

```
In [1]: 1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import itertools
6 %matplotlib inline
7 import warnings
8 warnings.filterwarnings('ignore')
```

```
In [2]: 1 pd.set_option('display.max_rows',500)
2 pd.set_option('display.max_columns',500)
3 pd.set_option('display.width',100)
```

```
In [3]: 1 app_data=pd.read_csv(r'S:\DOCS\dataset\application_data.csv')
```

```
In [4]: 1 app_data
```

Out[4]:

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR
0	100002	1	Cash loans	M	N
1	100003	0	Cash loans	F	N
2	100004	0	Revolving loans	M	Y
3	100006	0	Cash loans	F	N
4	100007	0	Cash loans	M	N
...
307506	456251	0	Cash loans	M	N
307507	456252	0	Cash loans	F	N
307508	456253	0	Cash loans	F	N
307509	456254	1	Cash loans	F	N
307510	456255	0	Cash loans	F	N

307511 rows × 122 columns

```
In [5]: 1 app_data.shape
```

Out[5]: (307511, 122)

```
In [6]: 1 app_data.info(verbose=True)
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 307511 entries, 0 to 307510
Data columns (total 122 columns):
 #   Column           Dtype  
 --- 
 0   SK_ID_CURR        int64  
 1   TARGET            int64  
 2   NAME_CONTRACT_TYPE object 
 3   CODE_GENDER        object 
 4   FLAG_OWN_CAR       object 
 5   FLAG_OWN_REALTY   object 
 6   CNT_CHILDREN       int64  
 7   AMT_INCOME_TOTAL  float64 
 8   AMT_CREDIT          float64 
 9   AMT_ANNUITY         float64 
 10  AMT_GOODS_PRICE    float64 
 11  NAME_TYPE_SUITE    object 
 12  NAME_INCOME_TYPE   object 
 13  NAME_EDUCATION_TYPE object 
 14  NAME_FAMILY_STATUS  object 
```

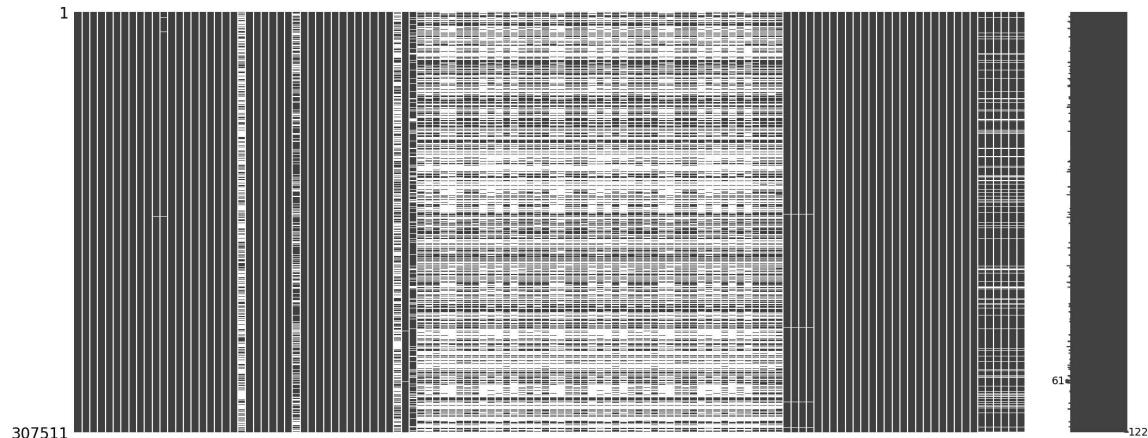
```
In [7]: 1 app_data.describe()
```

Out[7]:	SK_ID_CURR	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	A
count	307511.000000	307511.000000	307511.000000	3.075110e+05	3.075110e+05	3
mean	278180.518577	0.080729	0.417052	1.687979e+05	5.990260e+05	3
std	102790.175348	0.272419	0.722121	2.371231e+05	4.024908e+05	3
min	100002.000000	0.000000	0.000000	2.565000e+04	4.500000e+04	3
25%	189145.500000	0.000000	0.000000	1.125000e+05	2.700000e+05	3
50%	278202.000000	0.000000	0.000000	1.471500e+05	5.135310e+05	3
75%	367142.500000	0.000000	1.000000	2.025000e+05	8.086500e+05	3
max	456255.000000	1.000000	19.000000	1.170000e+08	4.050000e+06	3

```
In [8]: 1 import missingno as mn
```

In [9]: 1 mn.matrix(app_data)

Out[9]: <Axes: >



In [10]: 1 # many values are missing in dataset

In [11]: 1 round(app_data.isnull().sum()/app_data.shape[0] *100)

Out[11]:

SK_ID_CURR	0.0
TARGET	0.0
NAME_CONTRACT_TYPE	0.0
CODE_GENDER	0.0
FLAG_OWN_CAR	0.0
FLAG_OWN_REALTY	0.0
CNT_CHILDREN	0.0
AMT_INCOME_TOTAL	0.0
AMT_CREDIT	0.0
AMT_ANNUITY	0.0
AMT_GOODS_PRICE	0.0
NAME_TYPE_SUITE	0.0
NAME_INCOME_TYPE	0.0
NAME_EDUCATION_TYPE	0.0
NAME_FAMILY_STATUS	0.0
NAME_HOUSING_TYPE	0.0
REGION_POPULATION_RELATIVE	0.0
DAYS_BIRTH	0.0
DAYS_EMPLOYED	0.0
DAYS_REGISTRATION	0.0

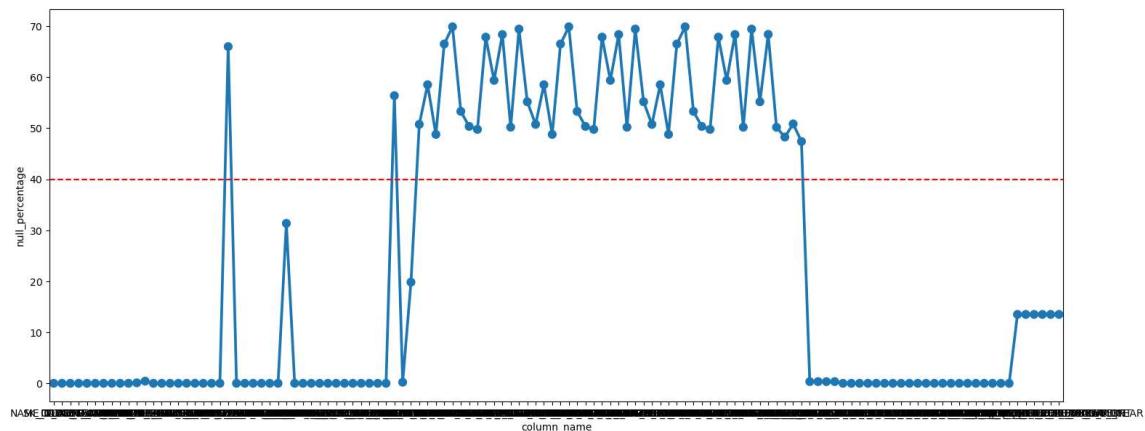
In [12]: 1 # many variables have missing value more than 40 %

In [13]: 1 null_data= (app_data.isnull().sum()/app_data.shape[0] *100).reset_index

In [14]: 1 null_data.columns=['column_name', 'null_percentage']

```
In [15]: 1 fig=plt.figure(figsize=(18,7))
2 ax=sns.pointplot(data=null_data, x='column_name', y='null_percentage')
3 ax.axhline(40, ls='--',color='red')
```

Out[15]: <matplotlib.lines.Line2D at 0x166806e82d0>



```
In [16]: 1 # the column point above red Line have missing more than 40%
```

```
In [17]: 1 null_value_40 = null_data[null_data['null_percentage']>=40]
```

```
In [18]: 1 len(null_value_40)
```

Out[18]: 49

```
In [19]: 1 source=app_data[['EXT_SOURCE_1','EXT_SOURCE_2','EXT_SOURCE_3','TARGET']]
```

```
In [20]: 1 ax=sns.heatmap(source.corr(), annot=True)
```



```
In [21]: 1 # no correlation between ext_sources and target
```

```
In [22]: 1 unwanted_app= null_value_40['column_name'].tolist() +['EXT_SOURCE_2', 'EX']
```

```
In [23]: 1 len(unwanted_app)
```

Out[23]: 51

In [24]:

```

1 flag_doc=['FLAG_DOCUMENT_2','FLAG_DOCUMENT_3','FLAG_DOCUMENT_4','FLAG_D
2 'FLAG_DOCUMENT_6','FLAG_DOCUMENT_7','FLAG_DOCUMENT_8','FLAG_DOCUMENT_9'
3 'FLAG_DOCUMENT_11','FLAG_DOCUMENT_12','FLAG_DOCUMENT_13','FLAG_DOCUMENT
4 'FLAG_DOCUMENT_16','FLAG_DOCUMENT_17','FLAG_DOCUMENT_18','FLAG_DOCUMENT
5 'FLAG_DOCUMENT_21']
6 flag_df=app_data[flag_doc + ['TARGET']]
7 flag_df['TARGET']=flag_df['TARGET'].replace({1:'defaulter', 0:'repayer'}
8 length=len(flag_doc)
9 fig=plt.figure(figsize=(21,24))
10
11 for i,j in itertools.zip_longest(flag_doc,range(length)):
12     plt.subplot(5,4,j+1)
13     ax = sns.countplot(data=flag_df,x=i, hue=flag_df["TARGET"],palette=

```



The above graph shows that in most of the loan application cases, clients who applied for loans has not submitted FLAG_DOCUMENT_X except FLAG_DOCUMENT_3. Thus, Except for FLAG_DOCUMENT_3, we can delete rest of the columns. Data shows if borrower has submitted FLAG_DOCUMENT_3 then there is a less chance of defaulting the loan.

```
In [25]: 1 flag_doc.remove('FLAG_DOCUMENT_3')
```

In [26]:

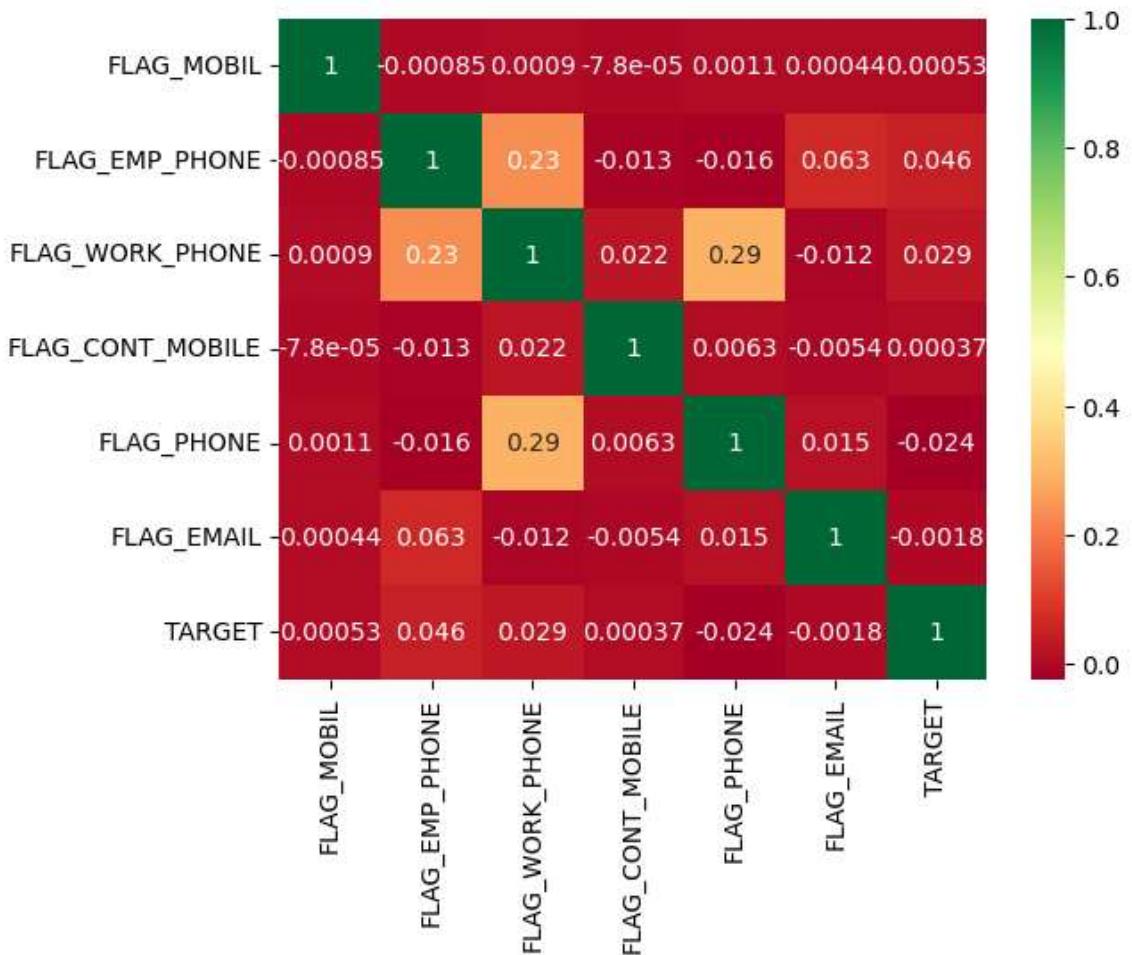
```
1 unwanted_app=unwanted_app + flag_doc  
2 unwanted_app
```

```
Out[26]: ['OWN_CAR_AGE',
  'EXT_SOURCE_1',
  'APARTMENTS_AVG',
  'BASEMENTAREA_AVG',
  'YEARS_BEGINEXPLUATATION_AVG',
  'YEARS_BUILD_AVG',
  'COMMONAREA_AVG',
  'ELEVATORS_AVG',
  'ENTRANCES_AVG',
  'FLOORSMAX_AVG',
  'FLOORSMIN_AVG',
  'LANDAREA_AVG',
  'LIVINGAPARTMENTS_AVG',
  'LIVINGAREA_AVG',
  'NONLIVINGAPARTMENTS_AVG',
  'NONLIVINGAREA_AVG',
  'APARTMENTS_MODE',
  'BASEMENTAREA_MODE',
  'YEARS_BEGINEXPLUATATION_MODE',
  'YEARS_BUILD_MODE',
  'COMMONAREA_MODE',
  'ELEVATORS_MODE',
  'ENTRANCES_MODE',
  'FLOORSMAX_MODE',
  'FLOORSMIN_MODE',
  'LANDAREA_MODE',
  'LIVINGAPARTMENTS_MODE',
  'LIVINGAREA_MODE',
  'NONLIVINGAPARTMENTS_MODE',
  'NONLIVINGAREA_MODE',
  'APARTMENTS_MEDI',
  'BASEMENTAREA_MEDI',
  'YEARS_BEGINEXPLUATATION_MEDI',
  'YEARS_BUILD_MEDI',
  'COMMONAREA_MEDI',
  'ELEVATORS_MEDI',
  'ENTRANCES_MEDI',
  'FLOORSMAX_MEDI',
  'FLOORSMIN_MEDI',
  'LANDAREA_MEDI',
  'LIVINGAPARTMENTS_MEDI',
  'LIVINGAREA_MEDI',
  'NONLIVINGAPARTMENTS_MEDI',
  'NONLIVINGAREA_MEDI',
  'FONDKAPREMONT_MODE',
  'HOUSETYPE_MODE',
  'TOTALAREA_MODE',
  'WALLSMATERIAL_MODE',
  'EMERGENCYSTATE_MODE',
  'EXT_SOURCE_2',
  'EXT_SOURCE_3',
  'FLAG_DOCUMENT_2',
  'FLAG_DOCUMENT_4',
  'FLAG_DOCUMENT_5',
  'FLAG_DOCUMENT_6',
  'FLAG_DOCUMENT_7',
  'FLAG_DOCUMENT_8',
  'FLAG_DOCUMENT_9',
  'FLAG_DOCUMENT_10',
  'FLAG_DOCUMENT_11',
  'FLAG_DOCUMENT_12']
```

```
'FLAG_DOCUMENT_13',
'FLAG_DOCUMENT_14',
'FLAG_DOCUMENT_15',
'FLAG_DOCUMENT_16',
'FLAG_DOCUMENT_17',
'FLAG_DOCUMENT_18',
'FLAG_DOCUMENT_19',
'FLAG_DOCUMENT_20',
'FLAG_DOCUMENT_21']
```

In [27]:

```
1 contact_col=['FLAG_MOBIL', 'FLAG_EMP_PHONE', 'FLAG_WORK_PHONE', 'FLAG_CONT_MOBILE',
2                 'FLAG_PHONE', 'FLAG_EMAIL', 'TARGET']
3
4 contact_corr=app_data[contact_col].corr()
5 ax=sns.heatmap(contact_corr,annot=True,cmap ="RdYlGn")
```



There is no correlation between flags of mobile phone, email etc with loan repayment; thus these columns can be deleted

In [28]:

```
1 contact_col.remove('TARGET')
```

In [29]:

```
1 unwanted_app=unwanted_app + contact_col
```

In [30]:

```
1 app_data.drop(unwanted_app,axis=1, inplace=True)
```

In [31]: 1 app_data.shape

Out[31]: (307511, 46)

In [32]: 1 app_data.info(verbose=True)

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 307511 entries, 0 to 307510
Data columns (total 46 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   SK_ID_CURR       307511 non-null   int64  
 1   TARGET           307511 non-null   int64  
 2   NAME_CONTRACT_TYPE 307511 non-null   object  
 3   CODE_GENDER      307511 non-null   object  
 4   FLAG_OWN_CAR     307511 non-null   object  
 5   FLAG_OWN_REALTY  307511 non-null   object  
 6   CNT_CHILDREN     307511 non-null   int64  
 7   AMT_INCOME_TOTAL 307511 non-null   float64 
 8   AMT_CREDIT        307511 non-null   float64 
 9   AMT_ANNUITY       307499 non-null   float64 
 10  AMT_GOODS_PRICE   307233 non-null   float64 
 11  NAME_TYPE_SUITE   306219 non-null   object  
 12  NAME_INCOME_TYPE  307511 non-null   object  
 13  NAME_EDUCATION_TYPE 307511 non-null   object  
 14  NAME_FAMILY_STATUS 307511 non-null   object  
 15  NAME_HOUSING_TYPE 307511 non-null   object  
 16  REGION_POPULATION_RELATIVE 307511 non-null   float64 
 17  DAYS_BIRTH        307511 non-null   int64  
 18  DAYS_EMPLOYED     307511 non-null   int64  
 19  DAYS_REGISTRATION 307511 non-null   float64 
 20  DAYS_ID_PUBLISH   307511 non-null   int64  
 21  OCCUPATION_TYPE    211120 non-null   object  
 22  CNT_FAM_MEMBERS   307509 non-null   float64 
 23  REGION_RATING_CLIENT 307511 non-null   int64  
 24  REGION_RATING_CLIENT_W_CITY 307511 non-null   int64  
 25  WEEKDAY_APPR_PROCESS_START 307511 non-null   object  
 26  HOUR_APPR_PROCESS_START 307511 non-null   int64  
 27  REG_REGION_NOT_LIVE_REGION 307511 non-null   int64  
 28  REG_REGION_NOT_WORK_REGION 307511 non-null   int64  
 29  LIVE_REGION_NOT_WORK_REGION 307511 non-null   int64  
 30  REG_CITY_NOT_LIVE_CITY 307511 non-null   int64  
 31  REG_CITY_NOT_WORK_CITY 307511 non-null   int64  
 32  LIVE_CITY_NOT_WORK_CITY 307511 non-null   int64  
 33  ORGANIZATION_TYPE   307511 non-null   object  
 34  OBS_30_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 35  DEF_30_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 36  OBS_60_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 37  DEF_60_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 38  DAYS_LAST_PHONE_CHANGE 307510 non-null   float64 
 39  FLAG_DOCUMENT_3     307511 non-null   int64  
 40  AMT_REQ_CREDIT_BUREAU_HOUR 265992 non-null   float64 
 41  AMT_REQ_CREDIT_BUREAU_DAY 265992 non-null   float64 
 42  AMT_REQ_CREDIT_BUREAU_WEEK 265992 non-null   float64 
 43  AMT_REQ_CREDIT_BUREAU_MON 265992 non-null   float64 
 44  AMT_REQ_CREDIT_BUREAU_QRT 265992 non-null   float64 
 45  AMT_REQ_CREDIT_BUREAU_YEAR 265992 non-null   float64 
dtypes: float64(18), int64(16), object(12)
memory usage: 107.9+ MB
```

Standardize Values

Convert DAYS_DECISION,DAYS_EMPLOYED, DAYS_REGISTRATION,DAYS_ID_PUBLISH from negative to positive as days cannot be negative.

```
In [33]: 1 date_col=['DAYS_BIRTH', 'DAYS_EMPLOYED', 'DAYS_REGISTRATION', 'DAYS_ID_PUE
```

```
In [34]: 1 for i in date_col:  
2     app_data[i] = abs(app_data[i])
```

```
In [35]: 1 # Binning Numerical Columns to create a categorical column
```

```
In [36]: 1 # Creating bins for income amount  
2 app_data['AMT_INCOME_TOTAL'] = app_data['AMT_INCOME_TOTAL'] / 100000  
3  
4 bins= [0,1,2,3,4,5,6,7,8,9,10,11]  
5 slot=['0-100k', '100k-200k', '200k-300k', '300k-400k', '400k-500k', '500k-600k',  
6     , '800k-900k', '900k-1M', 'above 1M']  
7  
8 app_data['AMT_INCOME_RANGE']= pd.cut(app_data['AMT_INCOME_TOTAL'], bins, labels=slot)
```

```
In [37]: 1 app_data['AMT_INCOME_RANGE'].value_counts(normalize=True)*100
```

```
Out[37]: 100k-200k    50.735000  
200k-300k    21.210691  
0-100k      20.729695  
300k-400k     4.776116  
400k-500k     1.744669  
500k-600k     0.356354  
600k-700k     0.282805  
800k-900k     0.096980  
700k-800k     0.052721  
900k-1M       0.009112  
above 1M       0.005858  
Name: AMT_INCOME_RANGE, dtype: float64
```

```
In [38]: 1 # Creating bins for Credit amount
```

```
In [39]: 1 app_data['AMT_CREDIT']=app_data['AMT_CREDIT']/100000  
2  
3 bins = [0,1,2,3,4,5,6,7,8,9,10,100]  
4 slots = ['0-100K', '100K-200K', '200K-300K', '300K-400K', '400K-500K', '500K-600K',  
5     , '800K-900K', '900K-1M', '1M Above']  
6  
7 app_data['AMT_CREDIT_RANGE']=pd.cut(app_data['AMT_CREDIT'], bins=bins, labels=slots)
```

In [40]: 1 app_data['AMT_CREDIT_RANGE'].value_counts(normalize=True)*100

Out[40]:

200k-300k	17.824728
1M Above	16.254703
500k-600k	11.131960
400k-500k	10.418489
100K-200K	9.801275
300k-400k	8.564897
600k-700k	7.820533
800k-900k	7.086576
700k-800k	6.241403
900k-1M	2.902986
0-100K	1.952450

Name: AMT_CREDIT_RANGE, dtype: float64

In [41]: 1 # 16 % people have taken Loan more than 1 million

In [42]: 1 # Creating bins for Age

```

2
3 app_data[ 'AGE' ] = app_data[ 'DAYS_BIRTH' ] // 365
4
5 bins=[0,20,30,40,50,100]
6 slots=['0-20','20-30','30-40','40-50','above 50']
7
8 app_data[ 'AGE_GROUP' ]=pd.cut(app_data[ 'AGE' ],bins=slots)

```

In [43]: 1 app_data['AGE_GROUP'].value_counts(normalize=True)*100

Out[43]:

above 50	31.604398
30-40	27.028952
40-50	24.194582
20-30	17.171743
0-20	0.000325

Name: AGE_GROUP, dtype: float64

31% loan applicants have age above 50 years. More than 55% of loan applicants have age over 40 years.

In [44]: 1 # Creating bins for Employment Time

```

2 app_data[ 'YEARS_EMPLOYED' ] = app_data[ 'DAYS_EMPLOYED' ] // 365
3 bins = [0,5,10,20,30,40,50,60,150]
4 slots = ['0-5','5-10','10-20','20-30','30-40','40-50','50-60','60 above'
5
6 app_data[ 'EMPLOYMENT_YEAR' ]=pd.cut(app_data[ 'YEARS_EMPLOYED' ],bins=bins)

```

```
In [45]: 1 app_data[ 'EMPLOYMENT_YEAR' ].value_counts(normalize=True)*100
```

```
Out[45]: 0-5      55.582363
5-10     24.966441
10-20    14.564315
20-30    3.750117
30-40    1.058720
40-50    0.078044
50-60    0.000000
60 above 0.000000
Name: EMPLOYMENT_YEAR, dtype: float64
```

More than 55% of the loan applicants have work experience within 0-5 years and almost 80% of them have less than 10 years of work experience

```
In [46]: 1 #Checking the number of unique values each column possess to identify c
```

```
In [47]: 1 app_data.nunique().sort_values()
```

```
Out[47]: LIVE_CITY_NOT_WORK_CITY                2
TARGET                                         2
NAME_CONTRACT_TYPE                           2
REG_REGION_NOT_LIVE_REGION                  2
FLAG_OWN_CAR                                2
FLAG_OWN_REALTY                             2
REG_REGION_NOT_WORK_REGION                 2
LIVE_REGION_NOT_WORK_REGION                2
FLAG_DOCUMENT_3                            2
REG_CITY_NOT_LIVE_CITY                     2
REG_CITY_NOT_WORK_CITY                    2
REGION_RATING_CLIENT                      3
CODE_GENDER                                 3
REGION_RATING_CLIENT_W_CITY               3
AMT_REQ_CREDIT_BUREAU_HOUR                5
NAME_EDUCATION_TYPE                       5
AGE_GROUP                                  5
NAME_FAMILY_STATUS                         6
NAME_HOUSING_TYPE                         6
EMPLOYMENT_YEAR                           6
WEEKDAY_APPR_PROCESS_START                7
NAME_TYPE_SUITE                            7
NAME_INCOME_TYPE                           8
AMT_REQ_CREDIT_BUREAU_WEEK                9
AMT_REQ_CREDIT_BUREAU_DAY                 9
DEF_60_CNT_SOCIAL_CIRCLE                  9
DEF_30_CNT_SOCIAL_CIRCLE                 10
AMT_CREDIT_RANGE                           11
AMT_INCOME_RANGE                          11
AMT_REQ_CREDIT_BUREAU_QRT                 11
CNT_CHILDREN                               15
CNT_FAM_MEMBERS                           17
OCCUPATION_TYPE                           18
HOUR_APPR_PROCESS_START                  24
AMT_REQ_CREDIT_BUREAU_MON                 24
AMT_REQ_CREDIT_BUREAU_YEAR                25
OBS_60_CNT_SOCIAL_CIRCLE                 33
OBS_30_CNT_SOCIAL_CIRCLE                 33
AGE                                      50
YEARS_EMPLOYED                            51
ORGANIZATION_TYPE                        58
REGION_POPULATION_RELATIVE                81
AMT_GOODS_PRICE                           1002
AMT_INCOME_TOTAL                          2548
DAYS_LAST_PHONE_CHANGE                   3773
AMT_CREDIT                                5603
DAYS_ID_PUBLISH                           6168
DAYS_EMPLOYED                             12574
AMT_ANNUITY                                13672
DAYS_REGISTRATION                         15688
DAYS_BIRTH                                 17460
SK_ID_CURR                                307511
dtype: int64
```

data type conversion

In [48]: 1 app_data.info()

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 307511 entries, 0 to 307510
Data columns (total 52 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   SK_ID_CURR       307511 non-null   int64  
 1   TARGET           307511 non-null   int64  
 2   NAME_CONTRACT_TYPE 307511 non-null   object  
 3   CODE_GENDER      307511 non-null   object  
 4   FLAG_OWN_CAR     307511 non-null   object  
 5   FLAG_OWN_REALTY  307511 non-null   object  
 6   CNT_CHILDREN     307511 non-null   int64  
 7   AMT_INCOME_TOTAL 307511 non-null   float64 
 8   AMT_CREDIT        307511 non-null   float64 
 9   AMT_ANNUITY       307499 non-null   float64 
 10  AMT_GOODS_PRICE   307233 non-null   float64 
 11  NAME_TYPE_SUITE   306219 non-null   object  
 12  NAME_INCOME_TYPE  307511 non-null   object  
 13  NAME_EDUCATION_TYPE 307511 non-null   object  
 14  NAME_FAMILY_STATUS 307511 non-null   object  
 15  NAME_HOUSING_TYPE 307511 non-null   object  
 16  REGION_POPULATION_RELATIVE 307511 non-null   float64 
 17  DAYS_BIRTH        307511 non-null   int64  
 18  DAYS_EMPLOYED     307511 non-null   int64  
 19  DAYS_REGISTRATION 307511 non-null   float64 
 20  DAYS_ID_PUBLISH   307511 non-null   int64  
 21  OCCUPATION_TYPE    211120 non-null   object  
 22  CNT_FAM_MEMBERS    307509 non-null   float64 
 23  REGION_RATING_CLIENT 307511 non-null   int64  
 24  REGION_RATING_CLIENT_W_CITY 307511 non-null   int64  
 25  WEEKDAY_APPR_PROCESS_START 307511 non-null   object  
 26  HOUR_APPR_PROCESS_START 307511 non-null   int64  
 27  REG_REGION_NOT_LIVE_REGION 307511 non-null   int64  
 28  REG_REGION_NOT_WORK_REGION 307511 non-null   int64  
 29  LIVE_REGION_NOT_WORK_REGION 307511 non-null   int64  
 30  REG_CITY_NOT_LIVE_CITY 307511 non-null   int64  
 31  REG_CITY_NOT_WORK_CITY 307511 non-null   int64  
 32  LIVE_CITY_NOT_WORK_CITY 307511 non-null   int64  
 33  ORGANIZATION_TYPE    307511 non-null   object  
 34  OBS_30_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 35  DEF_30_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 36  OBS_60_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 37  DEF_60_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 38  DAYS_LAST_PHONE_CHANGE 307510 non-null   float64 
 39  FLAG_DOCUMENT_3      307511 non-null   int64  
 40  AMT_REQ_CREDIT_BUREAU_HOUR 265992 non-null   float64 
 41  AMT_REQ_CREDIT_BUREAU_DAY 265992 non-null   float64 
 42  AMT_REQ_CREDIT_BUREAU_WEEK 265992 non-null   float64 
 43  AMT_REQ_CREDIT_BUREAU_MON 265992 non-null   float64 
 44  AMT_REQ_CREDIT_BUREAU_QRT 265992 non-null   float64 
 45  AMT_REQ_CREDIT_BUREAU_YEAR 265992 non-null   float64 
 46  AMT_INCOME_RANGE      307279 non-null   category 
 47  AMT_CREDIT_RANGE       307511 non-null   category 
 48  AGE                   307511 non-null   int64  
 49  AGE_GROUP             307511 non-null   category 
 50  YEARS_EMPLOYED        307511 non-null   int64  
 51  EMPLOYMENT_YEAR       224233 non-null   category 

dtypes: category(4), float64(18), int64(18), object(12)
memory usage: 113.8+ MB

```

```
In [49]: 1 #Numeric columns are already in int64 and float64 format. Hence proceed
```

```
In [50]: 1 categorical_columns = ['NAME_CONTRACT_TYPE', 'CODE_GENDER', 'NAME_TYPE_SUITE',
 2                           'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE', 'OCCUPATION_TYPE',
 3                           'ORGANIZATION_TYPE', 'FLAG_OWN_CAR', 'FLAG_OWN_REALTY',
 4                           'REG_CITY_NOT_LIVE_CITY', 'REG_CITY_NOT_WORK_CITY',
 5                           'LIVE_REGION_NOT_WORK_REGION', 'REGION_RATING_CLIENT',
 6                           'REGION_RATING_CLIENT_W_CITY']
 7
 8 for col in categorical_columns:
 9     app_data[col] = pd.Categorical(app_data[col])
```

```
In [51]: 1 app_data.info()
```

```

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 307511 entries, 0 to 307510
Data columns (total 52 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   SK_ID_CURR       307511 non-null   int64  
 1   TARGET           307511 non-null   int64  
 2   NAME_CONTRACT_TYPE 307511 non-null   category
 3   CODE_GENDER      307511 non-null   category
 4   FLAG_OWN_CAR     307511 non-null   category
 5   FLAG_OWN_REALTY  307511 non-null   category
 6   CNT_CHILDREN     307511 non-null   int64  
 7   AMT_INCOME_TOTAL 307511 non-null   float64 
 8   AMT_CREDIT        307511 non-null   float64 
 9   AMT_ANNUITY       307499 non-null   float64 
 10  AMT_GOODS_PRICE   307233 non-null   float64 
 11  NAME_TYPE_SUITE   306219 non-null   category
 12  NAME_INCOME_TYPE  307511 non-null   category
 13  NAME_EDUCATION_TYPE 307511 non-null   category
 14  NAME_FAMILY_STATUS 307511 non-null   category
 15  NAME_HOUSING_TYPE 307511 non-null   category
 16  REGION_POPULATION_RELATIVE 307511 non-null   float64 
 17  DAYS_BIRTH        307511 non-null   int64  
 18  DAYS_EMPLOYED     307511 non-null   int64  
 19  DAYS_REGISTRATION 307511 non-null   float64 
 20  DAYS_ID_PUBLISH   307511 non-null   int64  
 21  OCCUPATION_TYPE    211120 non-null   category
 22  CNT_FAM_MEMBERS   307509 non-null   float64 
 23  REGION_RATING_CLIENT 307511 non-null   category
 24  REGION_RATING_CLIENT_W_CITY 307511 non-null   category
 25  WEEKDAY_APPR_PROCESS_START 307511 non-null   category
 26  HOUR_APPR_PROCESS_START 307511 non-null   int64  
 27  REG_REGION_NOT_LIVE_REGION 307511 non-null   int64  
 28  REG_REGION_NOT_WORK_REGION 307511 non-null   category
 29  LIVE_REGION_NOT_WORK_REGION 307511 non-null   category
 30  REG_CITY_NOT_LIVE_CITY   307511 non-null   category
 31  REG_CITY_NOT_WORK_CITY  307511 non-null   category
 32  LIVE_CITY_NOT_WORK_CITY 307511 non-null   category
 33  ORGANIZATION_TYPE     307511 non-null   category
 34  OBS_30_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 35  DEF_30_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 36  OBS_60_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 37  DEF_60_CNT_SOCIAL_CIRCLE 306490 non-null   float64 
 38  DAYS_LAST_PHONE_CHANGE 307510 non-null   float64 
 39  FLAG_DOCUMENT_3       307511 non-null   int64  
 40  AMT_REQ_CREDIT_BUREAU_HOUR 265992 non-null   float64 
 41  AMT_REQ_CREDIT_BUREAU_DAY  265992 non-null   float64 
 42  AMT_REQ_CREDIT_BUREAU_WEEK 265992 non-null   float64 
 43  AMT_REQ_CREDIT_BUREAU_MON  265992 non-null   float64 
 44  AMT_REQ_CREDIT_BUREAU_QRT  265992 non-null   float64 
 45  AMT_REQ_CREDIT_BUREAU_YEAR 265992 non-null   float64 
 46  AMT_INCOME_RANGE      307279 non-null   category
 47  AMT_CREDIT_RANGE       307511 non-null   category
 48  AGE                   307511 non-null   int64  
 49  AGE_GROUP            307511 non-null   category
 50  YEARS_EMPLOYED       307511 non-null   int64  
 51  EMPLOYMENT_YEAR      224233 non-null   category
dtypes: category(23), float64(18), int64(11)
memory usage: 74.8 MB

```

```
In [52]: 1 # null value imputation
```

```
In [53]: 1 # checking% of null value
2
3 round(app_data.isnull().sum() / app_data.shape[0] *100,2)
```

Out[53]:

SK_ID_CURR	0.00
TARGET	0.00
NAME_CONTRACT_TYPE	0.00
CODE_GENDER	0.00
FLAG_OWN_CAR	0.00
FLAG_OWN_REALTY	0.00
CNT_CHILDREN	0.00
AMT_INCOME_TOTAL	0.00
AMT_CREDIT	0.00
AMT_ANNUITY	0.00
AMT_GOODS_PRICE	0.09
NAME_TYPE_SUITE	0.42
NAME_INCOME_TYPE	0.00
NAME_EDUCATION_TYPE	0.00
NAME_FAMILY_STATUS	0.00
NAME_HOUSING_TYPE	0.00
REGION_POPULATION_RELATIVE	0.00
DAYS_BIRTH	0.00
DAYS_EMPLOYED	0.00
DAYS_REGISTRATION	0.00
DAYS_ID_PUBLISH	0.00
OCCUPATION_TYPE	31.35
CNT_FAM_MEMBERS	0.00
REGION_RATING_CLIENT	0.00
REGION_RATING_CLIENT_W_CITY	0.00
WEEKDAY_APPR_PROCESS_START	0.00
HOUR_APPR_PROCESS_START	0.00
REG_REGION_NOT_LIVE_REGION	0.00
REG_REGION_NOT_WORK_REGION	0.00
LIVE_REGION_NOT_WORK_REGION	0.00
REG_CITY_NOT_LIVE_CITY	0.00
REG_CITY_NOT_WORK_CITY	0.00
LIVE_CITY_NOT_WORK_CITY	0.00
ORGANIZATION_TYPE	0.00
OBS_30_CNT_SOCIAL_CIRCLE	0.33
DEF_30_CNT_SOCIAL_CIRCLE	0.33
OBS_60_CNT_SOCIAL_CIRCLE	0.33
DEF_60_CNT_SOCIAL_CIRCLE	0.33
DAYS_LAST_PHONE_CHANGE	0.00
FLAG_DOCUMENT_3	0.00
AMT_REQ_CREDIT_BUREAU_HOUR	13.50
AMT_REQ_CREDIT_BUREAU_DAY	13.50
AMT_REQ_CREDIT_BUREAU_WEEK	13.50
AMT_REQ_CREDIT_BUREAU_MON	13.50
AMT_REQ_CREDIT_BUREAU_QRT	13.50
AMT_REQ_CREDIT_BUREAU_YEAR	13.50
AMT_INCOME_RANGE	0.08
AMT_CREDIT_RANGE	0.00
AGE	0.00
AGE_GROUP	0.00
YEARS_EMPLOYED	0.00
EMPLOYMENT_YEAR	27.08

dtype: float64

```
In [54]: 1 # impute name_type_suit
```

```
In [55]: 1 app_data['NAME_TYPE_SUITE'].describe()
```

```
Out[55]: count      306219  
unique       7  
top    Unaccompanied  
freq      248526  
Name: NAME_TYPE_SUITE, dtype: object
```

```
In [56]: 1 app_data['NAME_TYPE_SUITE'].fillna(app_data['NAME_TYPE_SUITE'].mode()[0])
```

```
In [57]: 1 round(app_data.isnull().sum() / app_data.shape[0] *100,2)
```

```
Out[57]: SK_ID_CURR           0.00  
TARGET              0.00  
NAME_CONTRACT_TYPE 0.00  
CODE_GENDER          0.00  
FLAG_OWN_CAR         0.00  
FLAG_OWN_REALTY     0.00  
CNT_CHILDREN         0.00  
AMT_INCOME_TOTAL    0.00  
AMT_CREDIT            0.00  
AMT_ANNUITY           0.00  
AMT_GOODS_PRICE       0.09  
NAME_TYPE_SUITE       0.00  
NAME_INCOME_TYPE     0.00  
NAME_EDUCATION_TYPE  0.00  
NAME_FAMILY_STATUS    0.00  
NAME_HOUSING_TYPE    0.00  
REGION_POPULATION_RELATIVE 0.00  
DAYS_BIRTH             0.00  
DAYS_EMPLOYED          0.00  
DAYS_REGISTRATION      0.00  
DAYS_ID_PUBLISH        0.00  
OCCUPATION_TYPE        31.35  
CNT_FAM_MEMBERS         0.00  
REGION_RATING_CLIENT   0.00  
REGION_RATING_CLIENT_W_CITY 0.00  
WEEKDAY_APPR_PROCESS_START 0.00  
HOUR_APPR_PROCESS_START 0.00  
REG_REGION_NOT_LIVE_REGION 0.00  
REG_REGION_NOT_WORK_REGION 0.00  
LIVE_REGION_NOT_WORK_REGION 0.00  
REG_CITY_NOT_LIVE_CITY 0.00  
REG_CITY_NOT_WORK_CITY 0.00  
LIVE_CITY_NOT_WORK_CITY 0.00  
ORGANIZATION_TYPE       0.00  
OBS_30_CNT_SOCIAL_CIRCLE 0.33  
DEF_30_CNT_SOCIAL_CIRCLE 0.33  
OBS_60_CNT_SOCIAL_CIRCLE 0.33  
DEF_60_CNT_SOCIAL_CIRCLE 0.33  
DAYS_LAST_PHONE_CHANGE 0.00  
FLAG_DOCUMENT_3          0.00  
AMT_REQ_CREDIT_BUREAU_HOUR 13.50  
AMT_REQ_CREDIT_BUREAU_DAY 13.50  
AMT_REQ_CREDIT_BUREAU_WEEK 13.50  
AMT_REQ_CREDIT_BUREAU_MON 13.50  
AMT_REQ_CREDIT_BUREAU_QRT 13.50  
AMT_REQ_CREDIT_BUREAU_YEAR 13.50  
AMT_INCOME_RANGE         0.08  
AMT_CREDIT_RANGE          0.00  
AGE                      0.00  
AGE_GROUP                0.00  
YEARS_EMPLOYED            0.00  
EMPLOYMENT_YEAR           27.08  
dtype: float64
```

```
In [58]: 1 app_data[ 'OCCUPATION_TYPE' ].describe()
```

```
Out[58]: count      211120
unique       18
top    Laborers
freq      55186
Name: OCCUPATION_TYPE, dtype: object
```

```
In [59]: 1 # Impute categorical variable 'OCCUPATION_TYPE' which has higher null p
2 # with a new category as assigning to any existing category
3
4 app_data[ 'OCCUPATION_TYPE' ] = app_data[ 'OCCUPATION_TYPE' ].cat.add_categories([Unknown])
5 app_data[ 'OCCUPATION_TYPE' ].fillna('Unknown', inplace =True)
```

Impute numerical variables with the median as there are no outliers that can be seen from results of describe() and mean() returns decimal values and these columns represent number of enquiries made which cannot be decimal:

```
In [60]: 1 app_data[ [ 'AMT_REQ_CREDIT_BUREAU_HOUR' , 'AMT_REQ_CREDIT_BUREAU_DAY' ,
2                               'AMT_REQ_CREDIT_BUREAU_WEEK' , 'AMT_REQ_CREDIT_BUREAU_MON'
3                               'AMT_REQ_CREDIT_BUREAU_QRT' , 'AMT_REQ_CREDIT_BUREAU_YEAR' ] ]
```

	AMT_REQ_CREDIT_BUREAU_HOUR	AMT_REQ_CREDIT_BUREAU_DAY	AMT_REQ_CREDIT_BUREAU_WEEK	AMT_REQ_CREDIT_BUREAU_MON	AMT_REQ_CREDIT_BUREAU_QRT	AMT_REQ_CREDIT_BUREAU_YEAR
count	265992.000000	265992.000000	265992.000000	265992.000000	265992.000000	265992.000000
mean	0.006402	0.007000	0.006402	0.007000	0.006402	0.007000
std	0.083849	0.110757	0.083849	0.110757	0.083849	0.110757
min	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	4.000000	9.000000	4.000000	9.000000	4.000000	9.000000

```
In [61]: 1 amount = [ 'AMT_REQ_CREDIT_BUREAU_HOUR' , 'AMT_REQ_CREDIT_BUREAU_DAY' , 'AMT_REQ_CREDIT_BUREAU_WEEK' , 'AMT_REQ_CREDIT_BUREAU_MON' , 'AMT_REQ_CREDIT_BUREAU_QRT' , 'AMT_REQ_CREDIT_BUREAU_YEAR' ]
```

```
In [62]: 1 for i in amount:
2     app_data[i].fillna(app_data[i].median(),inplace=True)
```

```
In [63]: 1 round(app_data.isnull().sum() / app_data.shape[0] *100,2)
```

```
Out[63]: SK_ID_CURR           0.00
TARGET              0.00
NAME_CONTRACT_TYPE 0.00
CODE_GENDER          0.00
FLAG_OWN_CAR         0.00
FLAG_OWN_REALTY     0.00
CNT_CHILDREN         0.00
AMT_INCOME_TOTAL    0.00
AMT_CREDIT            0.00
AMT_ANNUITY           0.00
AMT_GOODS_PRICE       0.09
NAME_TYPE_SUITE      0.00
NAME_INCOME_TYPE     0.00
NAME_EDUCATION_TYPE  0.00
NAME_FAMILY_STATUS    0.00
NAME_HOUSING_TYPE    0.00
REGION_POPULATION_RELATIVE 0.00
DAYS_BIRTH             0.00
DAYS_EMPLOYED          0.00
DAYS_REGISTRATION      0.00
DAYS_ID_PUBLISH        0.00
OCCUPATION_TYPE        0.00
CNT_FAM_MEMBERS        0.00
REGION_RATING_CLIENT   0.00
REGION_RATING_CLIENT_W_CITY 0.00
WEEKDAY_APPR_PROCESS_START 0.00
HOUR_APPR_PROCESS_START 0.00
REG_REGION_NOT_LIVE_REGION 0.00
REG_REGION_NOT_WORK_REGION 0.00
LIVE_REGION_NOT_WORK_REGION 0.00
REG_CITY_NOT_LIVE_CITY 0.00
REG_CITY_NOT_WORK_CITY 0.00
LIVE_CITY_NOT_WORK_CITY 0.00
ORGANIZATION_TYPE       0.00
OBS_30_CNT_SOCIAL_CIRCLE 0.33
DEF_30_CNT_SOCIAL_CIRCLE 0.33
OBS_60_CNT_SOCIAL_CIRCLE 0.33
DEF_60_CNT_SOCIAL_CIRCLE 0.33
DAYS_LAST_PHONE_CHANGE 0.00
FLAG_DOCUMENT_3          0.00
AMT_REQ_CREDIT_BUREAU_HOUR 0.00
AMT_REQ_CREDIT_BUREAU_DAY 0.00
AMT_REQ_CREDIT_BUREAU_WEEK 0.00
AMT_REQ_CREDIT_BUREAU_MON 0.00
AMT_REQ_CREDIT_BUREAU_QRT 0.00
AMT_REQ_CREDIT_BUREAU_YEAR 0.00
AMT_INCOME_RANGE         0.08
AMT_CREDIT_RANGE          0.00
AGE                     0.00
AGE_GROUP               0.00
YEARS_EMPLOYED           0.00
EMPLOYMENT_YEAR          27.08
dtype: float64
```

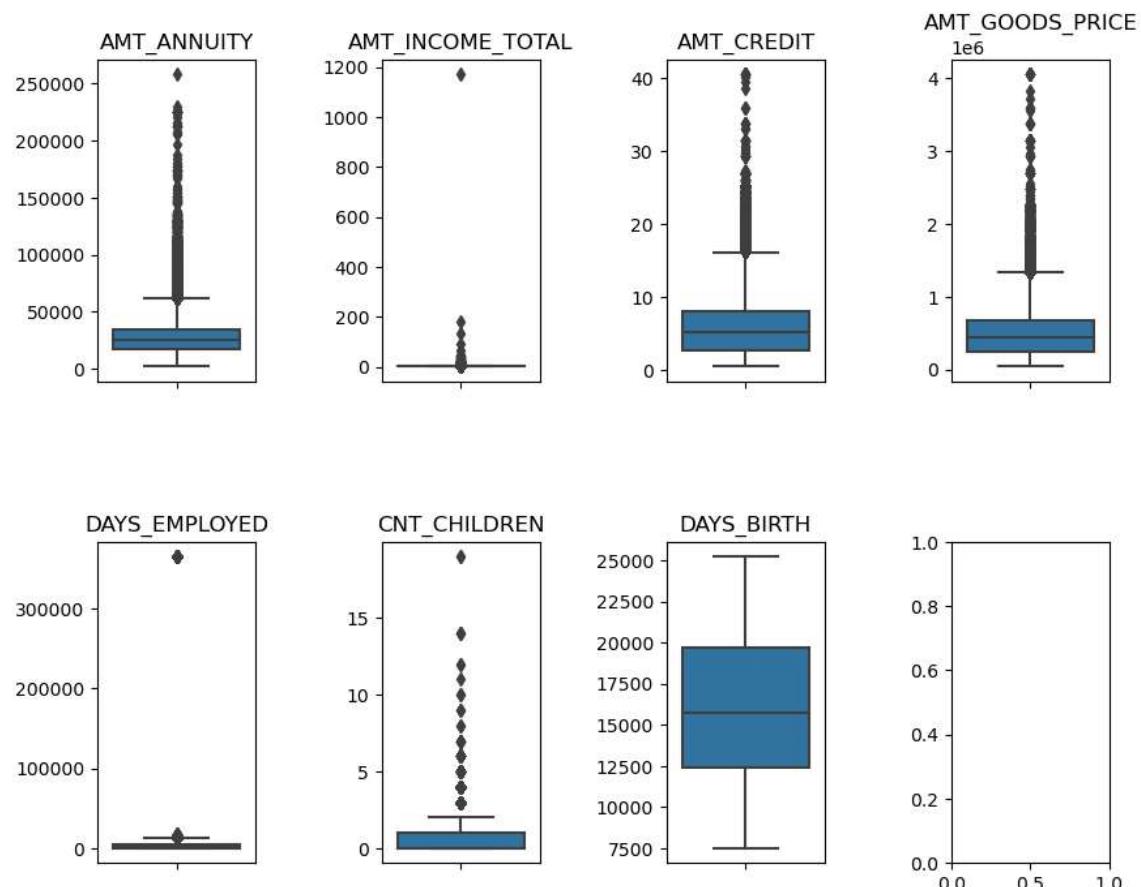
We still have few null values in the columns: AMT_GOODS_PRICE, OBS_30_CNT_SOCIAL_CIRCLE, DEF_30_CNT_SOCIAL_CIRCLE, OBS_60_CNT_SOCIAL_CIRCLE, DEF_60_CNT_SOCIAL_CIRCLE. We can ignore as this

percentage is very less.

```
In [64]: 1 # identifying outliers
```

```
In [65]: 1 app_outlier_col_1 = ['AMT_ANNUITY', 'AMT_INCOME_TOTAL', 'AMT_CREDIT', 'AMT_GOODS_PRICE']
2 app_outlier_col_2 = ['CNT_CHILDREN', 'DAYS_BIRTH']
```

```
In [77]: 1 fig, axes = plt.subplots(2, 4 , figsize=(10,8))
2 for i in enumerate(app_outlier_col_1):
3     plt.subplot(2,4,i[0]+1)
4     plt.subplots_adjust(wspace=0.8, hspace=0.5)
5     sns.boxplot(y=app_data[i[1]])
6     plt.title(i[1])
7     plt.ylabel('')
8
9 for i in enumerate(app_outlier_col_2):
10    plt.subplot(2,4,i[0]+6)
11    sns.boxplot(y=app_data[i[1]])
12    plt.title(i[1])
13    plt.ylabel("")
```



AMT_ANNUITY, AMT_CREDIT, AMT_GOODS_PRICE, CNT_CHILDREN have some number of outliers. AMT_INCOME_TOTAL has huge number of outliers which indicate that few of the loan applicants have high income when compared to the others. DAYS_BIRTH has no outliers which means the data available is reliable. DAYS_EMPLOYED has outlier values around 350000(days) which is around 958 years which is impossible and hence this has to be incorrect entry.

In [79]:

```
1 app_data[['AMT_ANNUITY', 'AMT_INCOME_TOTAL', 'AMT_CREDIT', 'AMT_GOODS_PRICE',
2 'DAYS_BIRTH', 'CNT_CHILDREN', 'DAYS_EMPLOYED']].describe()
```

Out[79]:

	AMT_ANNUITY	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_GOODS_PRICE	DAYS_BIR
count	307499.000000	307511.000000	307511.000000	3.072330e+05	307511.000000
mean	27108.573909	1.687979	5.990260	5.383962e+05	16036.995000
std	14493.737315	2.371231	4.024908	3.694465e+05	4363.988600
min	1615.500000	0.256500	0.450000	4.050000e+04	7489.000000
25%	16524.000000	1.125000	2.700000	2.385000e+05	12413.000000
50%	24903.000000	1.471500	5.135310	4.500000e+05	15750.000000
75%	34596.000000	2.025000	8.086500	6.795000e+05	19682.000000
max	258025.500000	1170.000000	40.500000	4.050000e+06	25229.000000

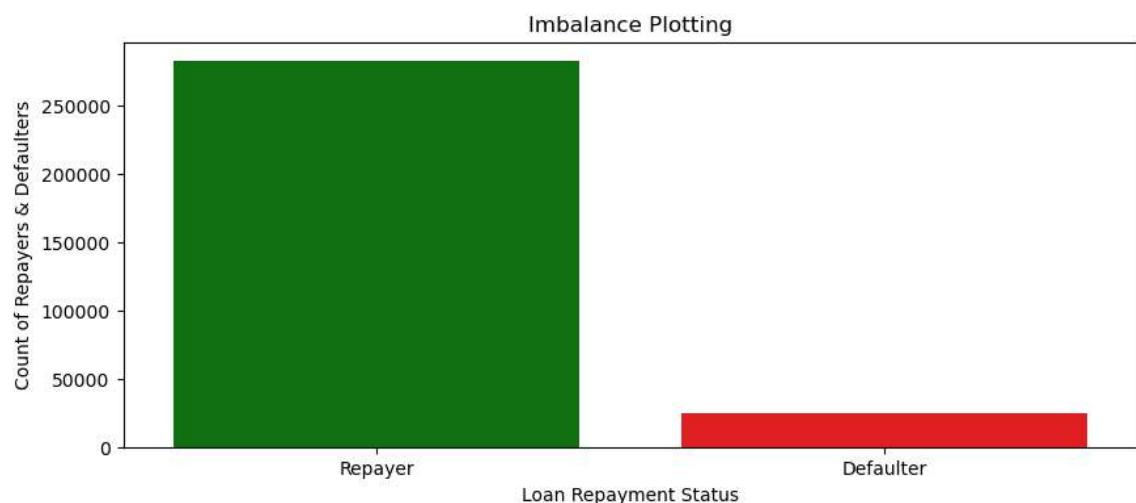
Strategy: The data analysis flow has been planned in following way : Imbalance in Data
Categorical Data Analysis Categorical segmented Univariate Analysis Categorical
Bi/Multivariate analysis

Numeric Data Analysis Bi-furcation of databased based on TARGET data Correlation Matrix
Numerical segmented Univariate Analysis Numerical Bi/Multivariate analysis

Imbalance Analysis

In [82]:

```
1 Imbalance = app_data["TARGET"].value_counts().reset_index()
2
3 plt.figure(figsize=(10,4))
4 x= ['Repayer', 'Defaulter']
5 sns.barplot(x=x, "TARGET", data = Imbalance, palette= ['g', 'r'])
6 plt.xlabel("Loan Repayment Status")
7 plt.ylabel("Count of Repayers & Defaulters")
8 plt.title("Imbalance Plotting")
9 plt.show()
```



In [86]:

```
1 count_0 = Imbalance.iloc[0]["TARGET"]
2 count_1 = Imbalance.iloc[1]["TARGET"]
3 count_0_perc = round(count_0 / (count_0 + count_1) * 100, 2)
4 count_1_perc = round(count_1 / (count_0 + count_1) * 100, 2)
5
6 print('Ratios of imbalance in percentage with respect to Repayer and De
7 print('Ratios of imbalance in relative with respect to Repayer and Defa
8
```

Ratios of imbalance in percentage with respect to Repayer and Defaulter data are: 91.93 and 8.07
Ratios of imbalance in relative with respect to Repayer and Defaulter data is 11.39 :1 (approx)

Plotting Functions

In [87]:

```

1 # function for plotting repetitive countplots in univariate categorical
2 # This function will create two subplots:
3 # 1. Count plot of categorical column w.r.t TARGET;
4 # 2. Percentage of defaulters within column
5
6 def univariate_categorical(feature,ylog=False,label_rotation=False,hori-
7   temp = app_data[feature].value_counts()
8   df1 = pd.DataFrame({feature: temp.index,'Number of contracts': temp
9
10  # Calculate the percentage of target=1 per category value
11  cat_perc = app_data[[feature, 'TARGET']].groupby([feature],as_index=False)
12  cat_perc['TARGET'] = cat_perc['TARGET']*100
13  cat_perc.sort_values(by='TARGET', ascending=False, inplace=True)
14
15  if(horizontal_layout):
16    fig, (ax1, ax2) = plt.subplots(ncols=2, figsize=(12,6))
17  else:
18    fig, (ax1, ax2) = plt.subplots(nrows=2, figsize=(20,24))
19
20  # 1. Subplot 1: Count plot of categorical column
21  # sns.set_palette("Set2")
22  s = sns.countplot(ax=ax1,
23    x = feature,
24    data=app_data,
25    hue ="TARGET",
26    order=cat_perc[feature],
27    palette=['g','r'])
28
29  # Define common styling
30  ax1.set_title(feature, fontdict={'fontsize' : 10, 'fontweight' : 3, 'fontstyle': 'italic'})
31  ax1.legend(['Repayer','Defaulter'])
32
33  # If the plot is not readable, use the log scale.
34  if ylog:
35    ax1.set_yscale('log')
36    ax1.set_ylabel("Count (log)",fontdict={'fontsize' : 10, 'fontweight' : 3, 'fontstyle': 'italic'})
37
38
39  if(label_rotation):
40    s.set_xticklabels(s.get_xticklabels(),rotation=90)
41
42  # 2. Subplot 2: Percentage of defaulters within the categorical col-
43  s = sns.barplot(ax=ax2,
44    x = feature,
45    y='TARGET',
46    order=cat_perc[feature],
47    data=cat_perc,
48    palette='Set2')
49
50  if(label_rotation):
51    s.set_xticklabels(s.get_xticklabels(),rotation=90)
52    plt.ylabel('Percent of Defaulters [%]', fontsize=10)
53    plt.tick_params(axis='both', which='major', labelsize=10)
54    ax2.set_title(feature + " Defaulter %", fontdict={'fontsize' : 15, 'fontstyle': 'italic'})
55
56  plt.show();

```

```
In [88]: 1 # function for plotting repetitive countplots in bivariate categorical  
2  
3 def bivariate_bar(x,y,df,hue,figsize):  
4  
5     plt.figure(figsize=figsize)  
6     sns.barplot(x=x,  
7                 y=y,  
8                 data=df,  
9                 hue=hue,  
10                palette =['g','r'])  
11  
12     # Defining aesthetics of Labels and Title of the plot using style a  
13     plt.xlabel(x,fontdict={'fontsize' : 10, 'fontweight' : 3, 'color' :  
14     plt.ylabel(y,fontdict={'fontsize' : 10, 'fontweight' : 3, 'color' :  
15     plt.title(col, fontdict={'fontsize' : 15, 'fontweight' : 5, 'color':  
16     plt.xticks(rotation=90, ha='right')  
17     plt.legend(labels = ['Repayer','Defaulter'])  
18     plt.show()
```

```
In [89]: 1 # function for plotting repetitive rel plots in bivariate numerical and categorical  
2  
3 def bivariate_rel(x,y,data, hue, kind, palette, legend,figsize):  
4  
5     plt.figure(figsize=figsize)  
6     sns.relplot(x=x,  
7                 y=y,  
8                 data=app_data,  
9                 hue="TARGET",  
10                kind=kind,  
11                palette = ['g','r'],  
12                legend = False)  
13     plt.legend(['Repayer','Defaulter'])  
14     plt.xticks(rotation=90, ha='right')  
15     plt.show()
```

```
In [90]: 1 #function for plotting repetitive countplots in univariate categorical
2
3 def univariate_merged(col,df,hue,palette,ylog,figsize):
4     plt.figure(figsize=figsize)
5     ax=sns.countplot(x=col,
6                         data=df,
7                         hue=hue,
8                         palette=palette,
9                         order=df[col].value_counts().index)
10
11
12 if ylog:
13     plt.yscale('log')
14     plt.ylabel("Count (log)",fontdict={'fontsize' : 10, 'fontweight': 'bold'})
15 else:
16     plt.ylabel("Count",fontdict={'fontsize' : 10, 'fontweight' : 3,
17
18     plt.title(col , fontdict={'fontsize' : 15, 'fontweight' : 5, 'color': 'red'})
19     plt.legend(loc = "upper right")
20     plt.xticks(rotation=90, ha='right')
21
22 plt.show()
```

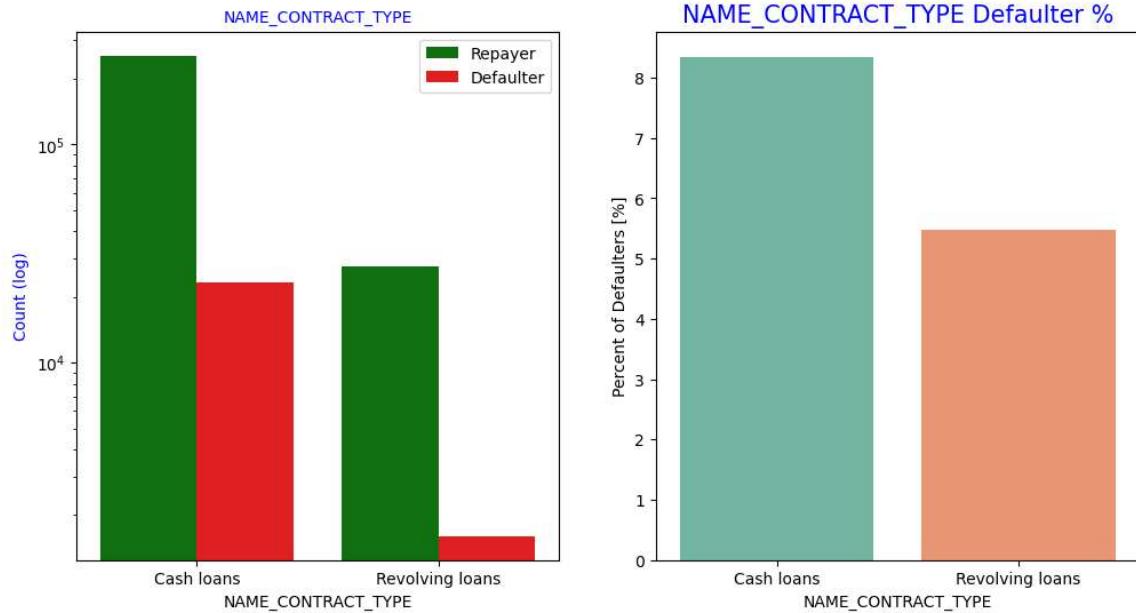
```
In [91]: 1 # Function to plot point plots on merged dataframe
2
3 def merged_pointplot(x,y):
4     plt.figure(figsize=(8,4))
5     sns.pointplot(x=x,
6                     y=y,
7                     hue="TARGET",
8                     data=loan_process_df,
9                     palette=['g','r'])
10    # plt.legend(['Repayer', 'Defaulter'])
```

Categorical Variables Analysis

```
In [92]: 1 # Segmented Univariate Analysis
```

In [93]:

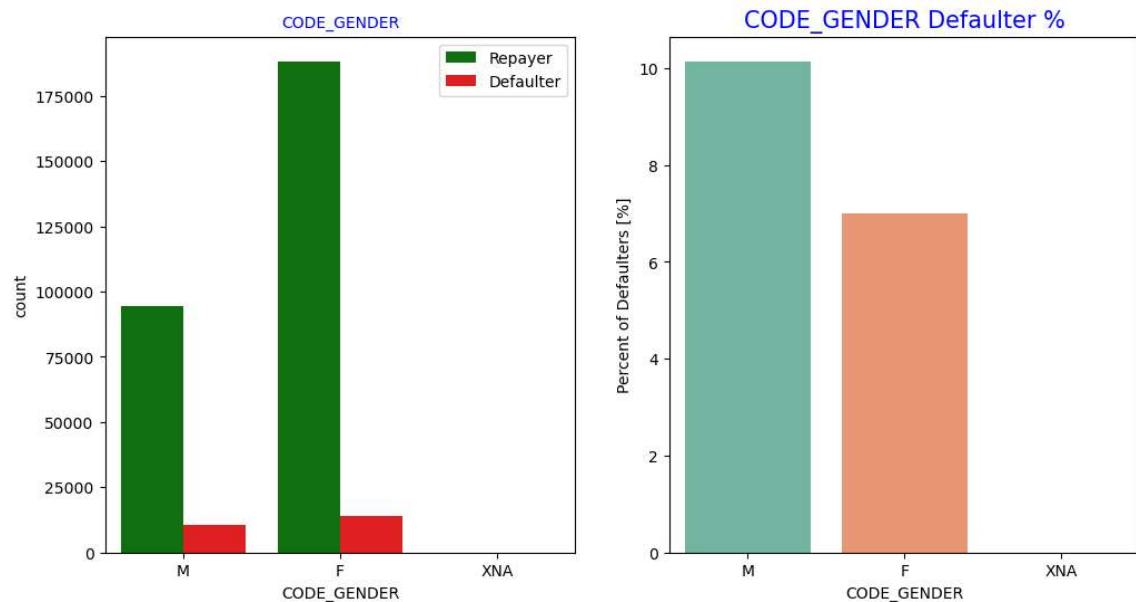
```
1 # Checking the contract type based on Loan repayment status
2 univariate_categorical('NAME_CONTRACT_TYPE',True)
```



Contract type: Revolving loans are just a small fraction (10%) from the total number of loans; in the same time, a larger amount of Revolving loans, comparing with their frequency, are not repaid.

In [94]:

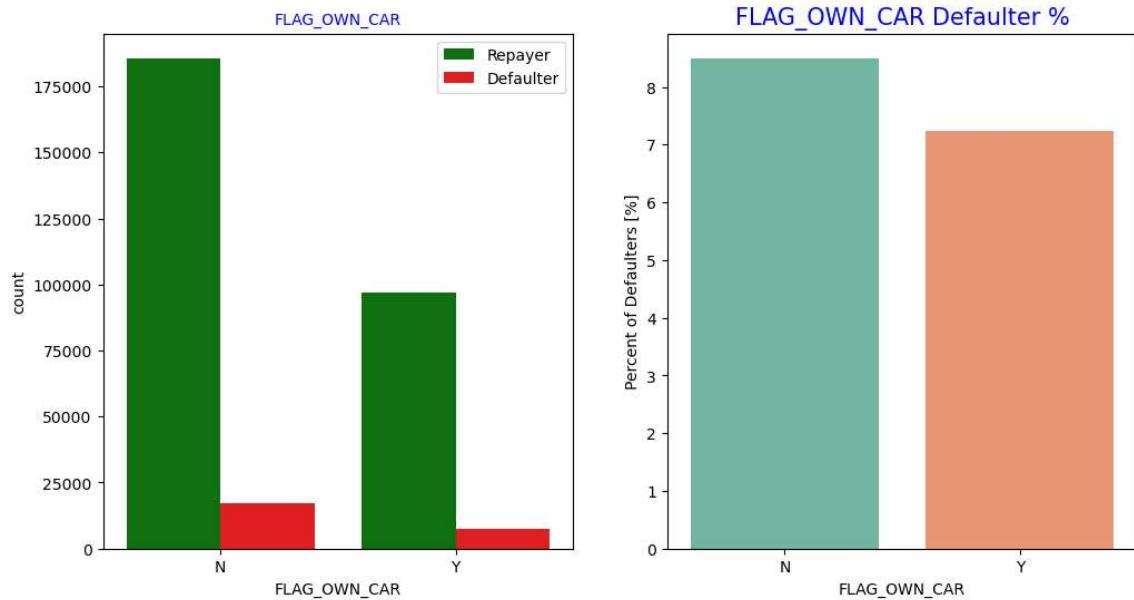
```
1 # Checking the type of Gender on Loan repayment status
2 univariate_categorical('CODE_GENDER')
```



Contract type: Revolving loans are just a small fraction (10%) from the total number of loans; in the same time, a larger amount of Revolving loans, comparing with their frequency, are not repaid.

In [95]:

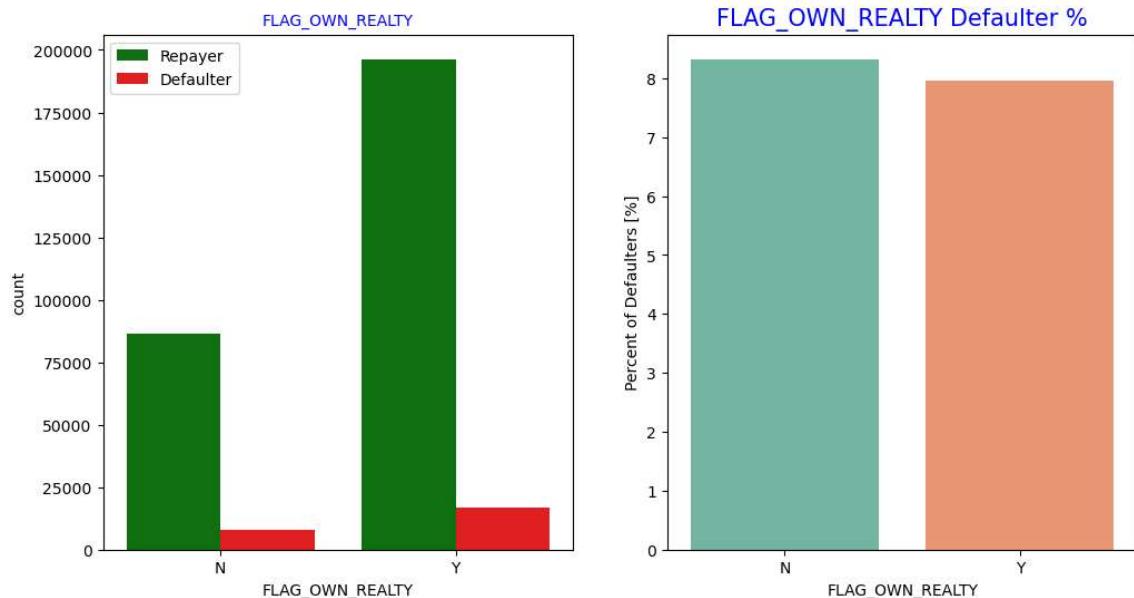
```
1 # Checking if owning a car is related to loan repayment status
2 univariate_categorical('FLAG_own_car')
```



Clients who own a car are half in number of the clients who don't own a car. But based on the percentage of default, there is no correlation between owning a car and loan repayment as in both cases the default percentage is almost same.

In [96]:

```
1 # Checking if owning a realty is related to loan repayment status
2 univariate_categorical('FLAG_own_realty')
```



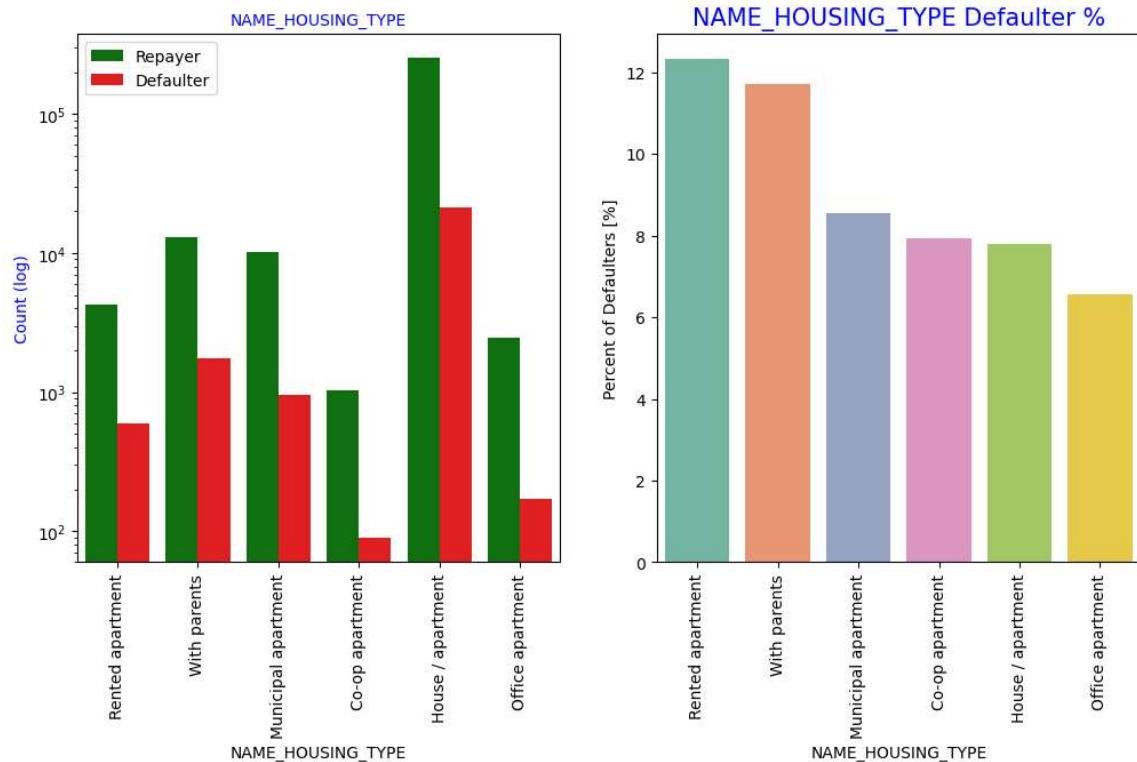
The clients who own real estate are more than double of the ones that don't own. But the defaulting rate of both categories are around the same (~8%). Thus there is no correlation between owning a realty and defaulting the loan.

In [97]:

```

1 # Analyzing Housing Type based on Loan repayment status
2 univariate_categorical("NAME_HOUSING_TYPE",True,True,True)

```



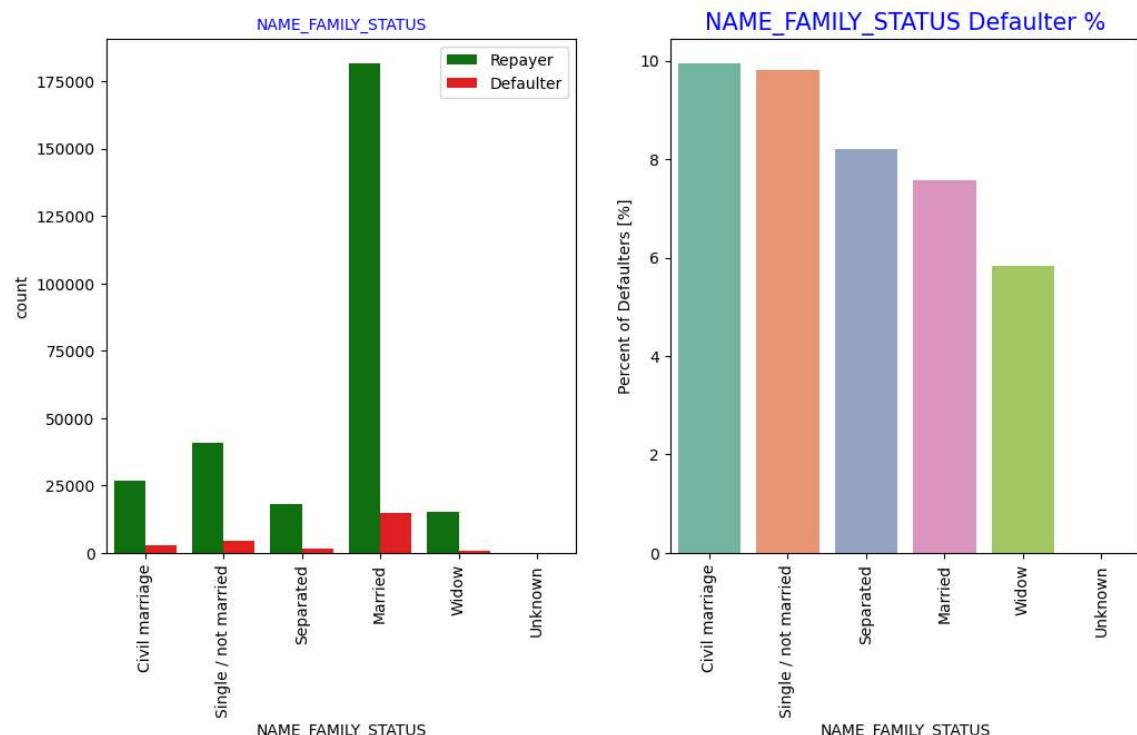
Majority of people live in House/apartment People living in office apartments have lowest default rate People living with parents (~11.5%) and living in rented apartments(>12%) have higher probability of defaulting

In [98]:

```

1 # Analyzing Family status based on Loan repayment status
2 univariate_categorical("NAME_FAMILY_STATUS",False,True,True)

```



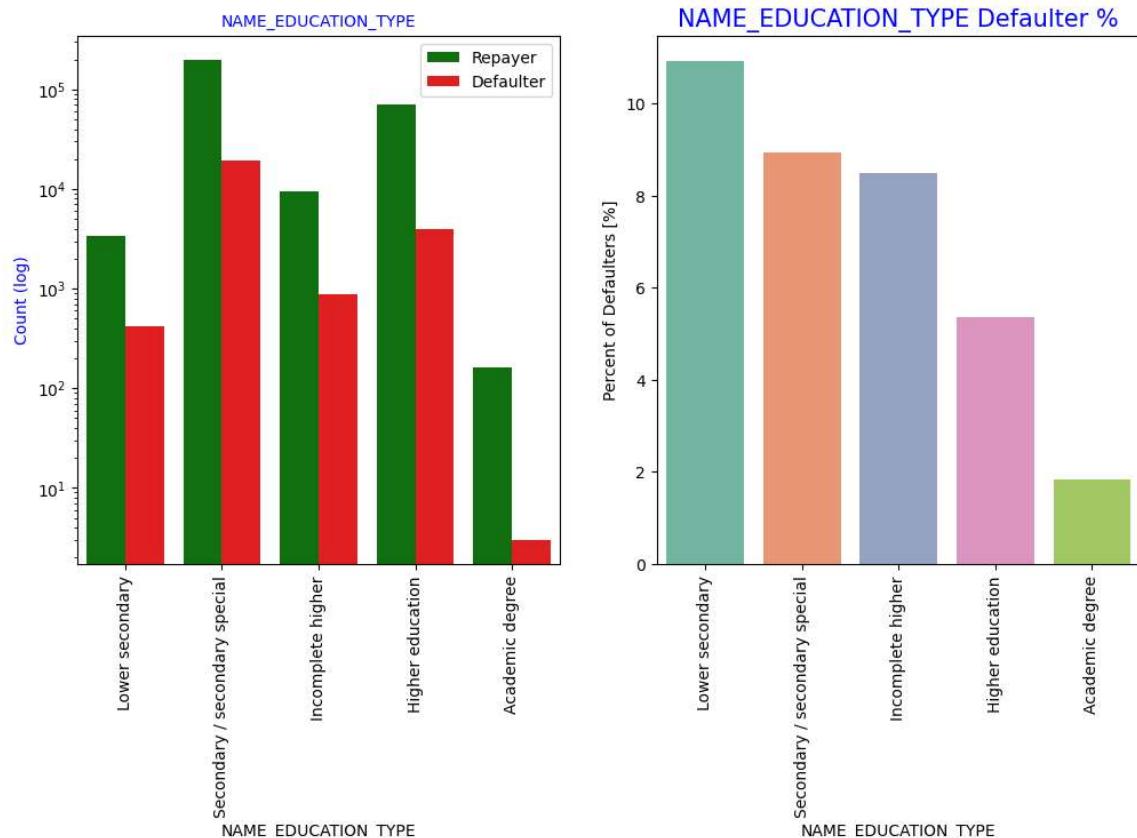
Most of the people who have taken loan are married, followed by Single/not married and civil marriage In terms of percentage of not repayment of loan, Civil marriage has the highest percent of not repayment (10%), with Widow the lowest (exception being Unknown).

In [99]:

```

1 # Analyzing Education Type based on Loan repayment status
2 univariate_categorical("NAME_EDUCATION_TYPE", True, True, True)

```



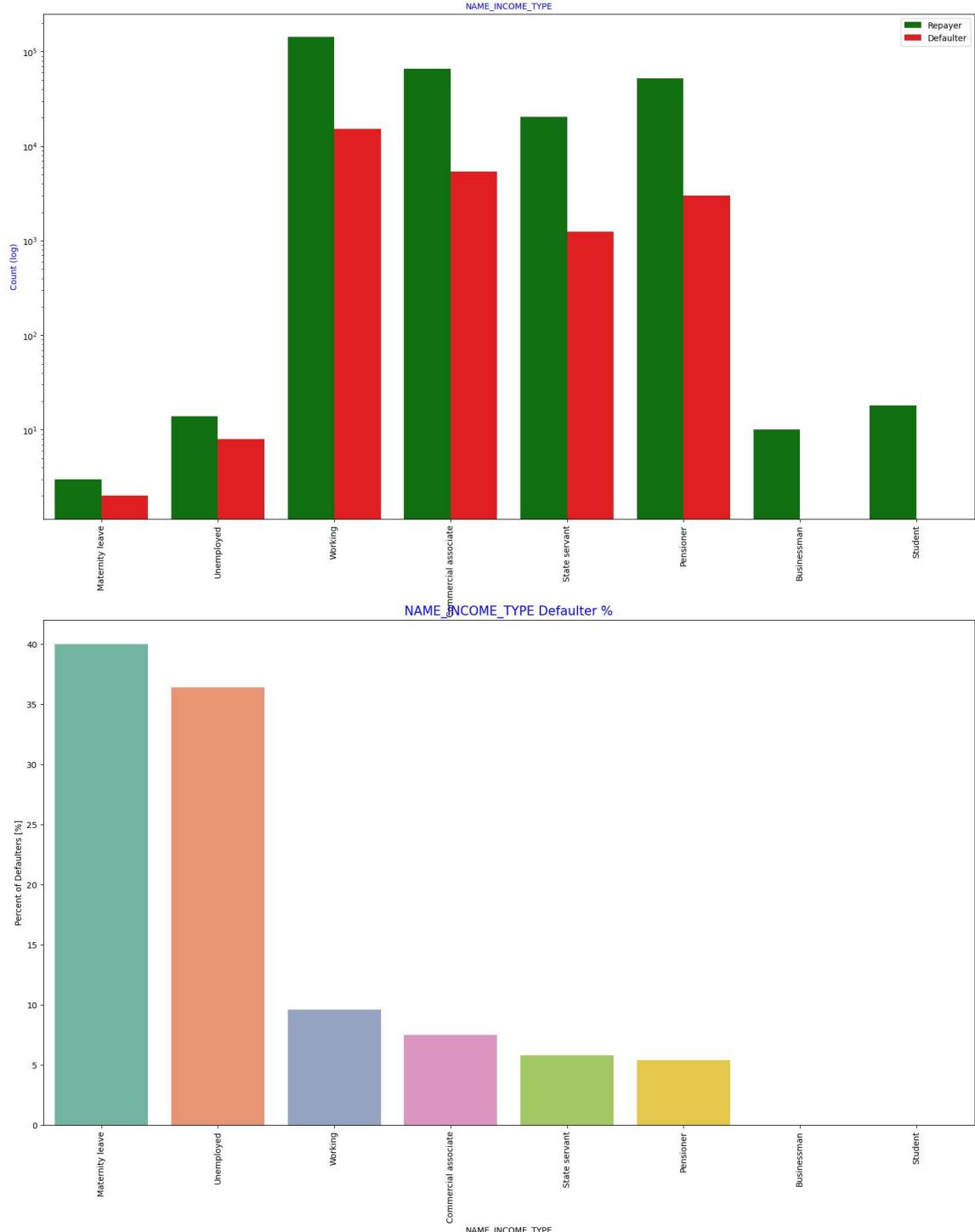
Majority of the clients have Secondary / secondary special education, followed by clients with Higher education. Only a very small number having an academic degree. The Lower secondary category, although rare, have the largest rate of not returning the loan (11%). The people with Academic degree have less than 2% defaulting rate.

In [100]:

```

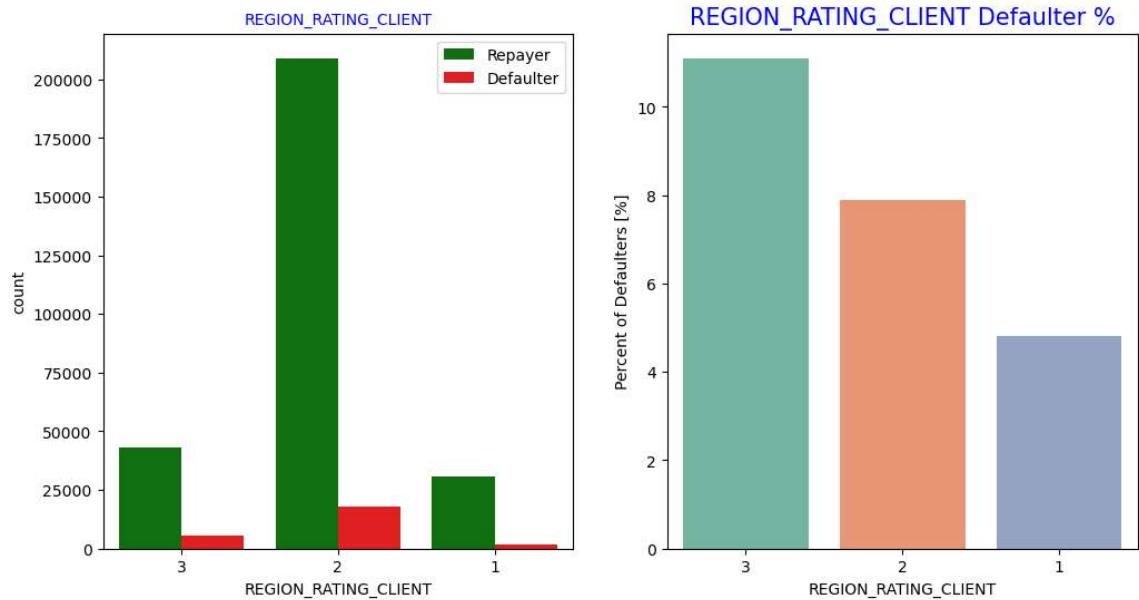
1 # Analyzing Income Type based on Loan repayment status
2 univariate_categorical("NAME_INCOME_TYPE",True,True,False)

```



Most of applicants for loans have income type as Working, followed by Commercial associate, Pensioner and State servant. The applicants with the type of income Maternity leave have almost 40% ratio of not returning loans, followed by Unemployed (37%). The rest of types of incomes are under the average of 10% for not returning loans. Student and Businessmen, though less in numbers do not have any default record. Thus these two category are safest for providing loan.

```
In [101]: 1 # Analyzing Region rating where applicant Lives based on Loan repayment
2 univariate_categorical("REGION_RATING_CLIENT",False,False,True)
```



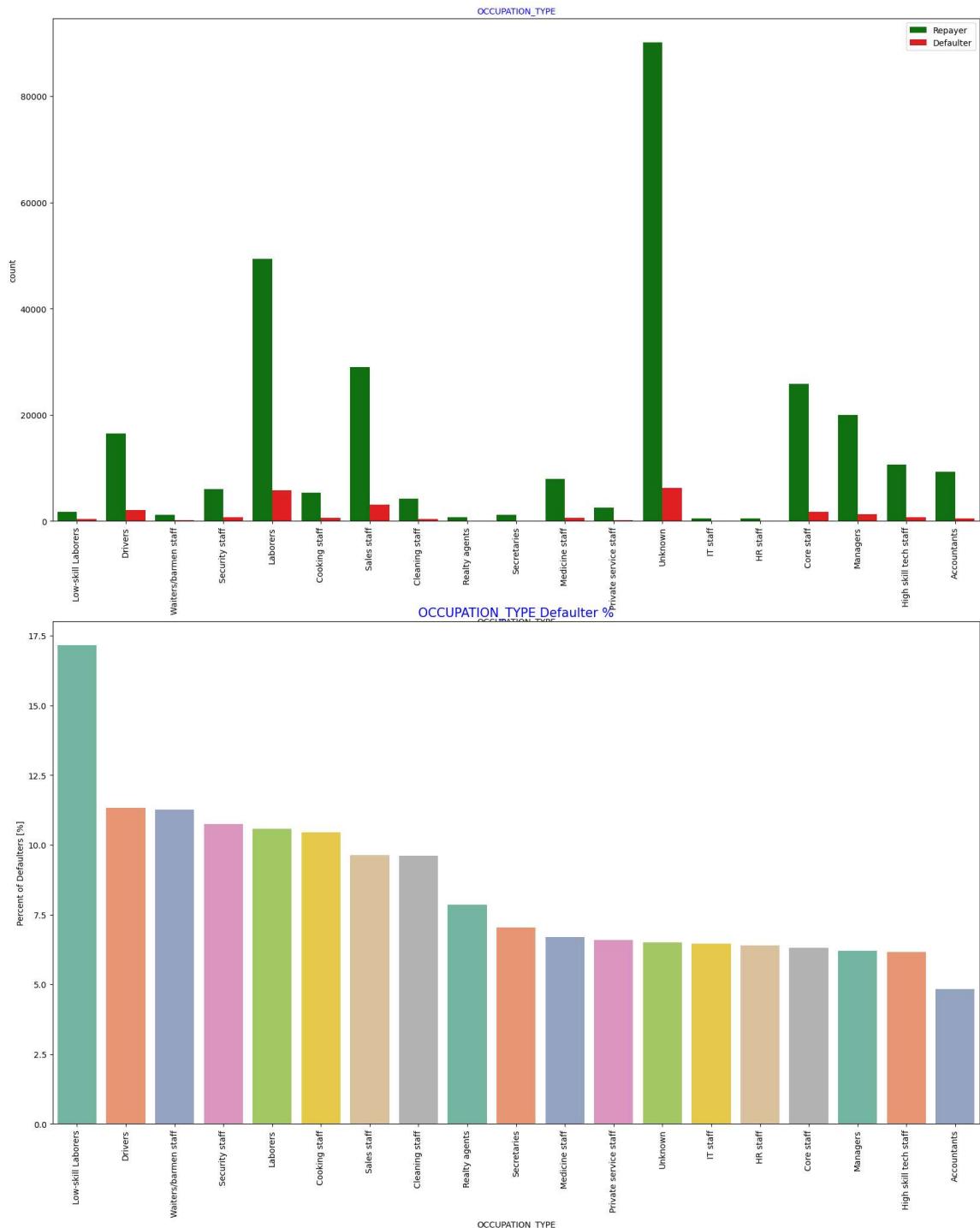
Most of the applicants are living in Region_Rating 2 place. Region Rating 3 has the highest default rate (11%) Applicant living in Region_Rating 1 has the lowest probability of defaulting, thus safer for approving loans

In [102]:

```

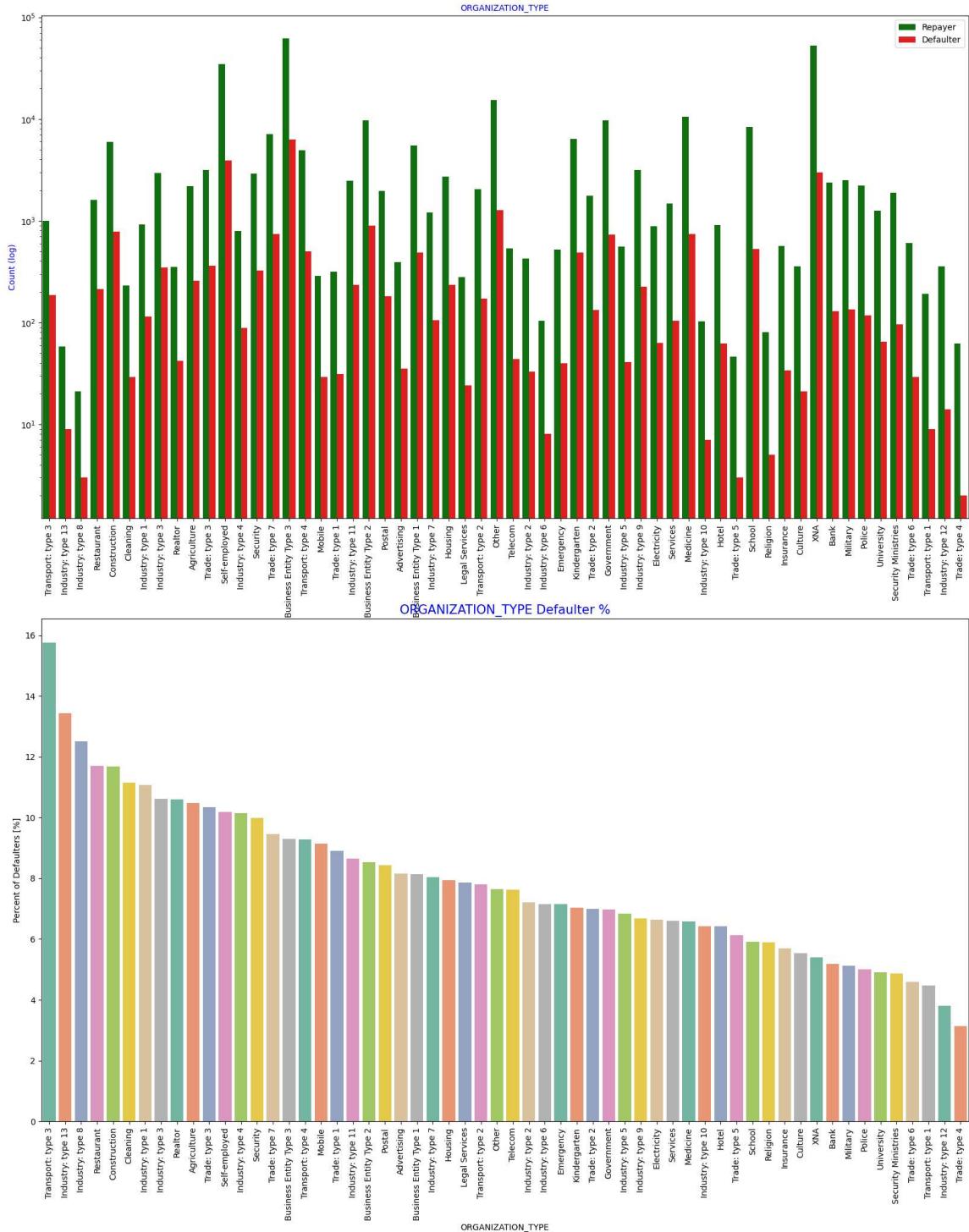
1 # Analyzing Occupation Type where applicant lives based on Loan repayment
2 univariate_categorical("OCCUPATION_TYPE", False, True, False)

```



Most of the loans are taken by Laborers, followed by Sales staff. IT staff take the lowest amount of loans. The category with highest percent of not repaid loans are Low-skill Laborers (above 17%), followed by Drivers and Waiters/barmen staff, Security staff, Laborers and Cooking staff.

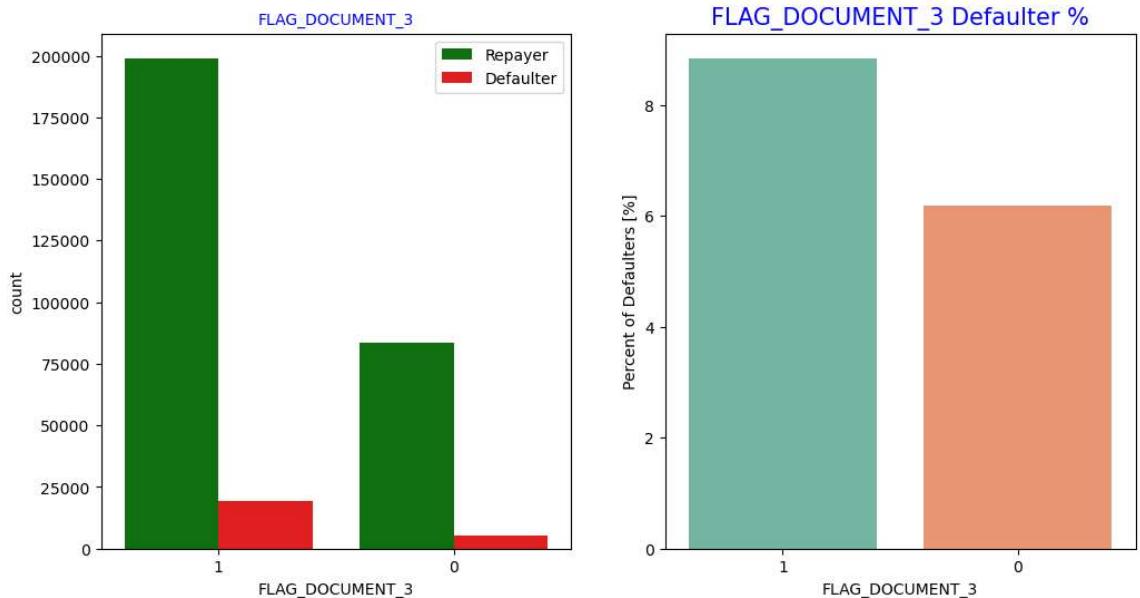
```
In [103]: 1 # Checking Loan repayment status based on Organization type  
2 univariate_categorical("ORGANIZATION_TYPE",True,True,False)
```



Organizations with highest percent of loans not repaid are Transport: type 3 (16%), Industry: type 13 (13.5%), Industry: type 8 (12.5%) and Restaurant (less than 12%). Self employed people have relative high defaulting rate, and thus should be avoided to be approved for loan or provide loan with higher interest rate to mitigate the risk of defaulting. Most of the people application for loan are from Business Entity Type 3 For a very high number of applications, Organization type information is unavailable(XNA) It can be seen that following category of organization type has lesser defaulters thus safer for providing loans: Trade Type 4 and 5 Industry type 8

In [105]:

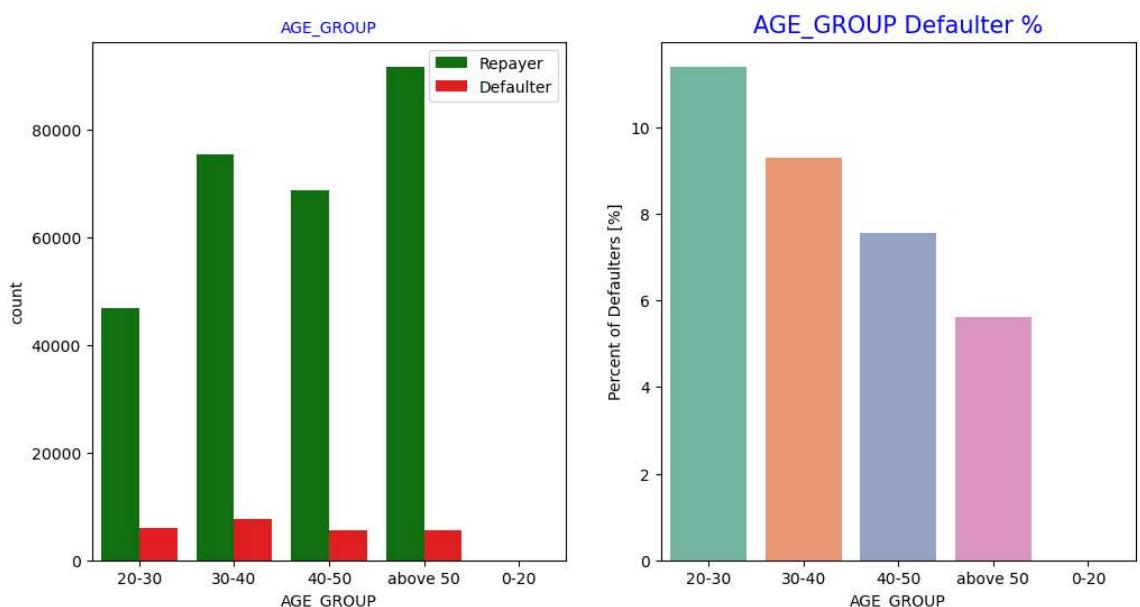
```
1 # Analyzing Flag_Doc_3 submission status based on Loan repayment status
2 univariate_categorical("FLAG_DOCUMENT_3",False,False,True)
```



There is no significant correlation between repayers and defaulters in terms of submitting document 3 as we see even if applicants have submitted the document, they have defaulted a slightly more (~9%) than who have not submitted the document (6%)

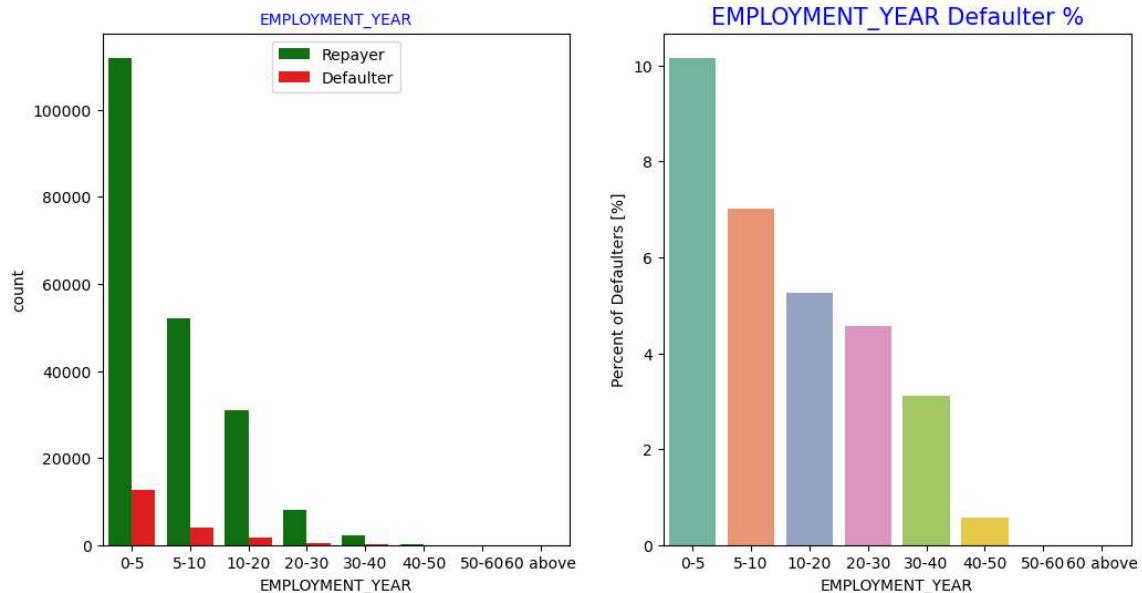
In [106]:

```
1 # Analyzing Age Group based on Loan repayment status
2 univariate_categorical("AGE_GROUP",False,False,True)
3
```



People in the age group range 20-40 have higher probability of defaulting People above age of 50 have low probability of defaulting

```
In [107]: 1 # Analyzing Employment_Year based on Loan repayment status
2 univariate_categorical("EMPLOYMENT_YEAR",False,False,True)
```



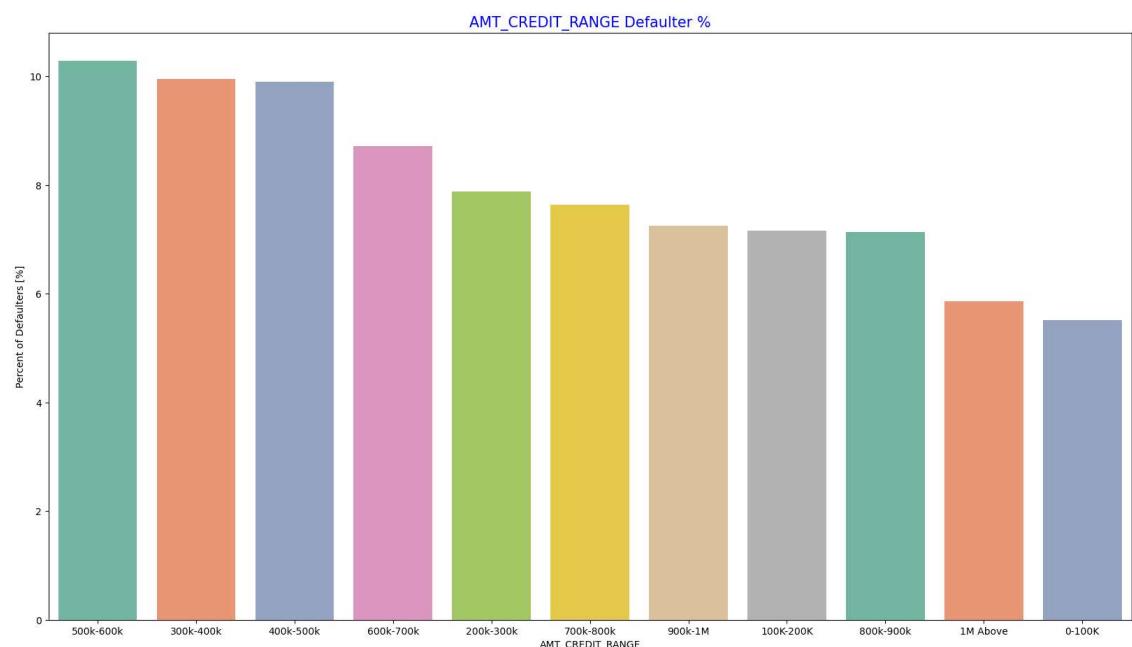
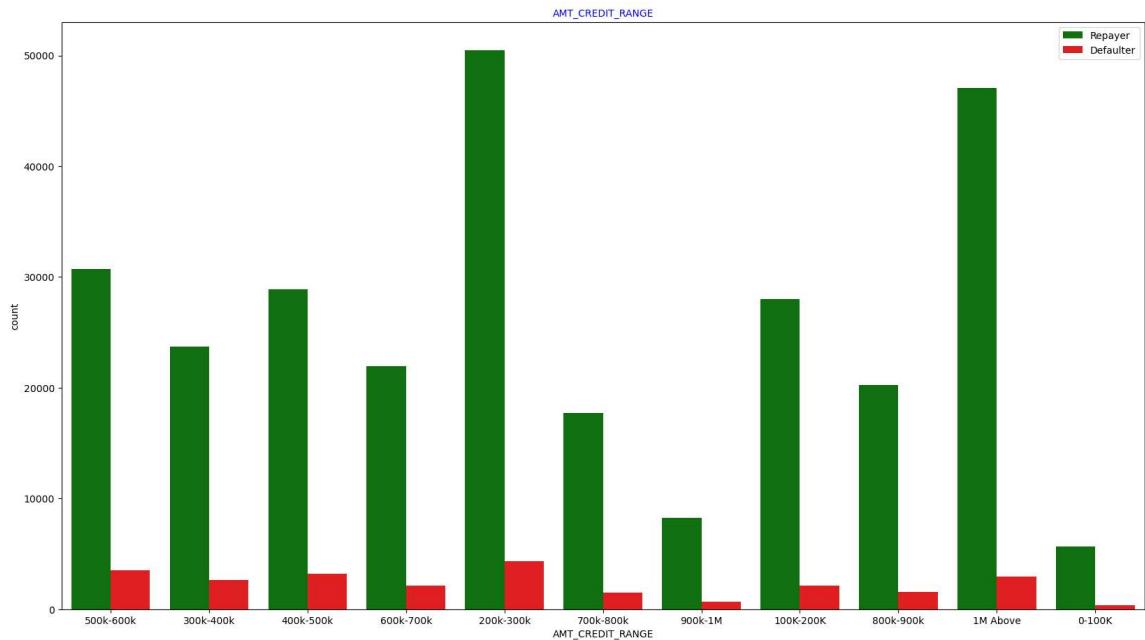
Majority of the applicants have been employed in between 0-5 years. The defaulting rating of this group is also the highest which is 10% With increase of employment year, defaulting rate is gradually decreasing with people having 40+ year experience having less than 1% default rate

In [108]:

```

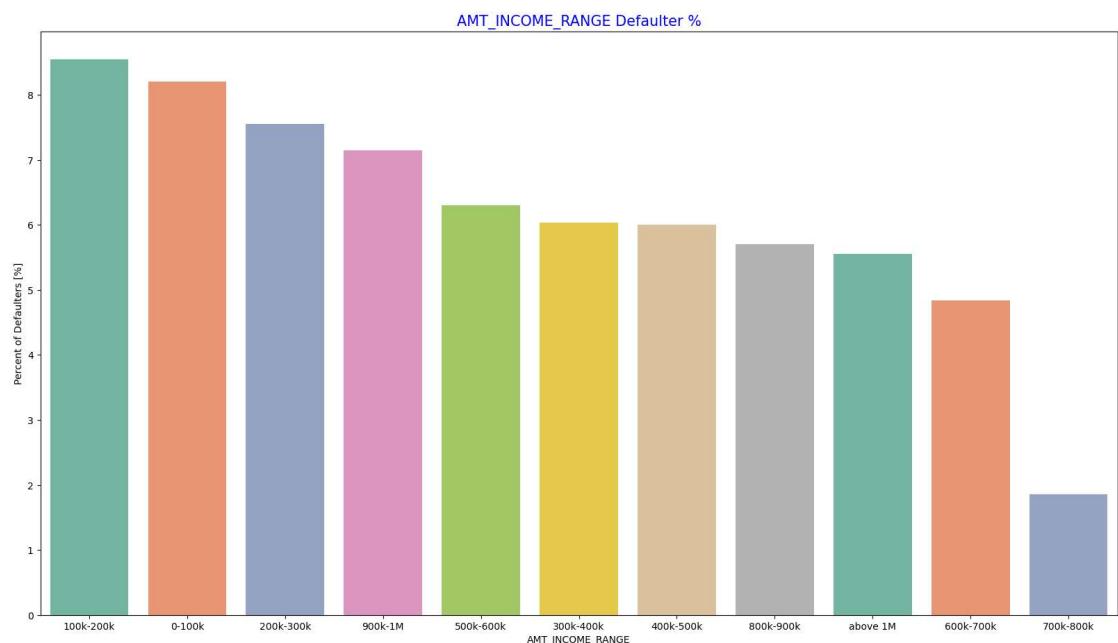
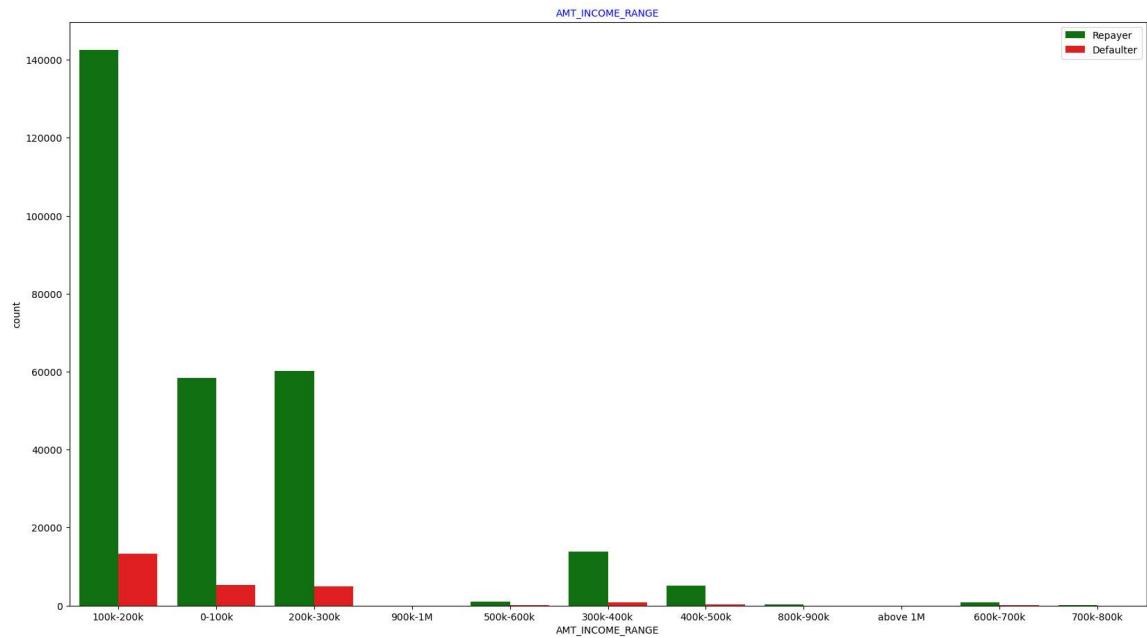
1 # Analyzing Amount_Credit based on loan repayment status
2 univariate_categorical("AMT_CREDIT_RANGE",False,False,False)

```



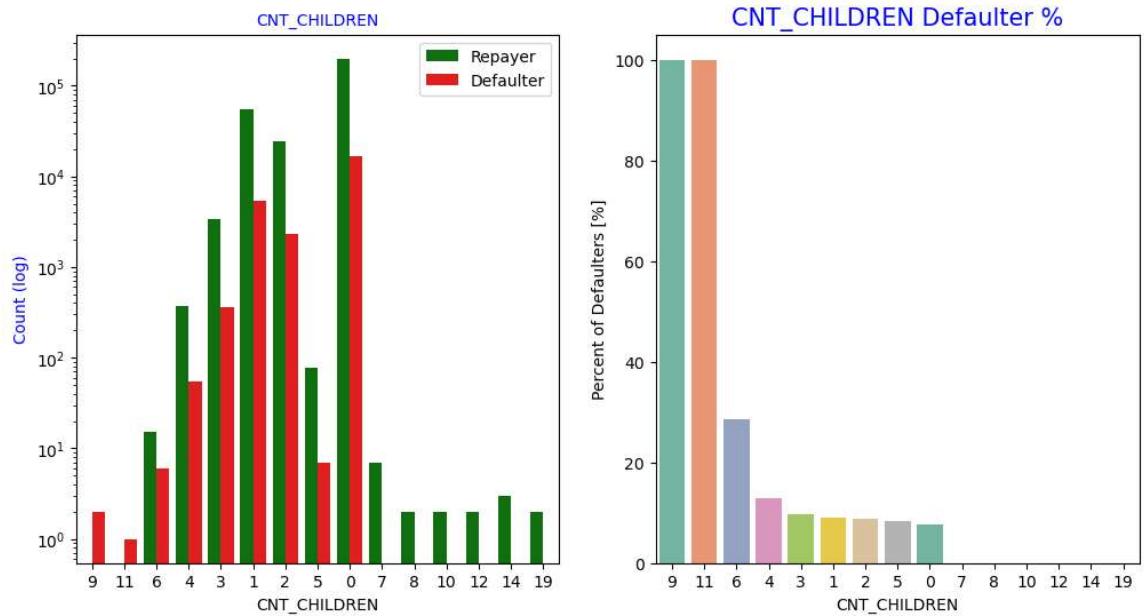
More than 80% of the loan provided are for amount less than 900,000 People who get loan for 300-600k tend to default more than others.

```
In [109]: 1 # Analyzing Amount_Income Range based on Loan repayment status
2 univariate_categorical("AMT_INCOME_RANGE",False,False,False)
```



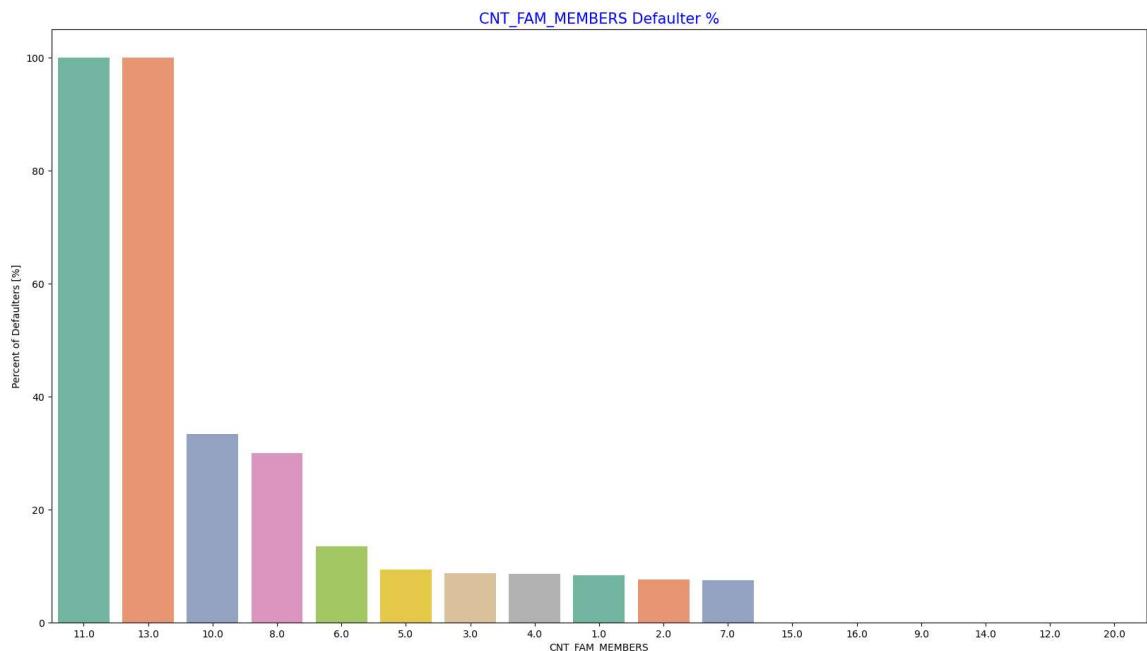
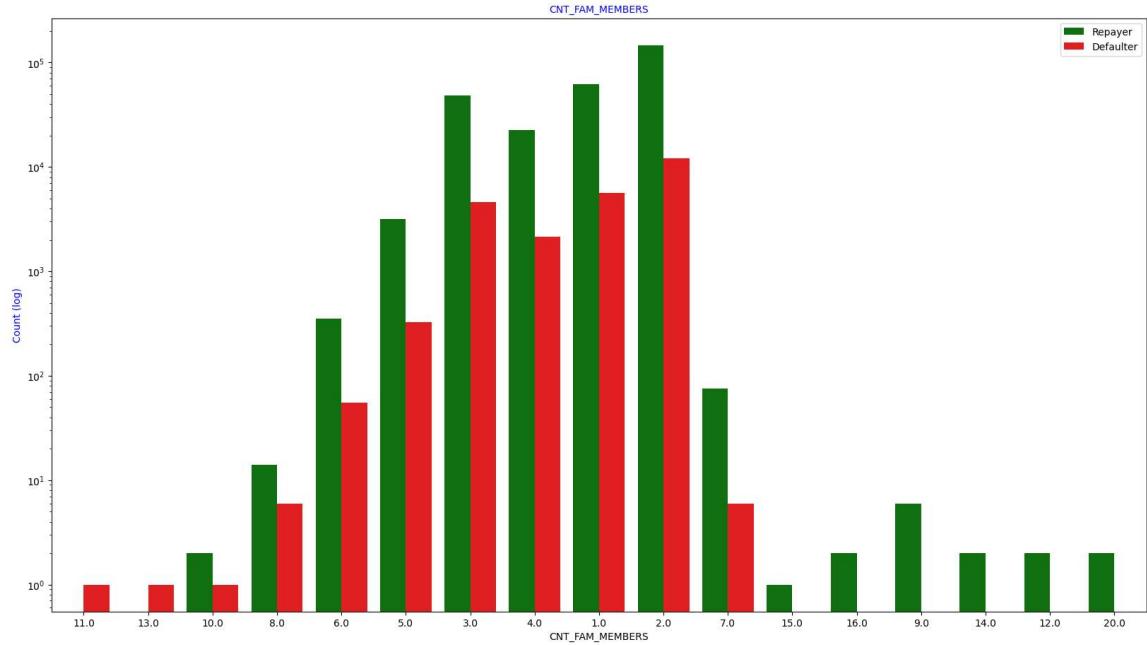
90% of the applications have Income total less than 300,000 Application with Income less than 300,000 has high probability of defaulting Applicant with Income more than 700,000 are less likely to default

```
In [110]: 1 # Analyzing Number of children based on Loan repayment status
2 univariate_categorical("CNT_CHILDREN",True)
```



Most of the applicants do not have children. Very few clients have more than 3 children.
 Client who have more than 4 children has a very high default rate with child count 9 and 11 showing 100% default rate.

```
In [111]: 1 # Analyzing Number of family members based on Loan repayment status
2 univariate_categorical("CNT_FAM_MEMBERS",True, False, False)
```



Family member follows the same trend as children where having more family members increases the risk of defaulting

Categorical Bi/Multivariate Analysis

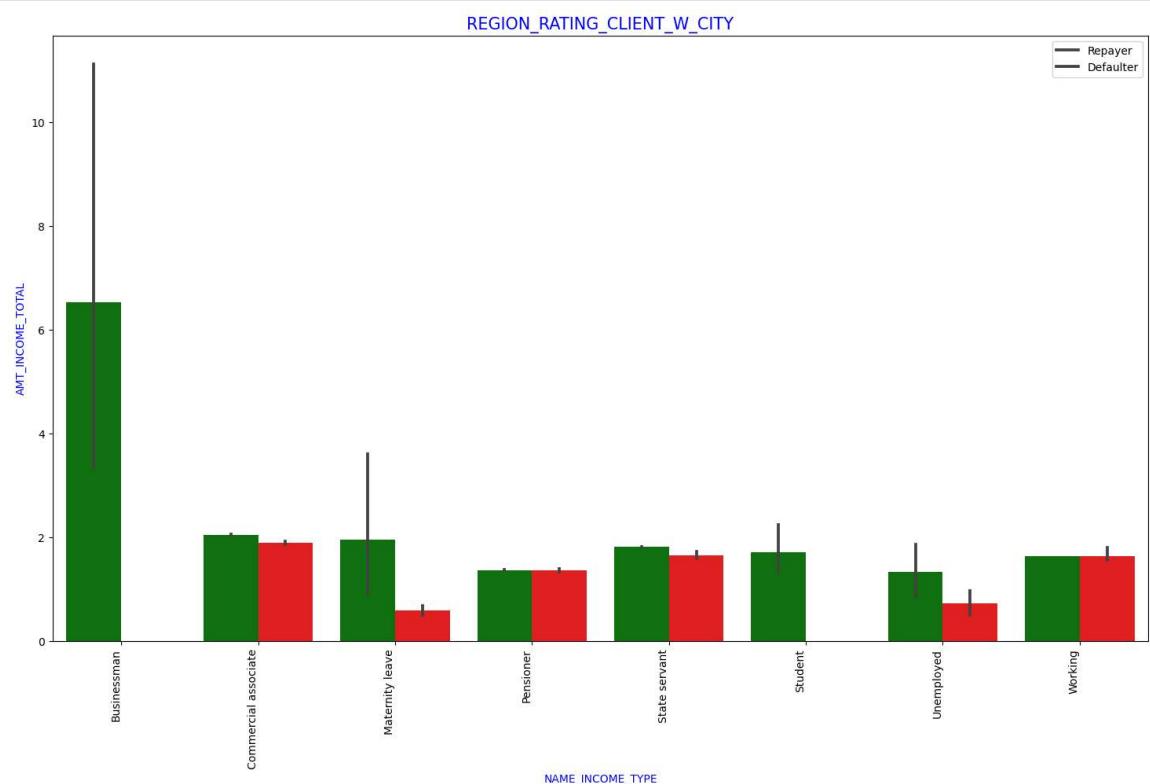
In [112]: 1 app_data.groupby('NAME_INCOME_TYPE')[['AMT_INCOME_TOTAL']].describe()

Out[112]:

NAME_INCOME_TYPE	count	mean	std	min	25%	50%	75%	max
Businessman	10.0	6.525000	6.272260	1.8000	2.250	4.9500	8.43750	22.5000
Commercial associate	71617.0	2.029553	1.479742	0.2655	1.350	1.8000	2.25000	180.0000
Maternity leave	5.0	1.404000	1.268569	0.4950	0.675	0.9000	1.35000	3.6000
Pensioner	55362.0	1.364013	0.766503	0.2565	0.900	1.1700	1.66500	22.5000
State servant	21703.0	1.797380	1.008806	0.2700	1.125	1.5750	2.25000	31.5000
Student	18.0	1.705000	1.066447	0.8100	1.125	1.5750	1.78875	5.6250
Unemployed	22.0	1.105364	0.880551	0.2655	0.540	0.7875	1.35000	3.3750
Working	158774.0	1.631699	3.075777	0.2565	1.125	1.3500	2.02500	1170.0000

In [113]: 1 # Income type vs Income Amount Range

2 bivariate_bar("NAME_INCOME_TYPE", "AMT_INCOME_TOTAL", app_data, "TARGET", (



It can be seen that business man's income is the highest and the estimated range with default 95% confidence level seem to indicate that the income of a business man could be in the range of slightly close to 4 lakhs and slightly above 10 lakhs

Numeric Variables Analysis

```
In [116]: 1 # Bifurcating the app-data dataframe based on Target value 0 and 1 for
2 app_data.columns
3
```

```
Out[116]: Index(['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', 'NAME_INCOME_TYPE', 'NAME_EDUCATION_TYPE',
   'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE', 'REGION_POPULATION_RELATIVE',
   'DAYS_BIRTH',
   'DAYS_EMPLOYED', 'DAYS_REGISTRATION', 'DAYS_ID_PUBLISH', 'OCCUPATION_TYPE',
   'CNT_FAM_MEMBERS', 'REGION_RATING_CLIENT', 'REGION_RATING_CLIENT_W_CITY',
   'WEEKDAY_APPR_PROCESS_START', 'HOUR_APPR_PROCESS_START', 'REG_REGION_NOT_LIVE_REGION',
   'REG_REGION_NOT_WORK_REGION', 'LIVE_REGION_NOT_WORK_REGION', 'REG_CITY_NOT_LIVE_CITY',
   'REG_CITY_NOT_WORK_CITY', 'LIVE_CITY_NOT_WORK_CITY', 'ORGANIZATION_TYPE',
   'OBS_30_CNT_SOCIAL_CIRCLE', 'DEF_30_CNT_SOCIAL_CIRCLE', 'OBS_60_CNT_SOCIAL_CIRCLE',
   'DEF_60_CNT_SOCIAL_CIRCLE', 'DAYS_LAST_PHONE_CHANGE', 'FLAG_DOCUMENT_3',
   'AMT_REQ_CREDIT_BUREAU_HOUR', 'AMT_REQ_CREDIT_BUREAU_DAY', 'AMT_REQ_CREDIT_BUREAU_WEEK',
   'AMT_REQ_CREDIT_BUREAU_MON', 'AMT_REQ_CREDIT_BUREAU_QRT', 'AMT_REQ_CREDIT_BUREAU_YEAR',
   'AMT_INCOME_RANGE', 'AMT_CREDIT_RANGE', 'AGE', 'AGE_GROUP', 'YEARS_EMPLOYED',
   'EMPLOYMENT_YEAR'],
  dtype='object')
```

```
In [118]: 1 # Bifurcating the app-data dataframe based on Target value 0 and 1 for
2 cols_for_correlation = ['NAME_CONTRACT_TYPE', 'CODE_GENDER', 'FLAG_OWN_
3   'CNT_CHILDREN', 'AMT_INCOME_TOTAL', 'AMT_CREDIT'
4   'NAME_TYPE_SUITE', 'NAME_INCOME_TYPE', 'NAME_EDU
5   'NAME_HOUSING_TYPE', 'REGION_POPULATION_RELATIV
6   'DAYS_REGISTRATION', 'DAYS_ID_PUBLISH', 'OCCUPA
7   'REGION_RATING_CLIENT_W_CITY', 'WEEKDAY_APPR_PR
8   'REG_REGION_NOT_LIVE_REGION', 'REG_REGION_NOT_W
9   'REG_CITY_NOT_LIVE_CITY', 'REG_CITY_NOT_WORK_CI
10  'OBS_60_CNT_SOCIAL_CIRCLE', 'DEF_60_CNT_SOCIAL_
11  'AMT_REQ_CREDIT_BUREAU_HOUR', 'AMT_REQ_CREDIT_B
12  'AMT_REQ_CREDIT_BUREAU_MON', 'AMT_REQ_CREDIT_BU
13
14
15 Repayer_df = app_data.loc[app_data['TARGET']==0, cols_for_correlation]
16 Defaulter_df = app_data.loc[app_data['TARGET']==1, cols_for_correlation]
```

Correlation between numeric variable

In [120]:

```

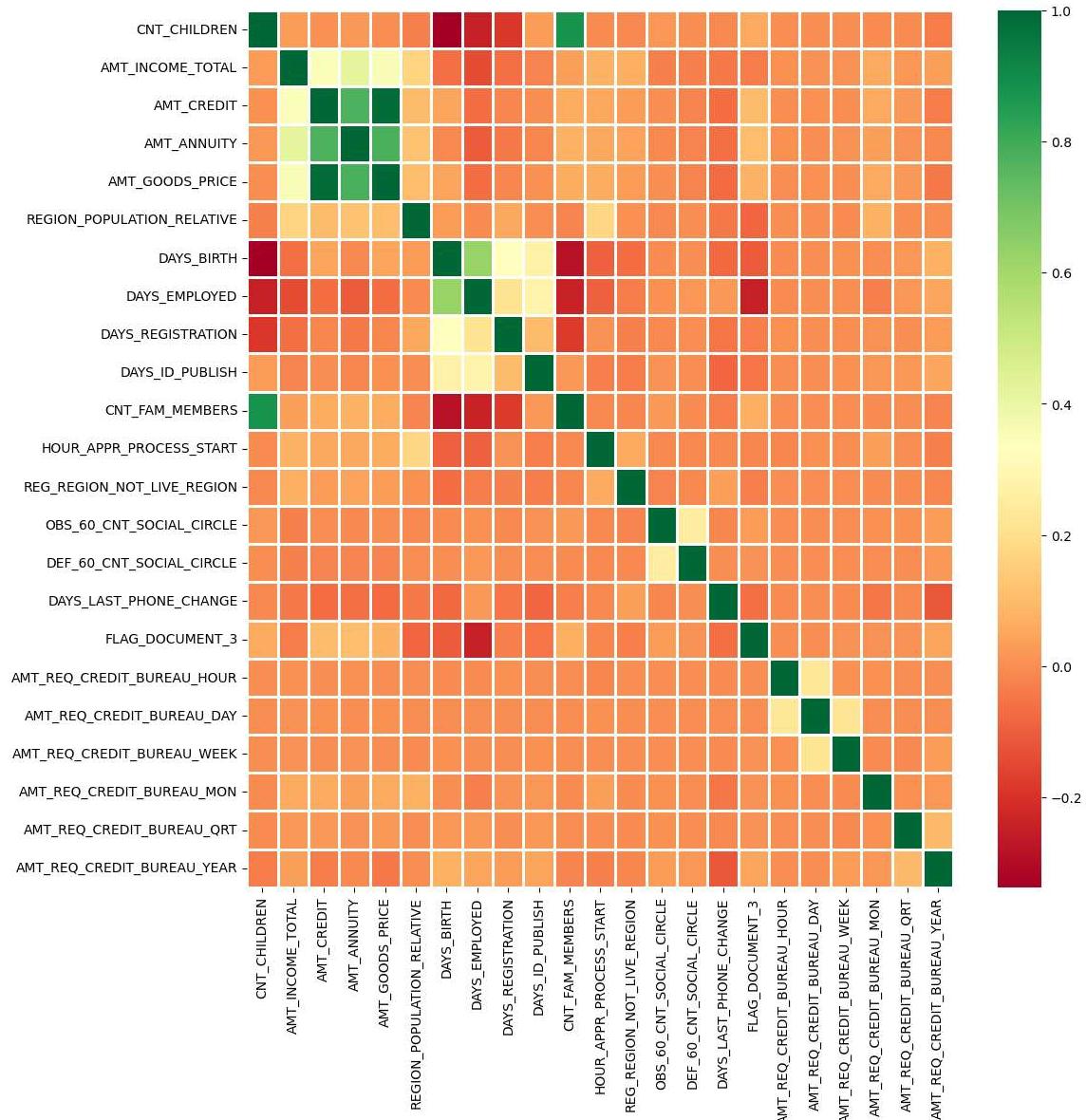
1 # Getting the top 10 correlation for the Repayers data
2 corr_repayer = Repayer_df.corr()
3 corr_repayer = corr_repayer.where(np.triu(np.ones(corr_repayer.shape),
4 corr_df_repayer = corr_repayer.unstack().reset_index()
5 corr_df_repayer.columns = ['VAR1', 'VAR2', 'Correlation']
6 corr_df_repayer.dropna(subset=["Correlation"], inplace=True)
7 corr_df_repayer["Correlation"] = corr_df_repayer["Correlation"].abs()
8 corr_df_repayer.sort_values(by='Correlation', ascending=False, inplace=True)
9 corr_df_repayer.head(10)

```

Out[120]:

	VAR1	VAR2	Correlation
94	AMT_GOODS_PRICE	AMT_CREDIT	0.987250
230	CNT_FAM_MEMBERS	CNT_CHILDREN	0.878571
95	AMT_GOODS_PRICE	AMT_ANNUITY	0.776686
71	AMT_ANNUITY	AMT_CREDIT	0.771309
167	DAYS_EMPLOYED	DAYS_BIRTH	0.626114
70	AMT_ANNUITY	AMT_INCOME_TOTAL	0.418953
93	AMT_GOODS_PRICE	AMT_INCOME_TOTAL	0.349462
47	AMT_CREDIT	AMT_INCOME_TOTAL	0.342799
138	DAYS_BIRTH	CNT_CHILDREN	0.336966
190	DAYS_REGISTRATION	DAYS_BIRTH	0.333151

```
In [121]: 1 fig = plt.figure(figsize=(12,12))
2 ax = sns.heatmap(Repayer_df.corr(), cmap="RdYlGn", annot=False, linewidths=1)
```



Correlating factors amongst repayers: Credit amount is highly correlated with amount of goods price loan annuity total income We can also see that repayers have high correlation in number of days employed.

In [123]:

```

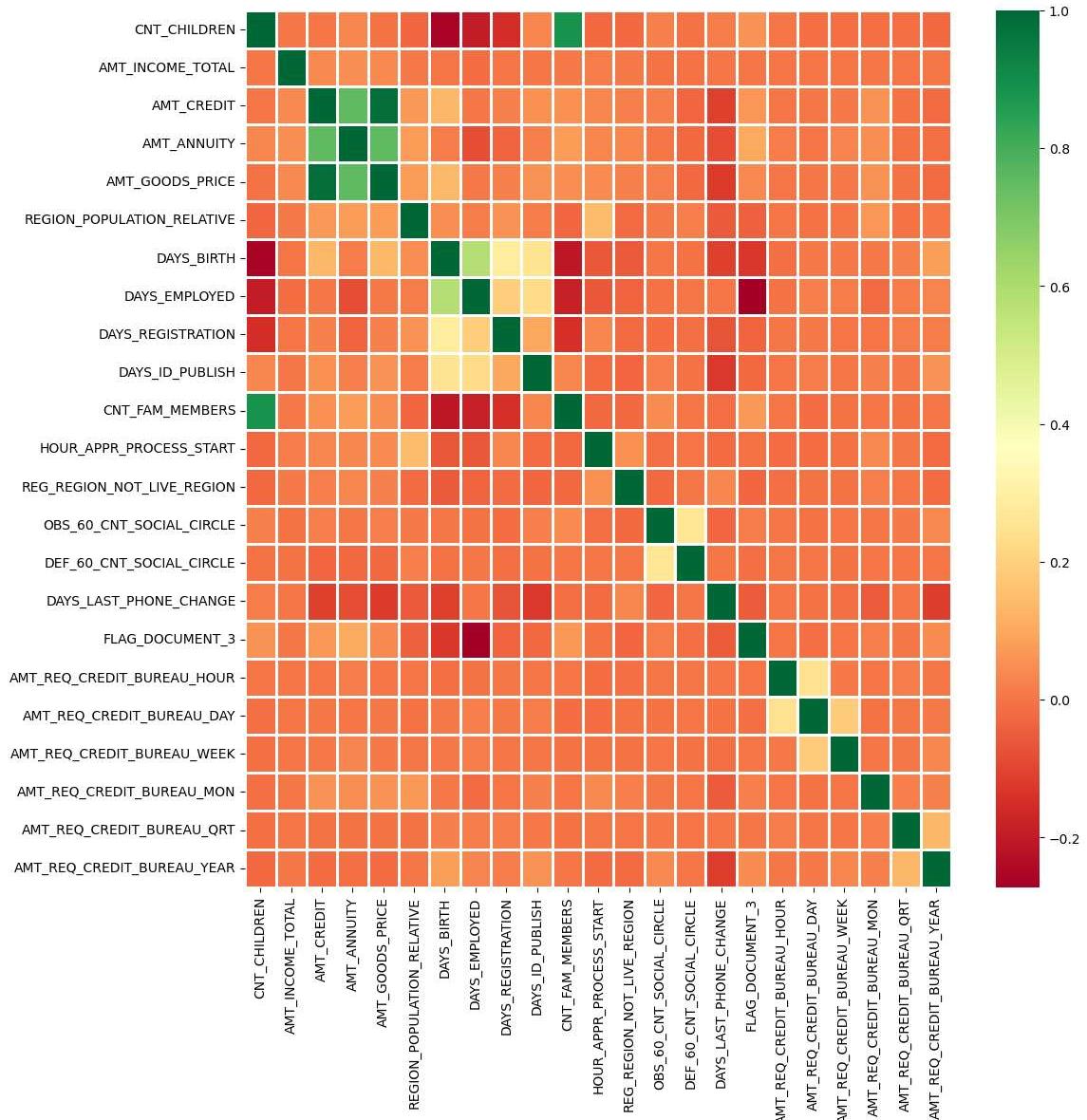
1 # Getting the top 10 correlation for the Defaulter data
2 corr_Defaulter = Defaulter_df.corr()
3 corr_Defaulter = corr_Defaulter.where(np.triu(np.ones(corr_Defaulter.shape), k=1).astype(bool))
4 corr_df_Defaulter = corr_Defaulter.unstack().reset_index()
5 corr_df_Defaulter.columns =['VAR1','VAR2','Correlation']
6 corr_df_Defaulter.dropna(subset = ["Correlation"], inplace = True)
7 corr_df_Defaulter["Correlation"] =corr_df_Defaulter["Correlation"].abs()
8 corr_df_Defaulter.sort_values(by='Correlation', ascending=False, inplace=True)
9 corr_df_Defaulter.head(10)

```

Out[123]:

	VAR1	VAR2	Correlation
94	AMT_GOODS_PRICE	AMT_CREDIT	0.983103
230	CNT_FAM_MEMBERS	CNT_CHILDREN	0.885484
95	AMT_GOODS_PRICE	AMT_ANNUITY	0.752699
71	AMT_ANNUITY	AMT_CREDIT	0.752195
167	DAY_S_EMPLOYED	DAY_S_BIRTH	0.582185
190	DAY_S_REGISTRATION	DAY_S_BIRTH	0.289114
375	FLAG_DOCUMENT_3	DAY_S_EMPLOYED	0.272169
335	DEF_60_CNT_SOCIAL_CIRCLE	OBS_60_CNT_SOCIAL_CIRCLE	0.264159
138	DAY_S_BIRTH	CNT_CHILDREN	0.259109
213	DAY_S_ID_PUBLISH	DAY_S_BIRTH	0.252863

```
In [124]: 1 fig = plt.figure(figsize=(12,12))
2 ax = sns.heatmap(Defaulter_df.corr(), cmap="RdYlGn", annot=False, linewidths=1)
```



Credit amount is highly correlated with amount of goods price which is same as repayers. But the loan annuity correlation with credit amount has slightly reduced in defaulters(0.75) when compared to repayers(0.77) We can also see that repayers have high correlation in number of days employed(0.62) when compared to defaulters(0.58). There is a severe drop in the correlation between total income of the client and the credit amount(0.038) amongst defaulters whereas it is 0.342 among repayers. Days_birth and number of children correlation has reduced to 0.259 in defaulters when compared to 0.337 in repayers. There is a slight increase in defaulted to observed count in social circle among defaulters(0.264) when compared to repayers(0.254)

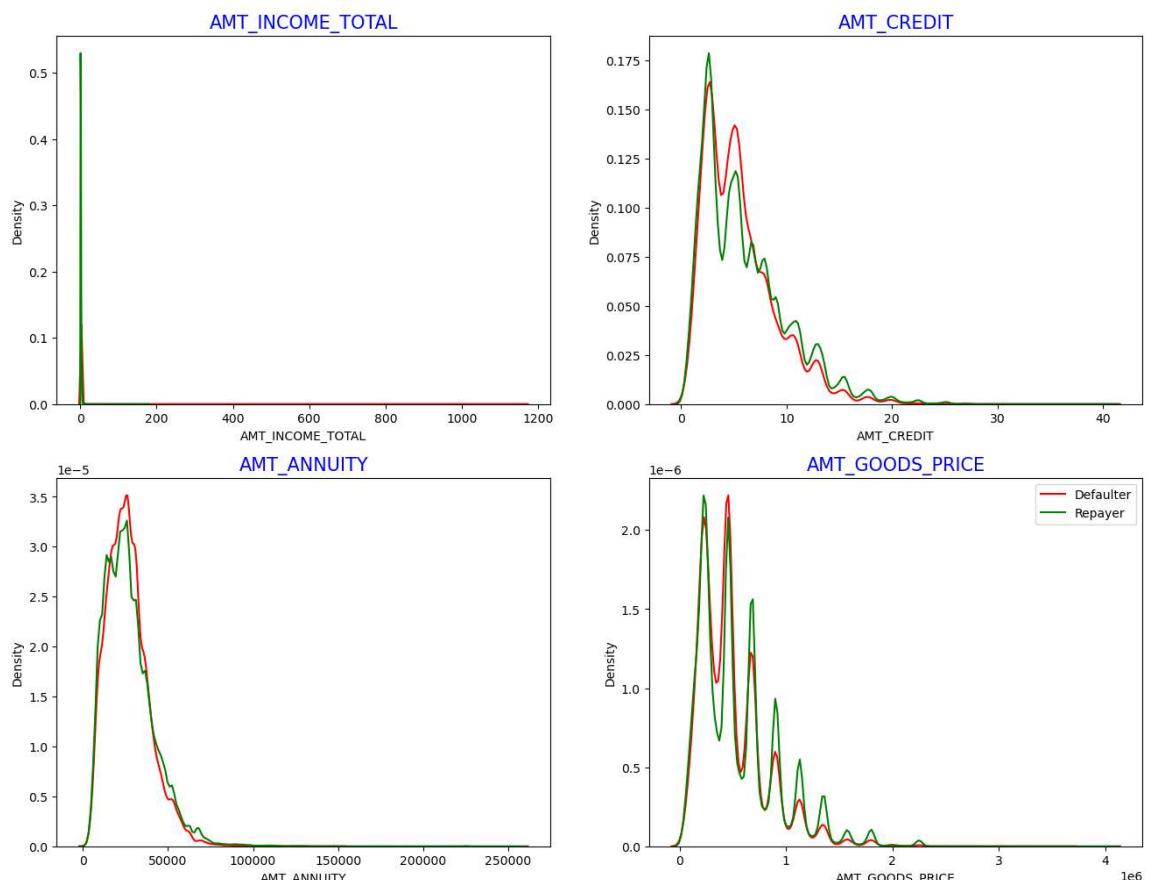
Numerical Univariate Analysis

In [125]:

```

1 # Plotting the numerical columns related to amount as distribution plot
2 amount = app_data[['AMT_INCOME_TOTAL','AMT_CREDIT','AMT_ANNUITY', 'AMT_GOODS_PRICE']]
3
4 fig = plt.figure(figsize=(16,12))
5
6 for i in enumerate(amount):
7     plt.subplot(2,2,i[0]+1)
8     sns.distplot(Defaulter_df[i[1]], hist=False, color='r',label ="Defaulter")
9     sns.distplot(Repayer_df[i[1]], hist=False, color='g', label ="Repayer")
10    plt.title(i[1], fontdict={'fontsize' : 15, 'fontweight' : 5, 'color': 'black'})
11
12    plt.legend()
13
14 plt.show()

```



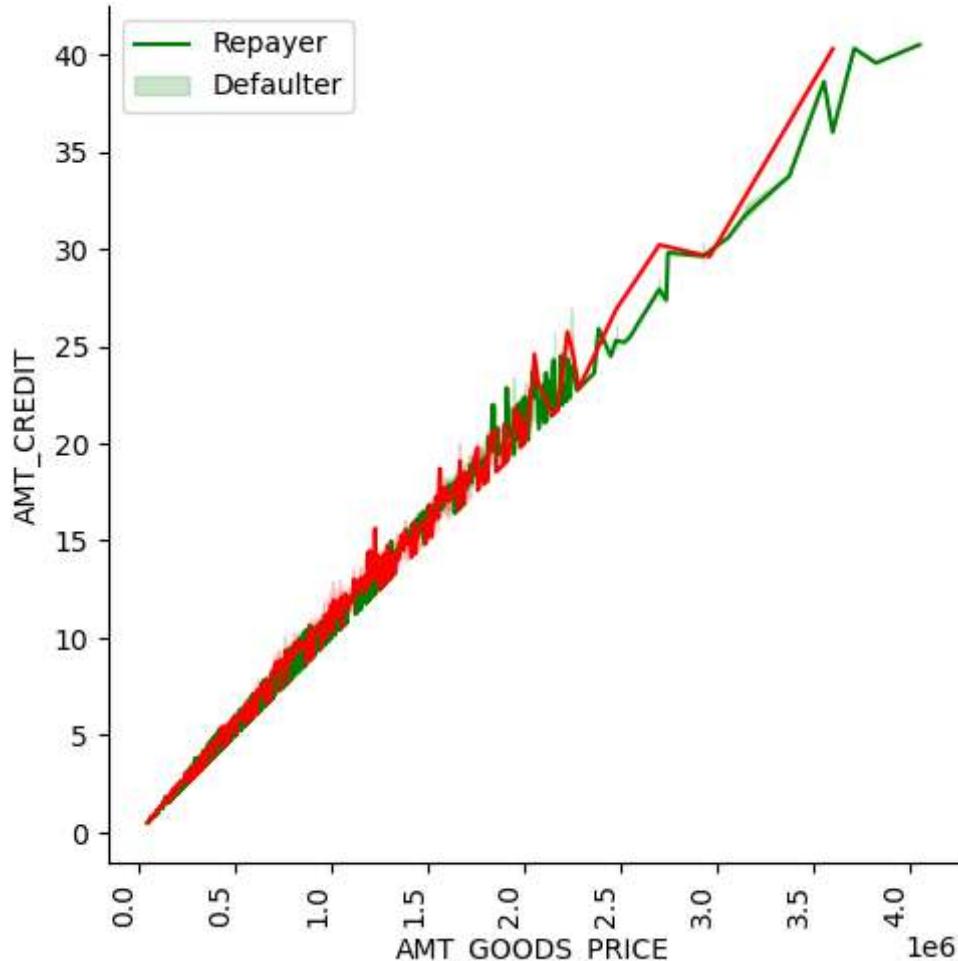
Most no of loans are given for goods price below 10 lakhs Most people pay annuity below 50000 for the credit loan Credit amount of the loan is mostly less then 10 lakhs The repayers and defaulters distribution overlap in all the plots and hence we cannot use any of these variables in isolation to make a decision

Numerical Bivariate Analysis

In [126]:

```
1 # Checking the relationship between Goods price and credit and comparin
2 bivariate_rel('AMT_GOODS_PRICE','AMT_CREDIT',app_data,"TARGET", "line",
3
```

<Figure size 1500x600 with 0 Axes>



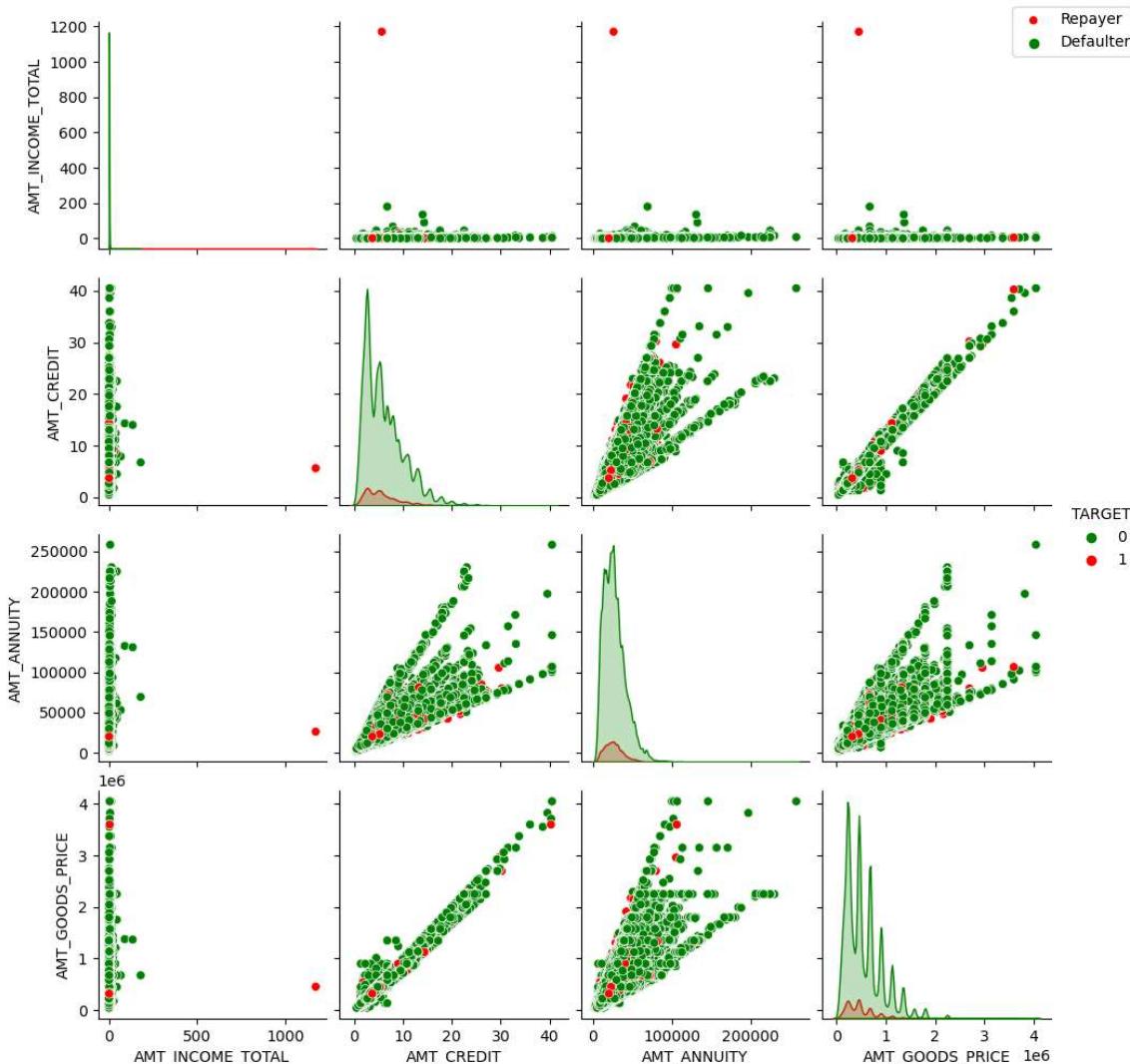
When the credit amount goes beyond 3M, there is an increase in defaulters.

In [127]:

```

1 # Plotting pairplot between amount variable to draw reference against L
2 amount = app_data[['AMT_INCOME_TOTAL', 'AMT_CREDIT',
3                     'AMT_ANNUITY', 'AMT_GOODS_PRICE', 'TARGET']]
4 amount = amount[(amount["AMT_GOODS_PRICE"].notnull()) & (amount["AMT_ANNUITY"].notnull())]
5 ax= sns.pairplot(amount,hue="TARGET",palette=["g","r"])
6 ax.fig.legend(labels=['Repayer', 'Defaulter'])
7 plt.show()

```



When amt_annuity >15000 amt_goods_price> 3M, there is a lesser chance of defaulters
 AMT_CREDIT and AMT_GOODS_PRICE are highly correlated as based on the scatterplot
 where most of the data are consolidated in form of a line There are very less defaulters for
 AMT_CREDIT >3M Inferences related to distribution plot has been already mentioned in
 previous distplot graphs inferences section

Merged Dataframes Analysis

```

1 Conclusions
2
3 Decisive Factor whether an applicant will be Repayer:
4 NAME_EDUCATION_TYPE: Academic degree has less defaults.
5 NAME_INCOME_TYPE: Student and Businessmen have no defaults.
6 REGION_RATING_CLIENT: RATING 1 is safer.
7 ORGANIZATION_TYPE: Clients with Trade Type 4 and 5 and Industry type
8 have defaulted less than 3%
8 DAYS_BIRTH: People above age of 50 have low probability of defaulting

```

9 DAYS_EMPLOYED: Clients with 40+ year experience having less than 1% default rate
10 AMT_INCOME_TOTAL: Applicant with Income more than 700,000 are less likely to default
11 NAME_CASH_LOAN_PURPOSE: Loans bought for Hobby, Buying garage are being repayed mostly.
12 CNT_CHILDREN: People with zero to two children tend to repay the loans.
13
14
15 Decisive Factor whether an applicant will be Defaulter:
16 CODE_GENDER: Men are at relatively higher default rate
17 NAME_FAMILY_STATUS : People who have civil marriage or who are single default a lot.
18 NAME_EDUCATION_TYPE: People with Lower Secondary & Secondary education
19 NAME_INCOME_TYPE: Clients who are either at Maternity leave OR Unemployed default a lot.
20 REGION_RATING_CLIENT: People who live in Rating 3 has highest defaults.
21 OCCUPATION_TYPE: Avoid Low-skill Laborers, Drivers and Waiters/barmen staff, Security staff, Laborers and Cooking staff as the default rate is huge.
22 ORGANIZATION_TYPE: Organizations with highest percent of loans not repaid are Transport: type 3 (16%),
23 Industry: type 13 (13.5%), Industry: type 8 (12.5%) and Restaurant (less than 12%). Self-employed people have relative high defaulting rate, and thus should be avoided to be approved for loan or provide loan with higher interest rate to mitigate the risk of defaulting.
24 DAYS_BIRTH: Avoid young people who are in age group of 20-40 as they have higher probability of defaulting
25 DAYS_EMPLOYED: People who have less than 5 years of employment have high default rate.
26 CNT_CHILDREN & CNT_FAM_MEMBERS: Client who have children equal to or more than 9 default 100% and hence their applications are to be rejected.
27 AMT_GOODS_PRICE: When the credit amount goes beyond 3M, there is an increase in defaulters.
28
29
30
31
32
33 NAME_HOUSING_TYPE: High number of loan applications are from the category of people who live in
34 Rented apartments & living with parents and hence offering the loan would mitigate the loss if any of those default.
35 AMT_CREDIT: People who get loan for 300-600k tend to default more than others and hence having higher interest specifically for this credit range would be ideal.
36 AMT_INCOME: Since 90% of the applications have Income total less than 300,000 and they have high probability of defaulting, they could be offered loan with higher interest compared to other income category.
37 CNT_CHILDREN & CNT_FAM_MEMBERS: Clients who have 4 to 8 children has a very high default rate and hence higher interest should be imposed on their loans.
38 NAME_CASH_LOAN_PURPOSE: Loan taken for the purpose of Repairs seems to have highest default rate.
39 high loan interest rate which is not feasible by the clients, thus they refuse the loan.

45 | The same approach could be followed in future as well.

In []:

1