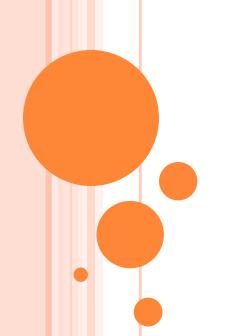
EDA CASE STUDY



SELECTING DATA

- As provided, we have worked on application_data.csv and previous_application.csv for our whole case study
- Firstly we worked on application_data.csv,
- Then previous_application.csv
- Finally we merged the two dataframes on their common column

		head()						
Out[3]:	SI	C_ID_CURR TA	ARGET	NAME_C	ONTRACT_TYPE	CODE_GENDER	FLAG_	_OWN_CA
	0	100002	1		Cash loans	M		
	1	100003	O		Cash loans	F		
	2	100004	O		Revolving loans	M		
	3	100006	0		Cash loans	F		
	4	100007	0		Cash loans	M		
	5 rows	s × 122 colum	ns					
	4							•
n [4]:	data.describe()							
Out[4]:								
		SK_ID_CURE	2	TARGET	CNT_CHILDREN	AMT_INCOME_T	OTAL	AMT_CRE
	count	307511.000000		TARGET 11.000000	CNT_CHILDREN 307511.000000	AMT_INCOME_T 3.075110		
	count		3075				0e+05	3.075110
		307511.000000	3075 ⁻	11.000000	307511.000000	3.075110	0e+05 9e+05	3.075110e 5.990260e
	mean	307511.000000 278180.51857	3075 ⁻ 7 3	11.000000 0.080729	307511.000000 0.417052	3.075110 1.687979	0e+05 9e+05 1e+05	3.075110e 5.990260e 4.024908e
	mean std	307511.000000 278180.51857 102790.175348	3075 ² 7 3	11.000000 0.080729 0.272419	307511.000000 0.417052 0.722121	3.075110 1.687979 2.37123	0e+05 9e+05 1e+05 0e+04	3.0751106 5.9902606 4.0249086 4.5000006
	mean std min	307511.000000 278180.518573 102790.175348 100002.000000	3075 ² 7 3 0	11.000000 0.080729 0.272419 0.000000	307511.000000 0.417052 0.722121 0.000000	3.075110 1.687979 2.37123 2.565000	0e+05 9e+05 1e+05 0e+04 0e+05	3.0751106 5.9902606 4.0249086 4.5000006
	mean std min 25%	307511.000000 278180.51857 102790.175348 100002.000000 189145.500000	3075	11.000000 0.080729 0.272419 0.000000 0.000000	307511.000000 0.417052 0.722121 0.000000 0.000000	3.075110 1.687979 2.37123 2.565000 1.125000	0e+05 9e+05 1e+05 0e+04 0e+05	3.0751106 5.9902606 4.0249086 4.5000006 2.7000006

- Initially we checked what file have and description of file.
- Shape of application_data is (307511,122)

 Now we have checked for null values in dataframe

• We dropped columns which have 50% and more null values

```
In [12]: data = data[data.columns[data.isnull().sum() < mid]]
    data.shape

Out[12]: (307511, 81)</pre>
```

• After dropping number of column reduced to 81 from 122

CATEGORICAL COLUMN VS NUMERIC COLUMN

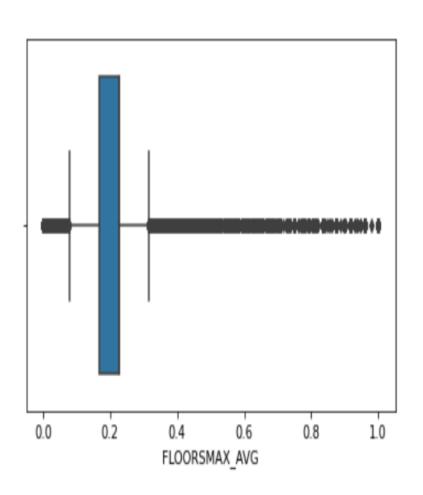
• This will list all the categorical column in dataframe

```
In [15]: # Categorical columns
list(set(data.columns) - set(data.describe().columns))
```

• This will list all the numeric column in dataframe

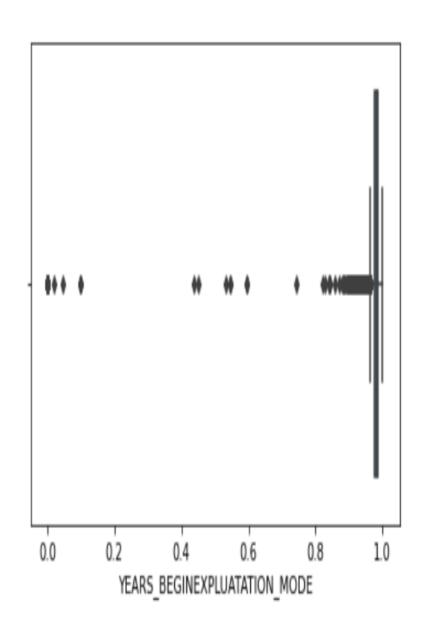
```
In [16]: #Numeric columns
data.describe().columns
```

FILL THE MISSING VALUES



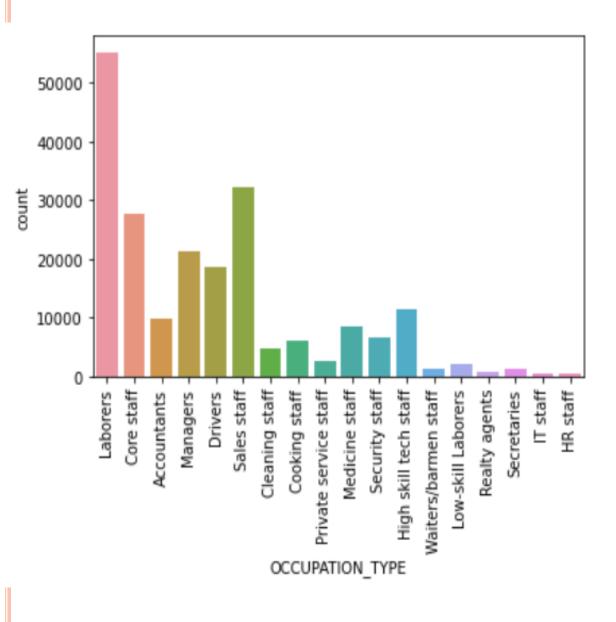
For FLOORMAX_AVG we will fillna with mean of FLOORMAX_AVG

Similarly for FLOORMAX_MEDI and FLOORMAX_MODE we filled missing with mean of respective column



For
YEARS_BEGINEXPLUATA
TION_MODE we will fillna
with median of
FLOORMAX_AVG
Because maximum value
belongs to their

Similarly for YEARS_BEGINEXPLUATA TION_AVG YEARS_BEGINEXPLUATA TION_MEDI with median of their respective column



Customers
having
Occupation
Type as
'Labourers' are
the highest
sector ging for
loans

HANDELLING COLUMN WITH UNIQUE VALUES HAVING NULL VALUES

• In this scenerio we will fillna() with mode of respective column and iloc[0]

• We have used this in:

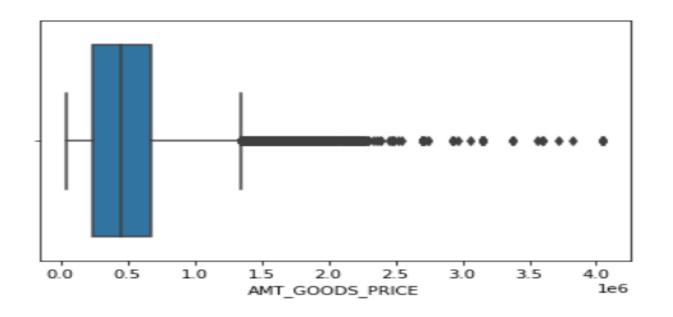
AMT_REQ_CREDIT_BUREAU_YEAR
AMT_REQ_CREDIT_BUREAU_QRT
AMT_REQ_CREDIT_BUREAU_MON
AMT_REQ_CREDIT_BUREAU_WEEK
AMT_REQ_CREDIT_BUREAU_DAY
AMT_REQ_CREDIT_BUREAU_HOUR

HANDELLING COLUMN WITH OBJECT TYPE VALUES HAVING NULL VALUES

- For such value we replaced all the missing value with most frequent value
- We done this in columns like:
 - NAME_TYPE_SUITE
 - EMERGENCYSTATE_MODE
 - OCCUPATION_TYPE

Using fillna(, inplace=True)

IDENTIFY OUTLIER USING PLOTS



- In this value after 400000 is clearly an oulier
- This is why we used mode to fill null values

• This is how we handeled all missing values using different methods to get clean and good dataframe

So, that we can carry our analysis in smooth way

• In this way we will have more accurate observations

OBSERVATION ON COLUMN STARTING FROM DATE

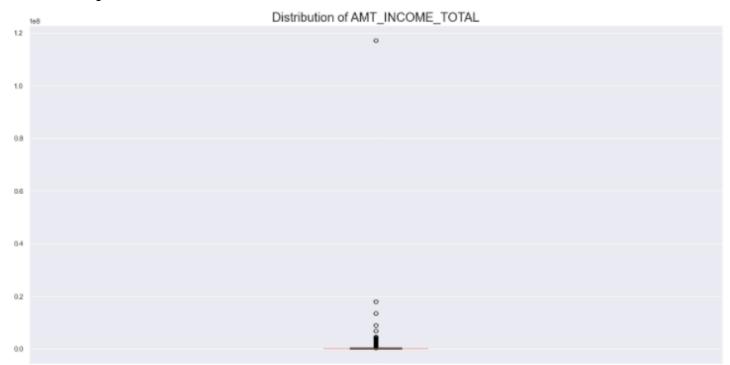
• With help of unique we know that there are some specific value for columns

```
In [174]: # Checking the values present in columns starting with 'DAYS'
    print(data['DAYS_BIRTH'].unique())
    print(data['DAYS_EMPLOYED'].unique())
    print(data['DAYS_REGISTRATION'].unique())
    print(data['DAYS_ID_PUBLISH'].unique())
    print(data['DAYS_LAST_PHONE_CHANGE'].unique())

[ -9461 -16765 -19046 ... -7951 -7857 -25061]
    [ -637 -1188 -225 ... -12971 -11084 -8694]
    [ -3648. -1186. -4260. ... -16396. -14558. -14798.]
    [ -2120 -291 -2531 ... -6194 -5854 -6211]
    [ -1134. -828. -815. ... -3988. -3899. -3538.]
```

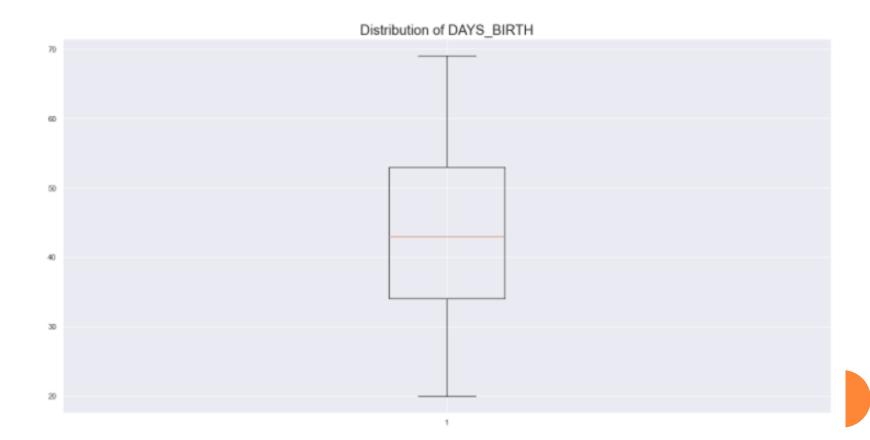
ANALYSING AMT_INCOME_TOTAL

• From graph we can surely say that 117 million is surely an outlier



ANALYSIS W.R.T TO CUSTOMER AGE

• No outlier and median in range 40-45

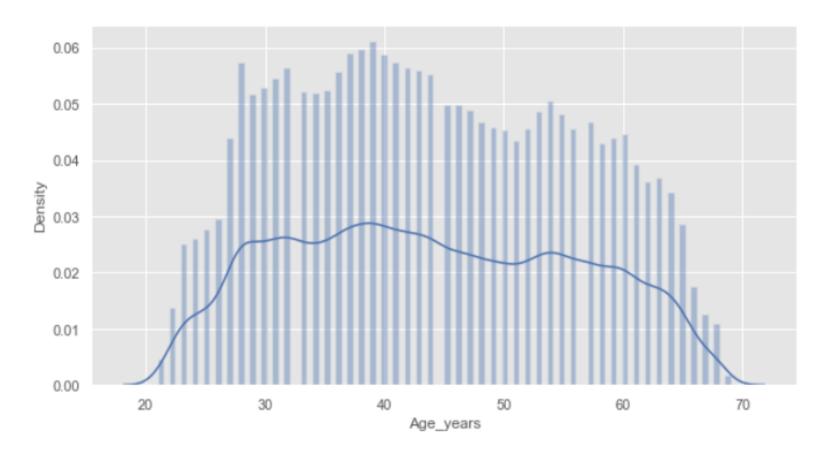


AGE GROUP ON TOP WHICH APPLY FOR LOAN

- For this we added 3 more column to dataset
 - TodayDate:
 - Date of current day
 - DateOfBirth:
 - Date of birth calculated from given data
 - Age_years
 - TodayDate DateOfBirth

This added another angel to observation

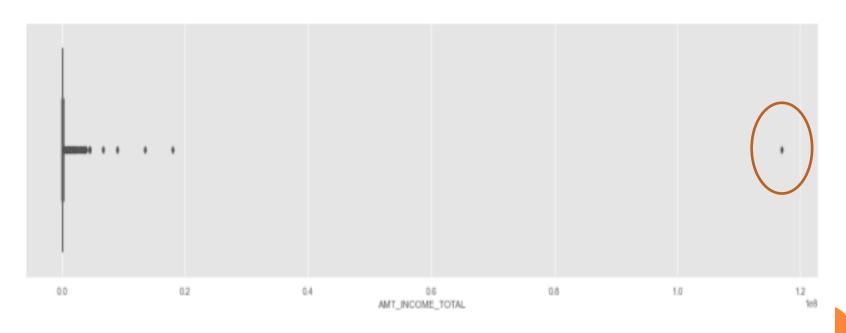
Age Wise Analysis



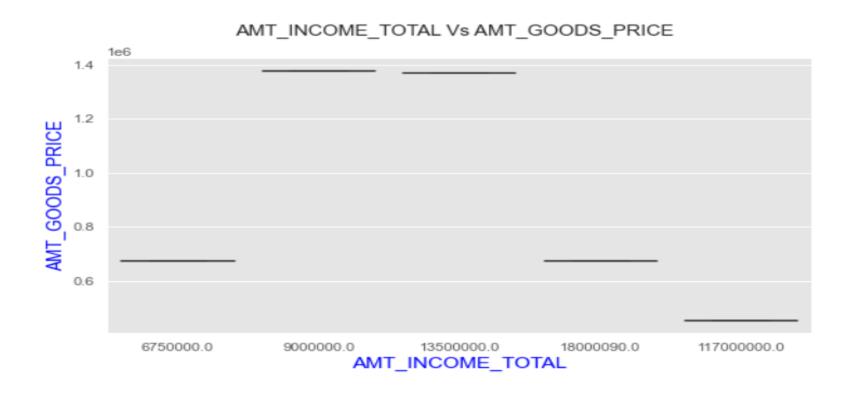
• Customer having 39 years of Age are highest in number (more than 6% of total) who has applied for loan data

ANALYSIS ON INCOME OF CLIENT

• Income in circle is clearly outlier, lets try to find fraud from income



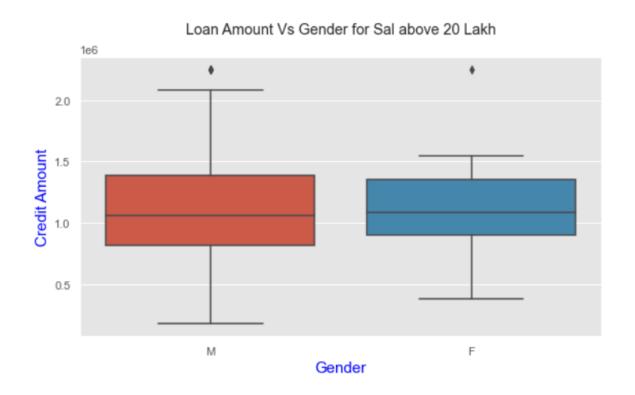
• To get more insight let's loan demand for salary more than 60lakh



• These are some strange cases where customer salary > 60 Lakhs, still they are going for some small loans, which can be checked

GENDER VS SALARY > 20 LAKH

• Median of both gender is similar but maximum value for both gender have huge differnce.



TARGET

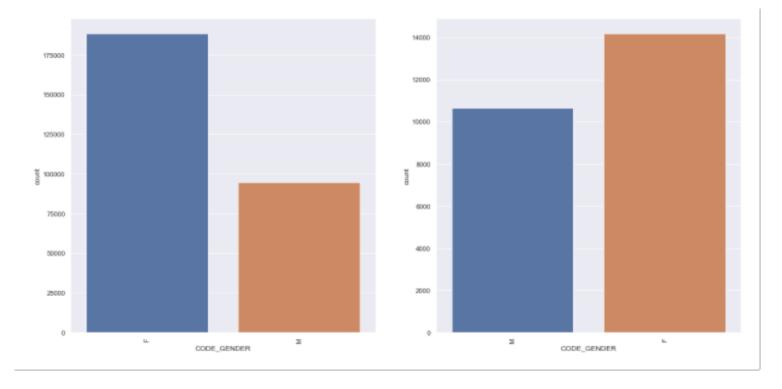
- Target ==1:
 - Client with payment difficulty
 - Count = 24825
- Target == 0:
 - Other client
 - Count= 282686

There is huge imbalance in target variable.

HANDLING IMBALANCE TARGET VARIABLE

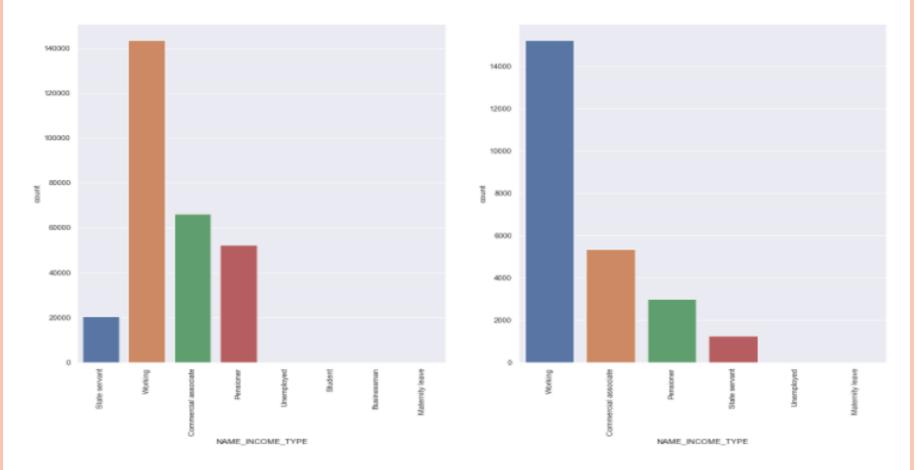
- After normalizing count:
 - 0 0.919271
 - 1 0.080729
 - We can observe that maximum loan holders don't have problem in repaying it.
- To resolve this imbalance we will divide target into 2 parts
- target0 = data.loc[data.TARGET == 0]
- target1 = data.loc[data.TARGET == 1]

CODE_GENDER VS TARGET0/TARGET1



Comaparing the Payment Difficulties and Non Payment Difficulties on the basis of Gender, we observe that Females are the majority in both the cases although there is an increase in the percentage in Male Payment Difficulties from Non-Payment Difficulties

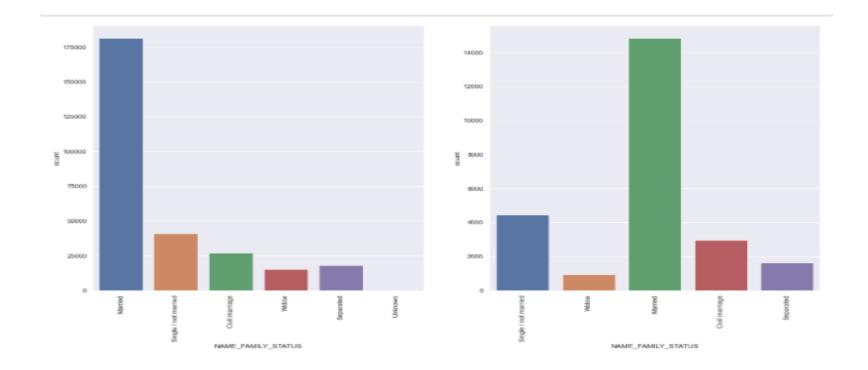
INCOME_SOURCE VS TARGET0/TARGET1



• We observe a decrease in the percentage of Payment Difficulties who are pentioners and an increase in the percentage of Payment Difficulties who are working when compared the percentages of both Payment Difficulties and non-Payment Difficulties.

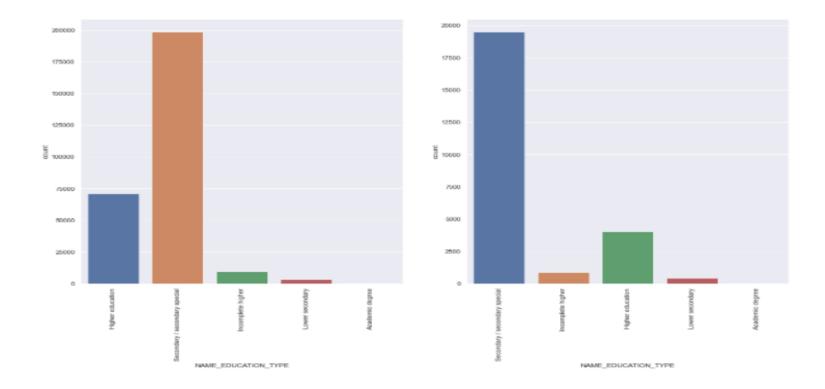


Family_status VS target0/target1



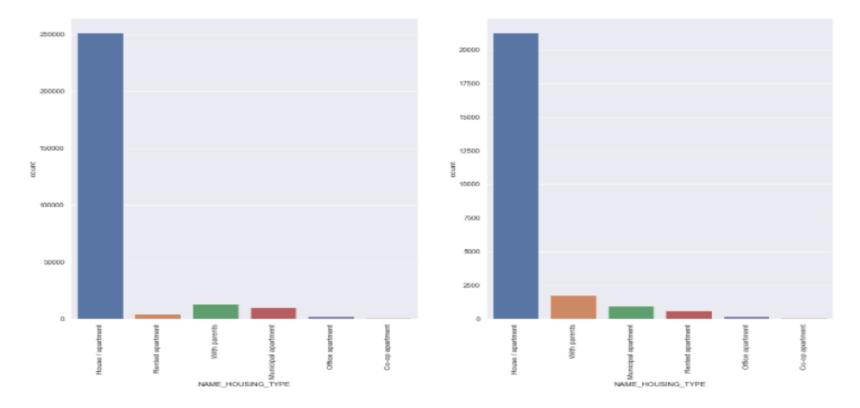
• We observe a decrease in the percentage of married and widowed with Loan Payment Difficulties and an increase in the the percentage of single and civil married with Loan Payment Difficulties when comapred with the percentages of both Loan Payment Difficulties and Loan Non-Payment Difficulties

EDUCATION VS TARGETO/TARGET1



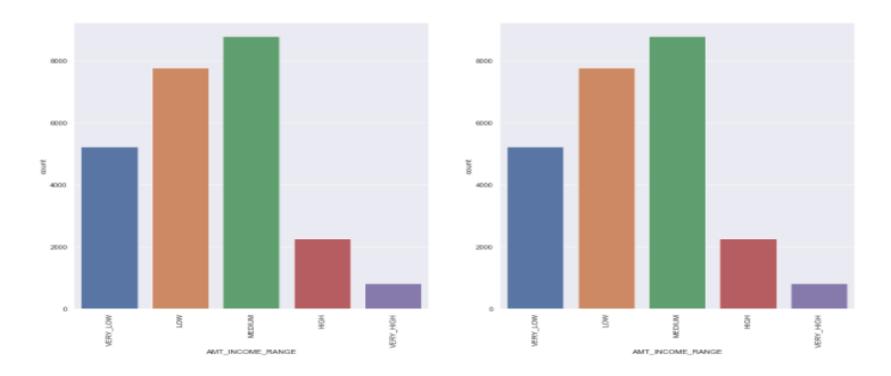
• We observe an increase in percentage of Loan Payment Difficulties whose educational qualifications are secondary/secondary special and a decrease in the percentage of Loan Payment Difficulties who have completed higher education when compared with the percentages of Loan Payment Difficulties and Loan Non-Payment Difficulties

HOUSE_TYPE VS TARGET0/TARGET1



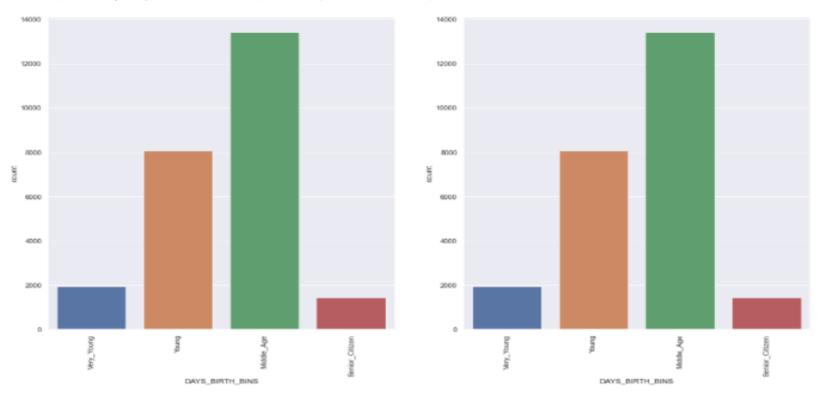
• We observe an increase in percentage of Loan Payment Difficulties whose educational qualifications are secondary/secondary special and a decrease in the percentage of Loan Payment Difficulties who have completed higher education when compared with the percentages of Loan Payment Difficulties and Loan Non-Payment Difficulties

INCOME_RANGE VS TARGET0/TARGET1



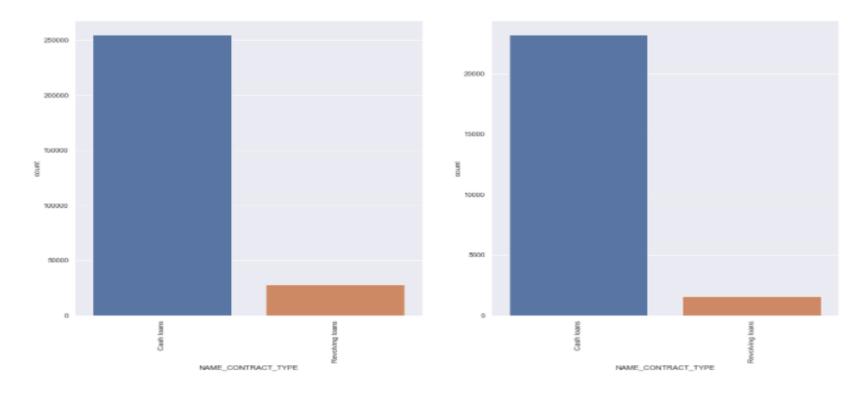
• We observe an increase in the percentage of Loan Payment Difficulties whose income is low when compared with the percentages of Payment Difficulties and Loan-Non Payment Difficulties

AGE VS TARGETO/TARGET1



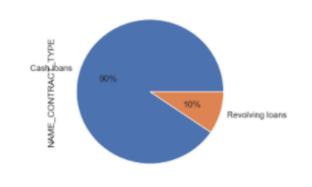
 We observe that there is an increase in the percentage of Loan Payment Difficulties who are young in age when compared to the percentages of Payment Difficulties and Loan-Non Payment Difficulties.

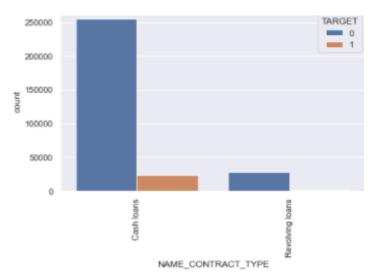
LOAN_TYPE VS TARGET0/TARGET1



• We can observe that cash loans are preffered by both Loan Payment Difficulties and Loan-Non Payment Difficulties although there is a decrease in the percentage of Payment Difficulties who opt for revolving loans.

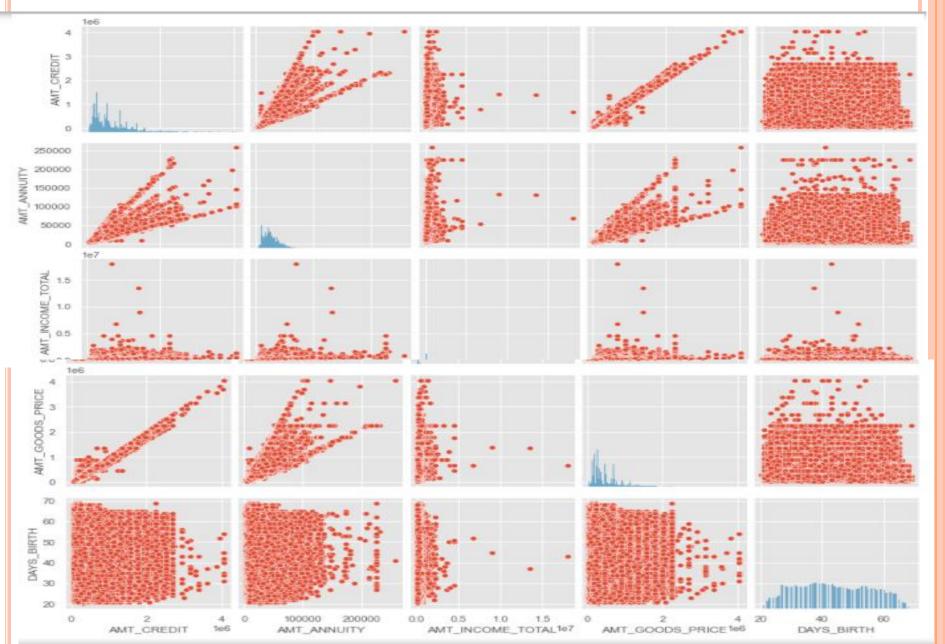
FUNCTION FOR ALL COLUMNS IN LIST





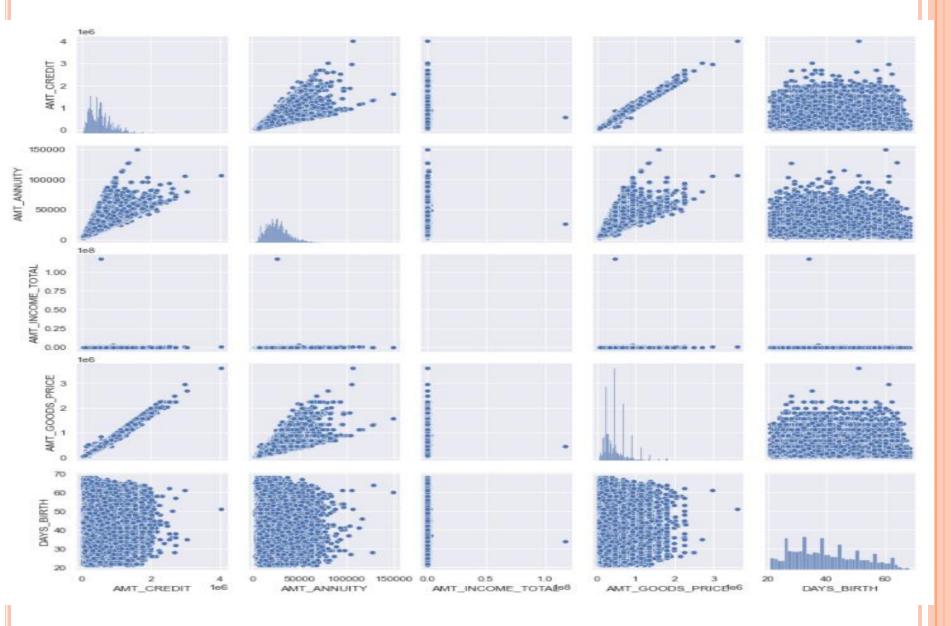


BIVARIENT ANALYSIS ON TARGETO



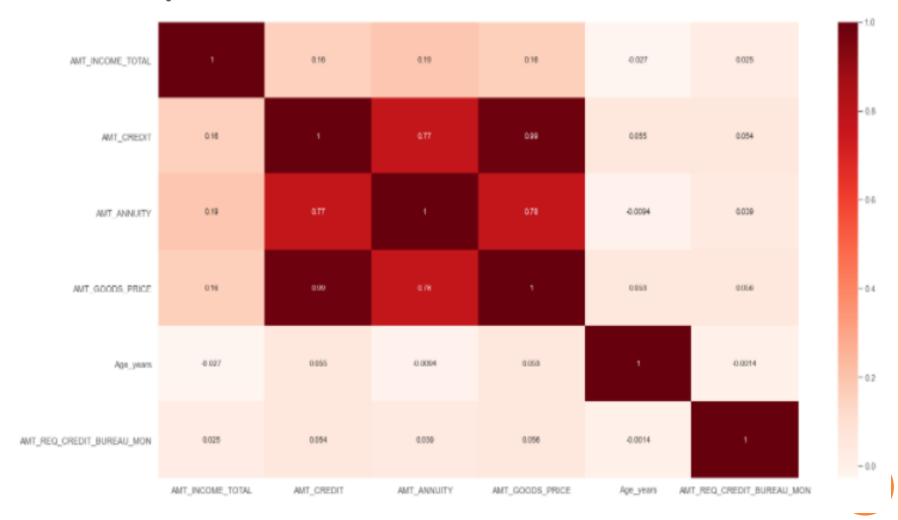
- AMT_CREDIT & AMT_ANNUITY are directly proportion to each other
- AMT_CREDIT & AMT_INCOME_TOTAL less income tends to more credit
- AMT_CREDIT & AMT_GOODS_PRIZE are directly proportional
- AMT_CREDIT & DAYS_BIRTH younger the person tends to take more loan

BIVARIENT ANALYSIS ON TARGET1

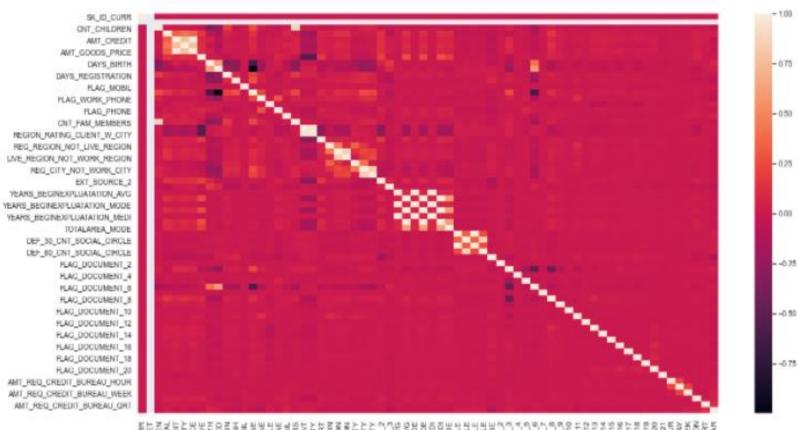


- AMT_CREDIT & AMT_ANNUITY are directly proportion to each other
- AMT_CREDIT & AMT_INCOME_TOTAL less income tends to more credit
- AMT_CREDIT & AMT_GOODS_PRIZE are directly proportional
- AMT_CREDIT & DAYS_BIRTH younger the person tends to take more loan

CORRELATION ON AMT_INCOME_TOTAL, AMT_CREDIT, AMT_ANNUITY, AMT_GOODS_PRICE, AGE_YEARS, AMT_REQ_CREDIT_BUREAU_MON



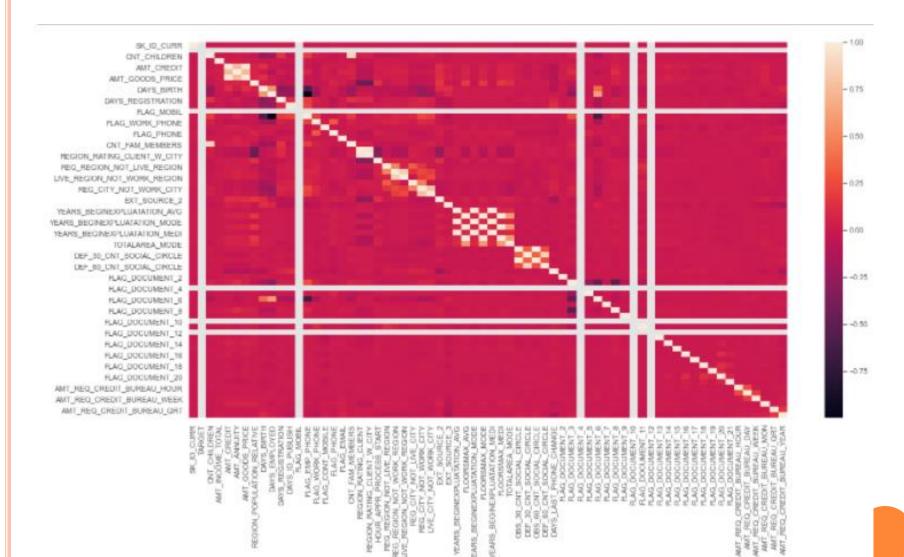
HEATMAP ON TARGETO



SK. JO. CURREN
AMT. MANCHEN
AMT. CHECKEN
AMT. ANACHTH
AMT

MAT PRO, CHEDIT BURBALI MAT PRO, CREDIT BURBALI MAT PRO, CREDIT BURBALI MAT REO, CREDIT BURBALI MAT REO, CREDIT BURBALI MAT PRO, CREDIT BURBALI MAT PRO, CREDIT BURBALI MAT PRO, CREDIT BURBALI MAT PRO, CREDIT BURBALI MAT

HEATMAP ON TARGET1

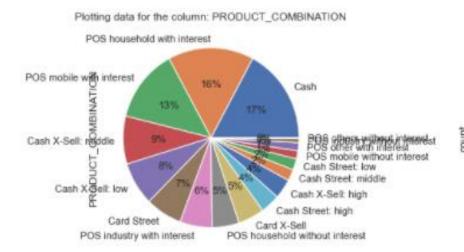


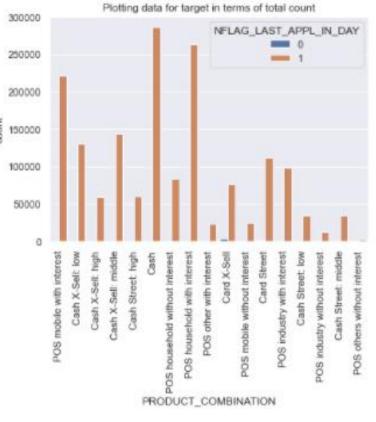
ANALYSIS ON PREVIOUS_APPLICATION.CSV

- After checking null values there are only 2 variable which have null values
- Product_combination and amt_credit
- We analysed it on bases of NFLAG_LAST_APPL_IN_DAY (as a target)

FUNCTION FOR UNIVARIENT ANALYSIS

Pltting PRODUCT COMBINATION





MERGING APPLICATION_DATA AND PREVIOUS_APPLICATION

 Both dataframe have common variable SK_ID_CURR

- We have merged both on SK_ID_CURR variable
- Used pandas function merge to merge both dataframe for futher analysis