

In [1]:

```
#importing all the important libraries
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline
import scipy.stats as stats
```

In [36]:

```
# to suppress warnings

import warnings
warnings.filterwarnings("ignore")
```

In [39]:

```
#notebook setting to display all the rows and columns to have better clarity on the data.

pd.set_option('display.max_rows', 500)
pd.set_option('display.max_columns', 500)
pd.set_option('display.width', 1000)
pd.set_option('display.expand_frame_repr', False)
```

In [23]:

```
#importing data
loan=pd.read_csv(r"C:\Users\santhosh\Videos\application_data.csv")
loan.head()
```

Out[23]:

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY
0	100002	1	Cash loans	M	N	N
1	100003	0	Cash loans	F	N	N
2	100004	0	Revolving loans	M	Y	N
3	100006	0	Cash loans	F	N	N
4	100007	0	Cash loans	M	N	N

5 rows × 122 columns

In [24]:

```
#checking rows and column
loan.shape
```

Out[24]:

(307511, 122)

In [29]:

```
#Checking the numeric variables
loan.describe()
```

Out[29]:

	SK_ID_CURR	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_TERM
count	307511.000000	307511.000000	307511.000000	3.075110e+05	3.075110e+05	307
mean	278180.518577	0.080729	0.417052	1.687979e+05	5.990260e+05	27
std	102790.175348	0.272419	0.722121	2.371231e+05	4.024908e+05	14
min	100002.000000	0.000000	0.000000	2.565000e+04	4.500000e+04	1
25%	189145.500000	0.000000	0.000000	1.125000e+05	2.700000e+05	16
50%	278202.000000	0.000000	0.000000	1.471500e+05	5.135310e+05	24
75%	367142.500000	0.000000	1.000000	2.025000e+05	8.086500e+05	34
max	456255.000000	1.000000	19.000000	1.170000e+08	4.050000e+06	258

8 rows × 106 columns

In [8]:

```
#checking how many null values are present in each of the columns

#creating a function to find null values for the dataframe
def null_values(df):
    return round((df.isnull().sum()*100/len(df)).sort_values(ascending = False),2)
```

In [41]:

null_values(loan)

Out[41]:

```
COMMONAREA_MEDI      69.87
COMMONAREA_AVG       69.87
COMMONAREA_MODE      69.87
NONLIVINGAPARTMENTS_MODE  69.43
NONLIVINGAPARTMENTS_AVG  69.43
NONLIVINGAPARTMENTS_MEDI  69.43
FONDKAPREMONT_MODE   68.39
LIVINGAPARTMENTS_MODE  68.35
LIVINGAPARTMENTS_AVG  68.35
LIVINGAPARTMENTS_MEDI  68.35
FLOORSMIN_AVG        67.85
FLOORSMIN_MODE       67.85
FLOORSMIN_MEDI       67.85
YEARS_BUILD_MEDI     66.50
YEARS_BUILD_MODE     66.50
YEARS_BUILD_AVG      66.50
OWN_CAR_AGE          65.99
LANDAREA_MEDI        59.38
```

In [42]:

```
#creating a variable null_col_50 for storing null columns having missing values more than 50%

null_col_50 = null_values(loan)[null_values(loan)>50]
```

In [43]:

```
#reviewing null_col_50

print(null_col_50)
print()
print("Num of columns having missing values more than 50% :",len(null_col_50))
```

COMMONAREA_MEDI	69.87
COMMONAREA_AVG	69.87
COMMONAREA_MODE	69.87
NONLIVINGAPARTMENTS_MODE	69.43
NONLIVINGAPARTMENTS_AVG	69.43
NONLIVINGAPARTMENTS_MEDI	69.43
FONDKAPREMONT_MODE	68.39
LIVINGAPARTMENTS_MODE	68.35
LIVINGAPARTMENTS_AVG	68.35
LIVINGAPARTMENTS_MEDI	68.35
FLOORSMIN_AVG	67.85
FLOORSMIN_MODE	67.85
FLOORSMIN_MEDI	67.85
YEARS_BUILD_MEDI	66.50
YEARS_BUILD_MODE	66.50
YEARS_BUILD_AVG	66.50
OWN_CAR_AGE	65.99
LANDAREA_MEDI	59.38
LANDAREA_MODE	59.38
LANDAREA_AVG	59.38
BASEMENTAREA_MEDI	58.52
BASEMENTAREA_AVG	58.52
BASEMENTAREA_MODE	58.52
EXT_SOURCE_1	56.38
NONLIVINGAREA_MODE	55.18
NONLIVINGAREA_AVG	55.18
NONLIVINGAREA_MEDI	55.18
ELEVATORS_MEDI	53.30
ELEVATORS_AVG	53.30
ELEVATORS_MODE	53.30
WALLSMATERIAL_MODE	50.84
APARTMENTS_MEDI	50.75
APARTMENTS_AVG	50.75
APARTMENTS_MODE	50.75
ENTRANCES_MEDI	50.35
ENTRANCES_AVG	50.35
ENTRANCES_MODE	50.35
LIVINGAREA_AVG	50.19
LIVINGAREA_MODE	50.19
LIVINGAREA_MEDI	50.19
HOUSETYPE_MODE	50.18

dtype: float64

Num of columns having missing values more than 50% : 41

In [44]:

```
null_col_50.index # Will drop all these columns
```

Out[44]:

```
Index(['COMMONAREA_MEDI', 'COMMONAREA_AVG', 'COMMONAREA_MODE', 'NONLIVINGAPARTMENTS_MODE', 'NONLIVINGAPARTMENTS_AVG', 'NONLIVINGAPARTMENTS_MEDI', 'FONDKAPREMONT_MODE', 'LIVINGAPARTMENTS_MODE', 'LIVINGAPARTMENTS_AVG', 'LIVINGAPARTMENTS_MEDI', 'FLOORSMIN_AVG', 'FLOORSMIN_MODE', 'FLOORSMIN_MEDI', 'YEARS_BUILD_MEDI', 'YEARS_BUILD_MODE', 'YEARS_BUILD_AVG', 'OWN_CAR_AGE', 'LANDAREA_MEDI', 'LANDAREA_MODE', 'LANDAREA_AVG', 'BASEMENTAREA_MEDI', 'BASEMENTAREA_AVG', 'BASEMENTAREA_MODE', 'EXT_SOURCE_1', 'NONLIVINGAREA_MODE', 'NONLIVINGAREA_AVG', 'NONLIVINGAREA_MEDI', 'ELEVATORS_MEDI', 'ELEVATORS_AVG', 'ELEVATORS_MODE', 'WALLSMATERIAL_MODE', 'APARTMENTS_MEDI', 'APARTMENTS_AVG', 'APARTMENTS_MODE', 'ENTRANCES_MEDI', 'ENTRANCES_AVG', 'ENTRANCES_MODE', 'LIVINGAREA_AVG', 'LIVINGAREA_MODE', 'LIVINGAREA_MEDI', 'HOUSETYPE_MODE'], dtype='object')
```

In [45]:

```
# Now Lets drop all the columns having missing values more than 50% that is 41 columns

loan.drop(columns = null_col_50.index, inplace = True)
```

In [56]:

```
loan.shape # Now there are 81 columns remaining
```

Out[56]:

```
(307511, 73)
```

In [57]:

```
# now we will deal with null values more than 15%

null_col_15 = null_values(loan)[null_values(loan)>15]
null_col_15
#removing 'OCCUPATION_TYPE', 'EXT_SOURCE_3' from "null_col_15" so that we can drop all other

null_col_15.drop(["OCCUPATION_TYPE", "EXT_SOURCE_3"], inplace = True)

print(null_col_15)
print()
print("No of columns having missing values more than 15% and are not reletable:", len(null_col_15))

Series([], dtype: float64)
```

```
No of columns having missing values more than 15% and are not reletable: 0
```

In [58]:

```
#thus removing columns having missing values more than 15% and which are not reletable to T
loan.drop(null_col_15.index,axis=1, inplace = True)
```

In [59]:

```
loan.shape # After dropping null_col_15, we have left with 73 columns
null_values(loan).head(10)
```

Out[59]:

```
OCCUPATION_TYPE      31.35
EXT_SOURCE_3         19.83
AMT_REQ_CREDIT_BUREAU_YEAR  13.50
AMT_REQ_CREDIT_BUREAU_QRT  13.50
AMT_REQ_CREDIT_BUREAU_MON  13.50
AMT_REQ_CREDIT_BUREAU_WEEK  13.50
AMT_REQ_CREDIT_BUREAU_DAY  13.50
AMT_REQ_CREDIT_BUREAU_HOUR  13.50
NAME_TYPE_SUITE       0.42
OBS_30_CNT_SOCIAL_CIRCLE  0.33
dtype: float64
```

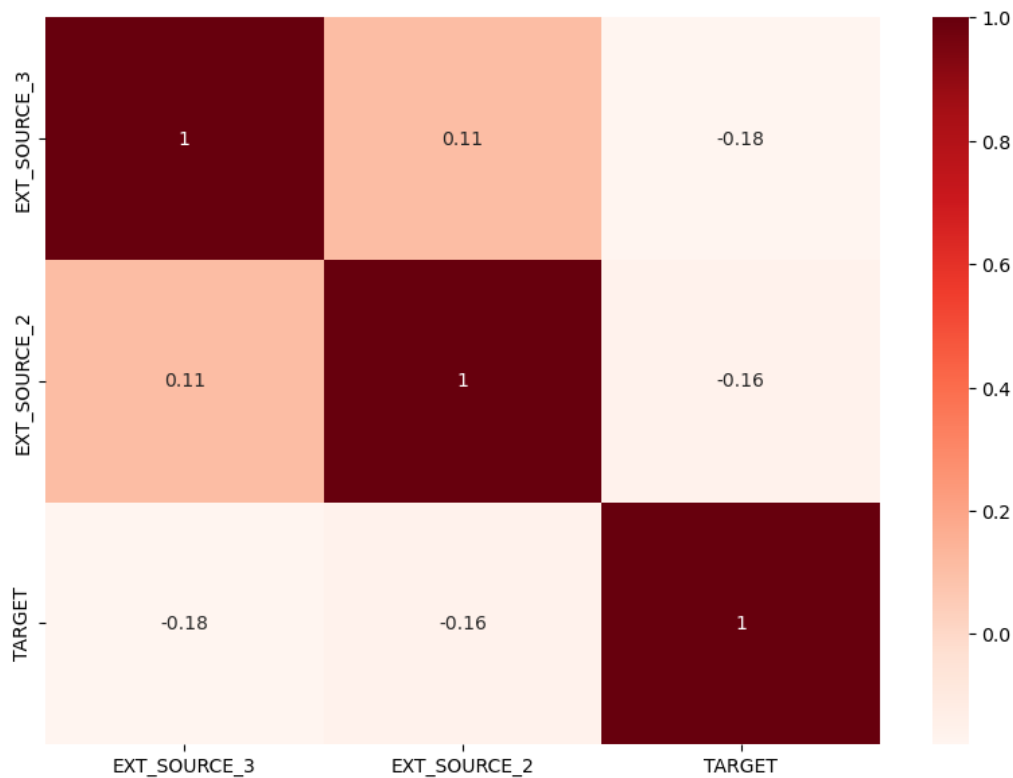
In [60]:

```
irrev = ["EXT_SOURCE_3", "EXT_SOURCE_2"] # putting irrevlent columns in varibale "irrev"
```

In [62]:

```
plt.figure(figsize= [10,7])
sns.heatmap(loan[irrev+["TARGET"]].corr(), cmap="Reds",annot=True)
plt.title("Correlation between EXT_SOURCE_3, EXT_SOURCE_2, TARGET", fontdict={"fontsize":20})
plt.show()
```

Correlation between EXT_SOURCE_3, EXT_SOURCE_2, TARGET

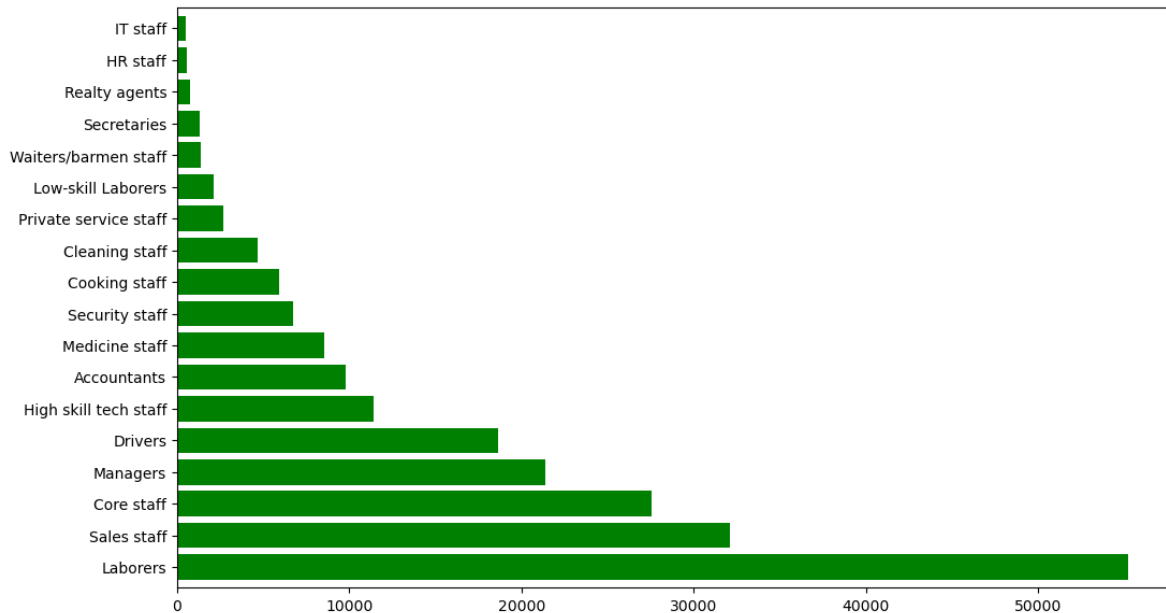


In [64]:

```
# Plotting a percentage graph having each category of "OCCUPATION_TYPE"

plt.figure(figsize = [12,7])
(loan["OCCUPATION_TYPE"].value_counts()).plot.barh(color= "green",width = .8)
plt.title("Percentage of Type of Occupations", fontdict={"fontsize":20}, pad =20)
plt.show()
```

Percentage of Type of Occupations



In [2]:

```
# importing previous_application.csv

prev_app1 = pd.read_csv(r"C:\Users\santhosh\Videos\previous_application.csv")
```

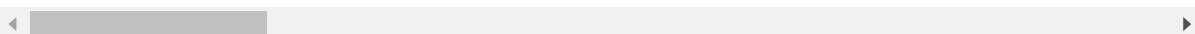
In [3]:

```
prev_app1.head()
```

Out[3]:

	SK_ID_PREV	SK_ID_CURR	NAME_CONTRACT_TYPE	AMT_ANNUITY	AMT_APPLICATION	AI
0	2030495	271877	Consumer loans	1730.430	17145.0	
1	2802425	108129	Cash loans	25188.615	607500.0	
2	2523466	122040	Cash loans	15060.735	112500.0	
3	2819243	176158	Cash loans	47041.335	450000.0	
4	1784265	202054	Cash loans	31924.395	337500.0	

5 rows × 37 columns



In [4]:

```
#Checking rows and columns of the raw data
prev_appl.shape
```

Out[4]:

(1670214, 37)

In [5]:

```
#Checking information of all the columns like data types
prev_appl.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1670214 entries, 0 to 1670213
Data columns (total 37 columns):
 #   Column                                Non-Null Count  Dtype
---  -
 0   SK_ID_PREV                            1670214 non-null int64
 1   SK_ID_CURR                            1670214 non-null int64
 2   NAME_CONTRACT_TYPE                   1670214 non-null object
 3   AMT_ANNUITY                          1297979 non-null float64
 4   AMT_APPLICATION                      1670214 non-null float64
 5   AMT_CREDIT                           1670213 non-null float64
 6   AMT_DOWN_PAYMENT                    774370 non-null float64
 7   AMT_GOODS_PRICE                     1284699 non-null float64
 8   WEEKDAY_APPR_PROCESS_START          1670214 non-null object
 9   HOUR_APPR_PROCESS_START             1670214 non-null int64
10   FLAG_LAST_APPL_PER_CONTRACT         1670214 non-null object
11   NFLAG_LAST_APPL_IN_DAY              1670214 non-null int64
12   RATE_DOWN_PAYMENT                   774370 non-null float64
13   RATE_INTEREST_PRIMARY               5951 non-null float64
14   RATE_INTEREST_SECOND                5951 non-null float64
```

In [6]:

```
# Checking the numeric variables of the dataframes
prev_appl.describe()
```

Out[6]:

	SK_ID_PREV	SK_ID_CURR	AMT_ANNUITY	AMT_APPLICATION	AMT_CREDIT	AMT_DOWN_PAYMENT
count	1.670214e+06	1.670214e+06	1.297979e+06	1.670214e+06	1.670213e+06	774370
mean	1.923089e+06	2.783572e+05	1.595512e+04	1.752339e+05	1.961140e+05	1.284699
std	5.325980e+05	1.028148e+05	1.478214e+04	2.927798e+05	3.185746e+05	1.670214
min	1.000001e+06	1.000010e+05	0.000000e+00	0.000000e+00	0.000000e+00	0.000000
25%	1.461857e+06	1.893290e+05	6.321780e+03	1.872000e+04	2.416050e+04	0.000000
50%	1.923110e+06	2.787145e+05	1.125000e+04	7.104600e+04	8.054100e+04	0.000000
75%	2.384280e+06	3.675140e+05	2.065842e+04	1.803600e+05	2.164185e+05	0.000000
max	2.845382e+06	4.562550e+05	4.180581e+05	6.905160e+06	6.905160e+06	1.670214

8 rows × 7 columns

In [9]:

```
#checking how many null values are present in each of the columns in percentage  
null_values(prev_appl)
```

Out[9]:

RATE_INTEREST_PRIVILEGED	99.64
RATE_INTEREST_PRIMARY	99.64
AMT_DOWN_PAYMENT	53.64
RATE_DOWN_PAYMENT	53.64
NAME_TYPE_SUITE	49.12
NFLAG_INSURED_ON_APPROVAL	40.30
DAYS_TERMINATION	40.30
DAYS_LAST_DUE	40.30
DAYS_LAST_DUE_1ST_VERSION	40.30
DAYS_FIRST_DUE	40.30
DAYS_FIRST_DRAWING	40.30
AMT_GOODS_PRICE	23.08
AMT_ANNUITY	22.29
CNT_PAYMENT	22.29
PRODUCT_COMBINATION	0.02
AMT_CREDIT	0.00
NAME_YIELD_GROUP	0.00
NAME_PORTFOLIO	0.00
NAME_SELLER_INDUSTRY	0.00
SELLERPLACE_AREA	0.00
CHANNEL_TYPE	0.00
NAME_PRODUCT_TYPE	0.00
SK_ID_PREV	0.00
NAME_GOODS_CATEGORY	0.00
NAME_CLIENT_TYPE	0.00
CODE_REJECT_REASON	0.00
SK_ID_CURR	0.00
DAYS_DECISION	0.00
NAME_CONTRACT_STATUS	0.00
NAME_CASH_LOAN_PURPOSE	0.00
NFLAG_LAST_APPL_IN_DAY	0.00
FLAG_LAST_APPL_PER_CONTRACT	0.00
HOURL_APPR_PROCESS_START	0.00
WEEKDAY_APPR_PROCESS_START	0.00
AMT_APPLICATION	0.00
NAME_CONTRACT_TYPE	0.00
NAME_PAYMENT_TYPE	0.00

dtype: float64

In [10]:

```
#creating a variable p_null_col_50 for storing null columns having missing values more than  
p_null_col_50 = null_values(prev_appl)[null_values(prev_appl)>50]
```

In [11]:

```
p_null_col_50 # There only 4 columns with missing valus more than 50%
```

Out[11]:

```
RATE_INTEREST_PRIVILEGED    99.64
RATE_INTEREST_PRIMARY       99.64
AMT_DOWN_PAYMENT           53.64
RATE_DOWN_PAYMENT          53.64
dtype: float64
```

In [12]:

```
#dropping null columns having missing values more than 50%
```

```
prev_appl.drop(columns = p_null_col_50.index, inplace = True)
```

In [14]:

```
#plotting a kdeplot to understand distribution of "AMT_ANNUITY"
```

```
plt.figure(figsize=(12,6))
sns.kdeplot(prev_appl['AMT_ANNUITY'])
plt.show()
```

