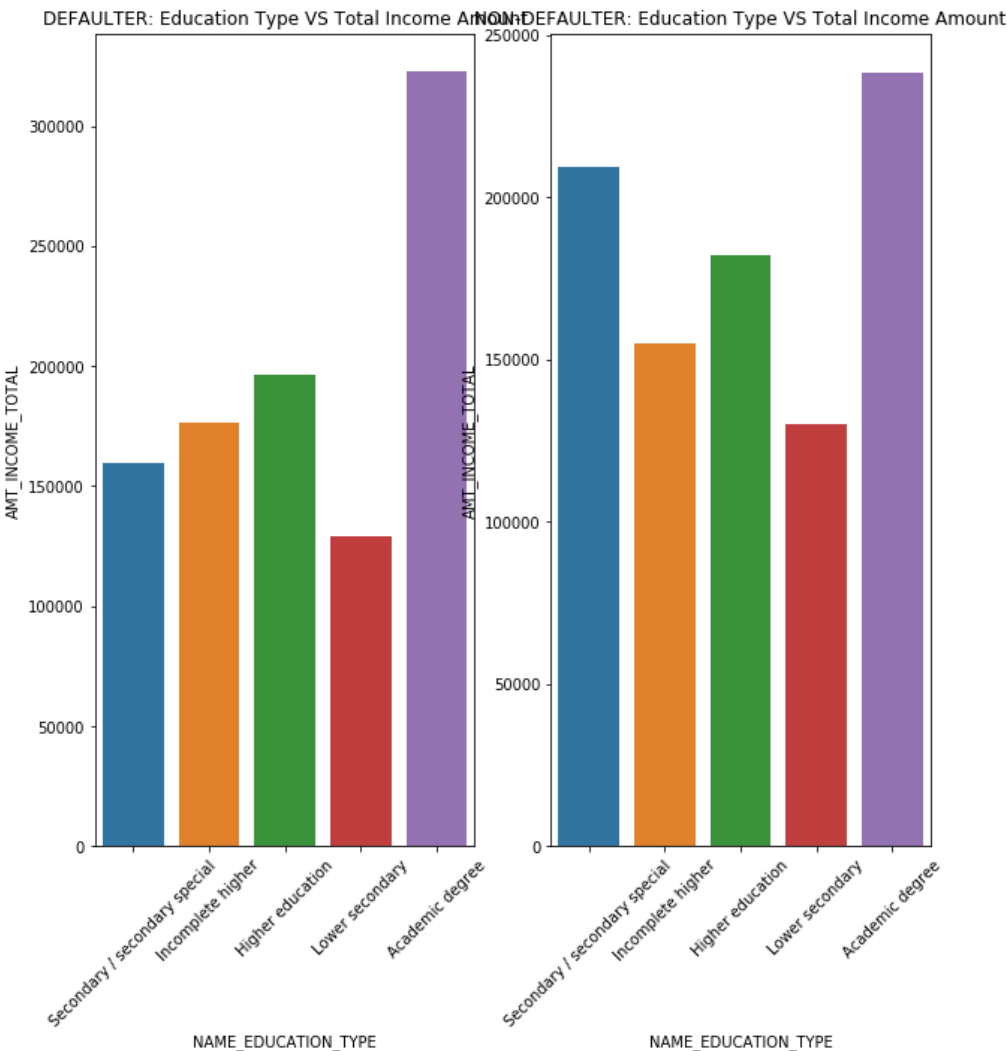
So As per heatmap analysis of both the Categories (Defaulter and Non-defaulter), it is confirmed that TOP-3 Correlated variables are same in both the Categories.

1. AMT\_GOODS\_PRICE and AMT\_CREDIT
2. AMT\_GOODS\_PRICE and AMT\_AMMUNITY
3. AMT\_AMMUNITY and AMT\_CREDIT

# Bivariate analysis for variables for both Target 0 and 1.

Below are the Comparison of the target variable across categories of variables. Left plot is for defaulter in all plots.

- Education Type VS Total Income: For Defaulter and Non Defaulter Academic Degree people are having more total Income compare to other educated person.
- Contract Type VS Amount Credit: For Cash loan types defaulter have Less Credited amount compare to Non-defaulter.

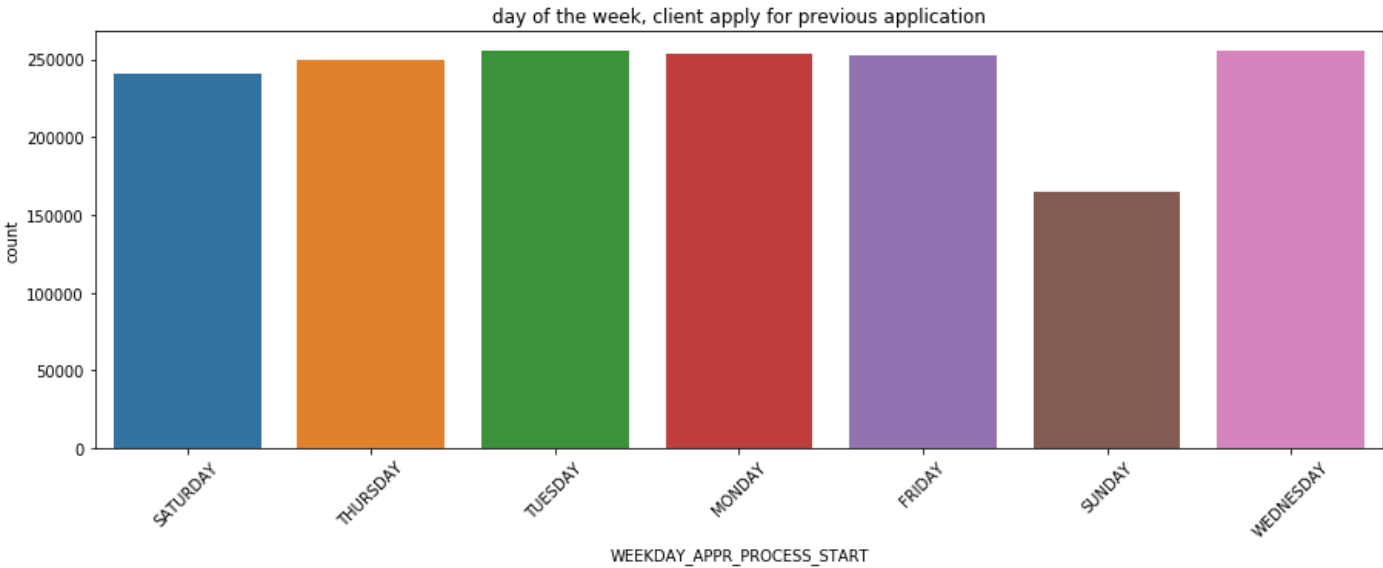




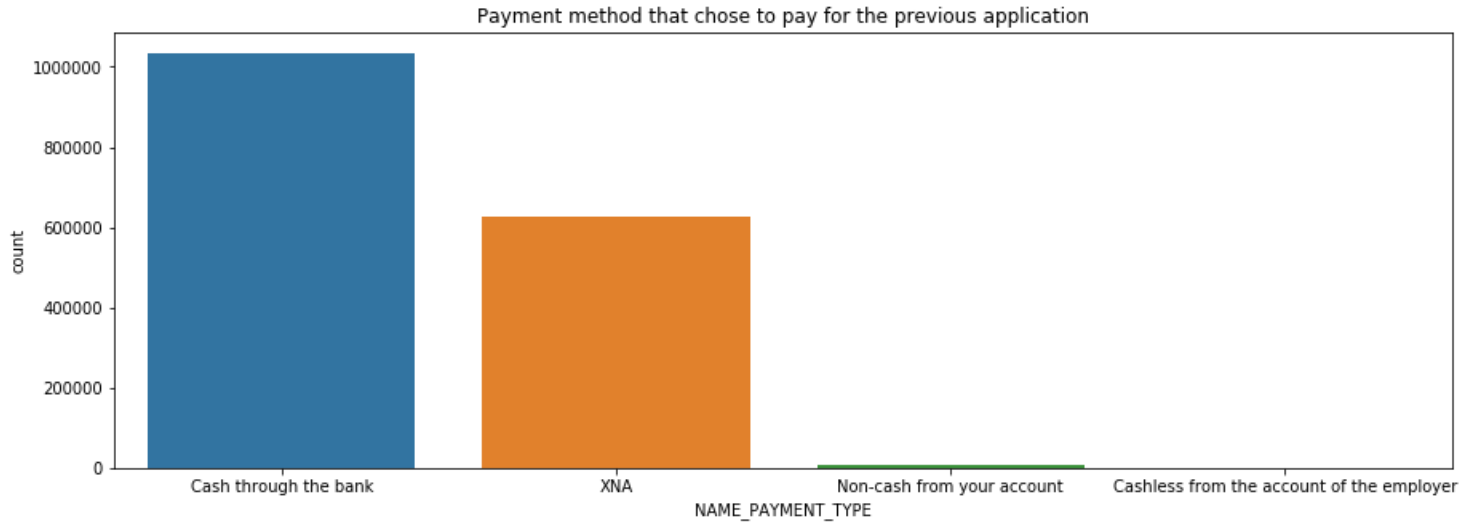
# Univariate analysis for variables of Previous Application Data

Below are the analysis of variables for Previous Application Data.

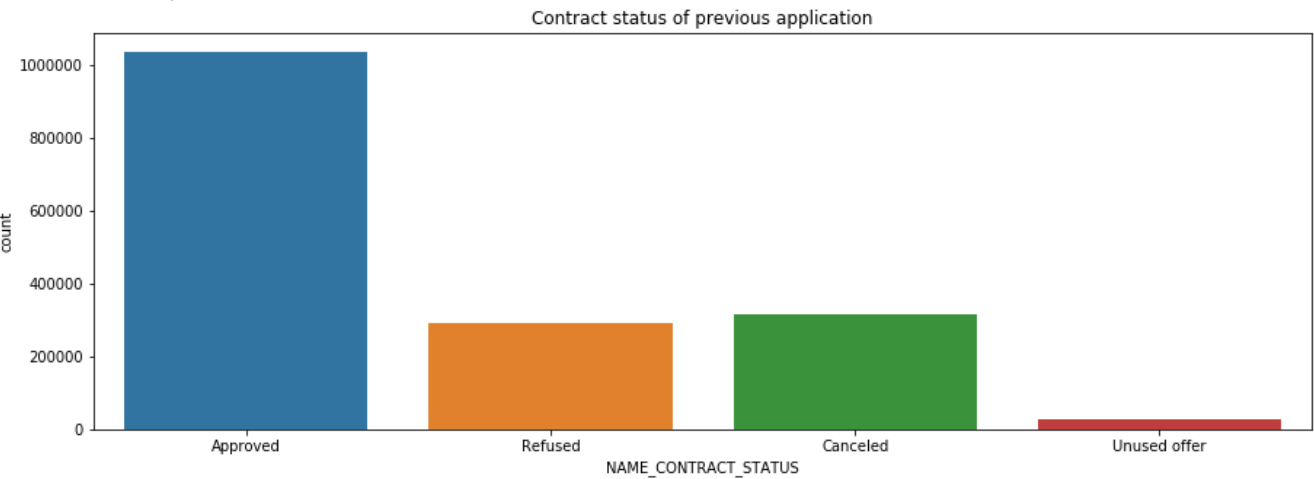
- WEEKDAY\_APPR\_PROCESS\_START: Sunday was the day when less application applied in previous applicaiton.



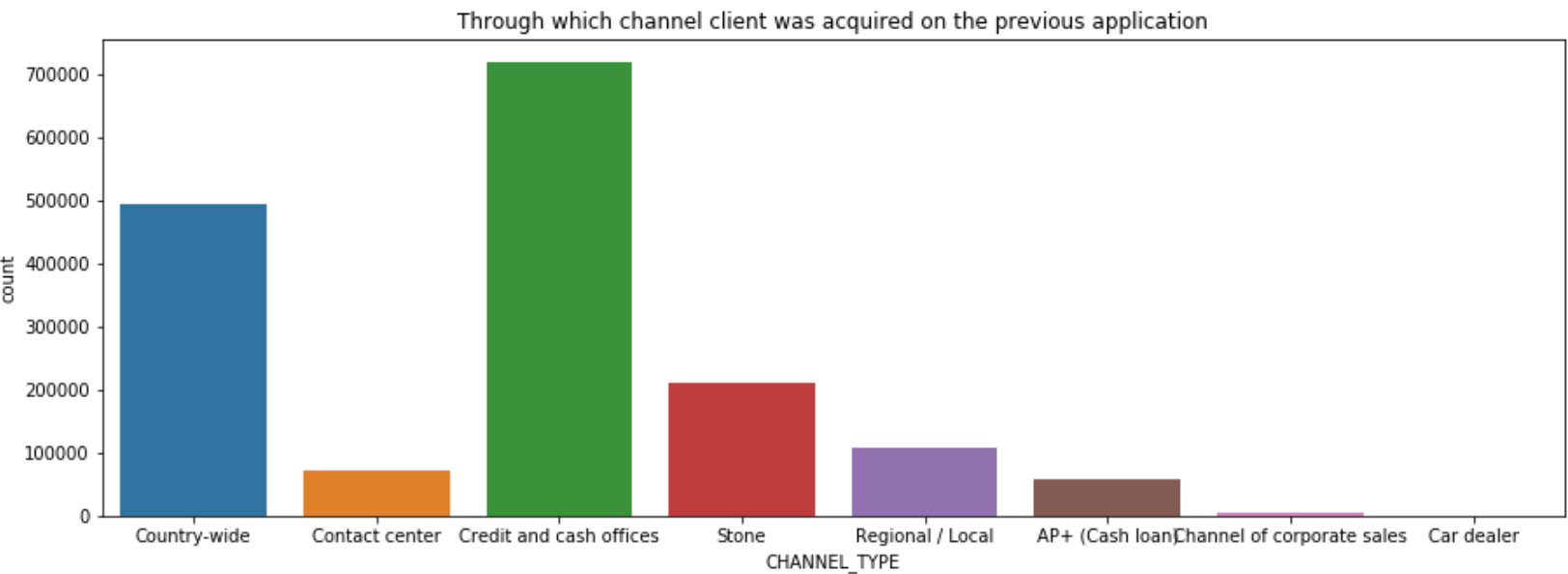
- NAME\_PAYMENT\_TYPE: As per shown in graph, mostly customers payed their previous application amount in Cash and mostly others didnt declare their payment method.



NAME\_CONTRACT\_STATUS: As per shown in graph, mostly application approved last time.  
But compare to refused, Cancelled were more.



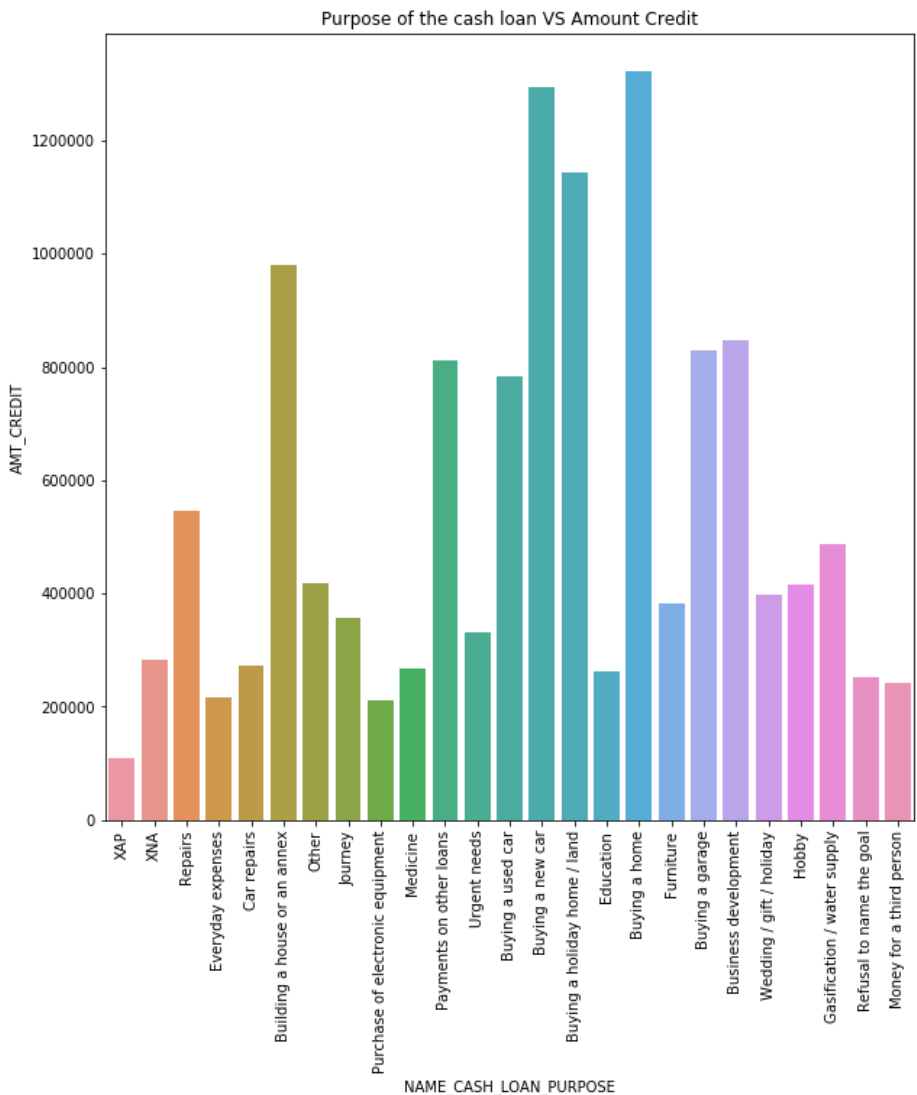
CHANNEL\_TYPE: As per shown in graph, mostly customers aquired on the previous application through Credit and Cash offices



# Bivariate analysis for variables of Previous Application Data

Below are the analysis of variables for Previous Application Data.

Purpose of the cash loan VS Amount Credit: Amount credited for Most of the Cash loans for the Buying a home.



Contract Type VS Application Amount: Application amount for refused applications are more high.

