

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')

In [2]: cd('C:\Users\758449\Downloads')
C:\Users\758449\Downloads

In [3]: #read the file(data)
data = pd.read_csv('application_data.csv')
data

Out[3]:
```

SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT
0	1	Cash loans	M	N	Y	0	0
1	0	Cash loans	F	N	N	N	0
2	0	Revolving loans	M	Y	Y	Y	0
3	0	Cash loans	F	N	Y	Y	0
4	0	Cash loans	M	N	Y	Y	0
...	...	...	...	...	...	...	...
307506	456251	Cash loans	M	N	N	N	0
307507	456252	Cash loans	F	N	Y	Y	0
307508	456253	Cash loans	F	N	Y	Y	0
307509	456254	Cash loans	F	N	Y	Y	0
307510	456255	Cash loans	F	N	N	N	0

307511 rows x 122 columns

```
In [4]: #GET SOME Info from dataset like how many numerical features and categorical features and description
data.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 307511 entries, 0 to 307510
Columns: 122 entries, SK_ID_CURR to AMT_REQ_CREDIT_BUREAU_YEAR
dtypes: float64(65), int64(41), object(16)
memory usage: 286.2+ MB

from above dataset we found onefingh total 122 entries and floats are 65 and int are 41 and object are 16 and data have 307511 rows

In [5]: pd.set_option('float_format', '{:f}'.format)
data.describe()

Out[5]:
```

	SK_ID_CURR	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION
count	307511.000000	307511.000000	307511.000000	307511.000000	307511.000000	307499.000000	307233.000000	
mean	27180.158577	0.060729	0.470522	108787.919297	599025.999706	27108.973609	538316.204329	
std	102790.175348	0.272419	0.722121	237123.146279	402490.769966	14493.737315	369446.469440	
min	100002.000000	0.000000	0.000000	25650.000000	45000.000000	1615.000000	40500.000000	
25%	189145.000000	0.000000	0.000000	112500.000000	270000.000000	16524.000000	238500.000000	
50%	278202.000000	0.000000	0.000000	147150.000000	513531.000000	24903.000000	450000.000000	
75%	367142.000000	0.000000	1.000000	202500.000000	808650.000000	34596.000000	679500.000000	
max	456255.000000	1.000000	19.000000	117000000.000000	4050000.000000	258025.500000	4050000.000000	

8 rows x 100 columns

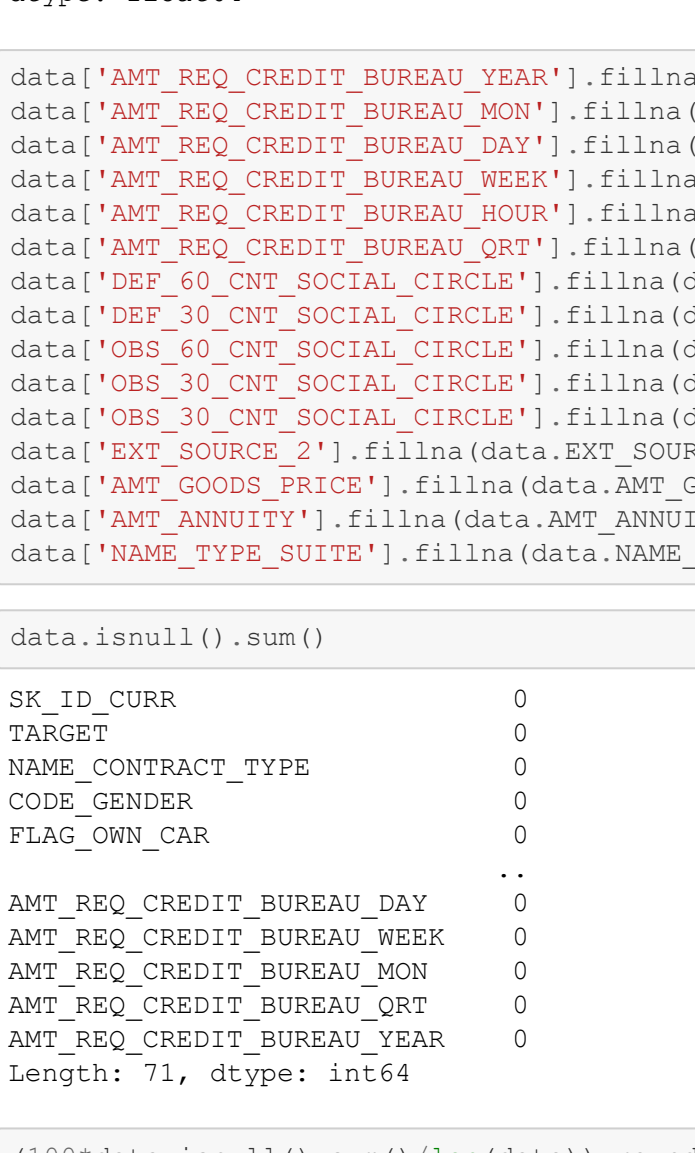
```
In [6]: #handling the missing values(replace with median)
pd.set_option('display.max_columns', None)
data.isnull().sum()

Out[6]:
```

	SK_ID_CURR	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION
AMT_REQ_CREDIT_BUREAU_DAY				41519				
AMT_REQ_CREDIT_BUREAU_WEEK				41519				
AMT_REQ_CREDIT_BUREAU_MON				41519				
AMT_REQ_CREDIT_BUREAU_QRT				41519				
AMT_REQ_CREDIT_BUREAU_YEAR				41519				
length:	122,	dtype:	int64					

```
In [7]: #100% data.isnull().sum()//len(data).round(2).sort_values(ascending=False)).head(40)
sns.distplot((100% data.isnull().sum()//len(data)).round(2).sort_values(ascending=False)).head(60)

Out[7]:
```



```
In [8]: (100% data.isnull().sum()//len(data)).round(2).sort_values(ascending=False).head(60).index
Index(['COMMONAREA_MEDI', 'COMMONAREA_AVG', 'COMMONAREA_MODE',
'NONLIVINGAPARTMENTS_MODE', 'NONLIVINGAPARTMENTS_MEDI',
'NONLIVINGAPARTMENTS_AVG', 'FLOORSAPARTMENTS_MODE',
'LIVINGAPARTMENTS_MEDI', 'LIVINGAPARTMENTS_MODE',
'LIVINGAPARTMENTS_AVG', 'FLOORSMIN_MEDI', 'FLOORSMIN_MODE',
'FLOORSMIN_AVG', 'YEARS_BUILD_MEDI', 'YEARS_BUILD_AVG',
'YEARS_BUILD_MODE', 'OWN_CAR_AGE', 'LANDAREA_MEDI', 'LANDAREA_AVG',
'LANDAREA_MODE', 'BASEMENTAREA_MEDI', 'BASEMENTAREA_AVG',
'BASEMENTAREA_MODE', 'EXT_SOURCE_1', 'NONLIVINGAREA_MEDI',
'NONLIVINGAREA_AVG', 'NONLIVINGAREA_MODE', 'ELEVATORS_MODE',
'ELEVATORS_AVG', 'ELEVATORS_MEDI', 'WALLSMATERIAL_MODE',
'APARTMENTS_MODE', 'APARTMENTS_AVG', 'APARTMENTS_MEDI',
'ENTRANCES_MEDI', 'ENTRANCES_MODE', 'ENTRANCES_AVG', 'LIVINGAREA_MEDI',
'LIVINGAREA_MODE', 'LIVINGAREA_AVG', 'HOUSETYPE_MODE', 'FLOORSMAX_MEDI',
'FLOORSMAX_AVG', 'FLOORSMAX_MODE', 'YEARS_BEGINEXPLUATATION_MEDI',
'YEARS_BEGINEXPLUATATION_AVG', 'YEARS_BEGINEXPLUATATION_MODE',
'TOTALAREA_MODE', 'EMERGENCYSITE_MODE', 'OCCUPATION_TYPE',
'EXT_SOURCE_3'],axis = 1,inplace = True)

In [9]: data.drop(['COMMONAREA_MEDI', 'COMMONAREA_AVG', 'COMMONAREA_MODE',
'NONLIVINGAPARTMENTS_MODE', 'NONLIVINGAPARTMENTS_MEDI',
'NONLIVINGAPARTMENTS_AVG', 'FLOORSAPARTMENTS_MODE',
'LIVINGAPARTMENTS_MEDI', 'LIVINGAPARTMENTS_MODE',
'LIVINGAPARTMENTS_AVG', 'FLOORSMIN_MEDI', 'FLOORSMIN_MODE',
'FLOORSMIN_AVG', 'YEARS_BUILD_MEDI', 'YEARS_BUILD_AVG',
'YEARS_BUILD_MODE', 'OWN_CAR_AGE', 'LANDAREA_MEDI', 'LANDAREA_AVG',
'LANDAREA_MODE', 'BASEMENTAREA_MEDI', 'BASEMENTAREA_AVG',
'BASEMENTAREA_MODE', 'EXT_SOURCE_1', 'NONLIVINGAREA_MEDI',
'NONLIVINGAREA_AVG', 'NONLIVINGAREA_MODE', 'ELEVATORS_MODE',
'ELEVATORS_AVG', 'ELEVATORS_MEDI', 'WALLSMATERIAL_MODE',
'APARTMENTS_MODE', 'APARTMENTS_AVG', 'APARTMENTS_MEDI',
'ENTRANCES_MEDI', 'ENTRANCES_MODE', 'ENTRANCES_AVG', 'LIVINGAREA_MEDI',
'LIVINGAREA_MODE', 'LIVINGAREA_AVG', 'HOUSETYPE_MODE', 'FLOORSMAX_MEDI',
'FLOORSMAX_AVG', 'FLOORSMAX_MODE', 'YEARS_BEGINEXPLUATATION_MEDI',
'YEARS_BEGINEXPLUATATION_AVG', 'YEARS_BEGINEXPLUATATION_MODE',
'TOTALAREA_MODE', 'EMERGENCYSITE_MODE', 'OCCUPATION_TYPE',
'EXT_SOURCE_3'],axis = 1,inplace = True)

In [10]: (100% data.isnull().sum()//len(data)).round(2).sort_values(ascending=False).head(20)

Out[10]:
```

	AMT_REQ_CREDIT_BUREAU_YEAR	AMT_REQ_CREDIT_BUREAU_MON	AMT_REQ_CREDIT_BUREAU_WEEK	AMT_REQ_CREDIT_BUREAU_DAY	AMT_REQ_CREDIT_BUREAU_QRT	NAME_TYPE_SUITE	DEF_60_CNT_SOCIAL_CIRCLE	OB6_60_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	OB3_30_CNT_SOCIAL_CIRCLE	EXT_SOURCE_2	AMT_GOODS_PRICE	DAYS_PUBLISH	FLAG_MOBIL	FLAG_WORK_PHONE	FLAG_EMP_PHONE	DAYS_EMPLOYED	FLAG_CONT_MOBILE	FLAG_PHONE
count	13.500000	13.500000	13.500000	13.500000	13.500000	0.420000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
mean	13.500000	13.500000	13.500000	13.500000	13.500000	0.420000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
std	13.500000	13.500000	13.500000	13.500000	13.500000	0.420000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
min	13.500000	13.500000	13.500000	13.500000	13.500000	0.420000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	13.500000	13.500000	13.500000	13.500000	13.500000	0.420000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	13.500000	13.500000	13.500000	13.500000	13.500000	0.420000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	13.500000	13.500000	13.500000	13.500000	13.500000	0.420000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	13.500000	13.500000	13.500000	13.500000	13.500000	0.420000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000

```
In [11]: data['AMT_REQ_CREDIT_BUREAU_YEAR'].fillna(data.AMT_REQ_CREDIT_BUREAU_YEAR.median(),inplace = True)
data['AMT_REQ_CREDIT_BUREAU_MON'].fillna(data.AMT_REQ_CREDIT_BUREAU_MON.median(),inplace = True)
data['AMT_REQ_CREDIT_BUREAU_WEEK'].fillna(data.AMT_REQ_CREDIT_BUREAU_WEEK.median(),inplace = True)
data['AMT_REQ_CREDIT_BUREAU_DAY'].fillna(data.AMT_REQ_CREDIT_BUREAU_DAY.median(),inplace = True)
data['AMT_REQ_CREDIT_BUREAU_QRT'].fillna(data.AMT_REQ_CREDIT_BUREAU_QRT.median(),inplace = True)
data['DEF_60_CNT_SOCIAL_CIRCLE'].fillna(data.OBS_60_CNT_SOCIAL_CIRCLE.median(),inplace = True)
data['OBS_60_CNT_SOCIAL_CIRCLE'].fillna(data.OBS_60_CNT_SOCIAL_CIRCLE.median(),inplace = True)
data['DEF_30_CNT_SOCIAL_CIRCLE'].fillna(data.OBS_30_CNT_SOCIAL_CIRCLE.median(),inplace = True)
data['OBS_30_CNT_SOCIAL_CIRCLE'].fillna(data.EXT_SOURCE_2.median(),inplace = True)
data['EXT_SOURCE_2'].fillna(data.EXT_SOURCE_2.median(),inplace = True)
data['NAME_TYPE_SUITE'].fillna(data.AMT_GOODS_PRICE.median(),inplace = True)
data['AMT_ANNUITY'].fillna(data.AMT_ANNUITY.median(),inplace = True)
data['NAME_TYPE_SUITE'].fillna(data.NAME_TYPE_SUITE.mode(),inplace = True)

In [12]: data.isnull().sum()

Out[12]:
```

	SK_ID_CURR	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION
AMT_REQ_CREDIT_BUREAU_DAY				0				
AMT_REQ_CREDIT_BUREAU_WEEK				0				
AMT_REQ_CREDIT_BUREAU_MON				0				
AMT_REQ_CREDIT_BUREAU_QRT				0				
AMT_REQ_CREDIT_BUREAU_YEAR				0				
length:	71,	dtype:	int64					

```
In [13]: (100% data.isnull().sum()//len(data)).round(2).sort_values(ascending=False).head(20)

Out[13]:
```

	NAME_TYPE_SUITE	AMT_REQ_CREDIT_BUREAU_YEAR	AMT_REQ_CREDIT_BUREAU_MON	AMT_REQ_CREDIT_BUREAU_WEEK	AMT_REQ_CREDIT_BUREAU_DAY	AMT_REQ_CREDIT_BUREAU_QRT	DEF_60_CNT_SOCIAL_CIRCLE	OB6_60_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	OB3_30_CNT_SOCIAL_CIRCLE	EXT_SOURCE_2	AMT_GOODS_PRICE	DAYS_PUBLISH	FLAG_MOBIL	FLAG_WORK_PHONE	FLAG_EMP_PHONE	DAYS_EMPLOYED	FLAG_CONT_MOBILE	FLAG_PHONE
count	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
mean	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
std	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
min	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000

```
In [14]: data['NAME_TYPE_SUITE'].value_counts()

Out[14]:
```

	NAME_TYPE_SUITE	count
Unaccompanied	148526	
Spouse	45119	
Spouse, partner	11370	
Children	3267	
Other	1770	
Other_A	866	
Group of people	271	
Name: NAME_TYPE_SUITE, dtype: int64		

```
In [15]: data['NAME_TYPE_SUITE'].fillna(data.NAME_TYPE_SUITE.mode(),inplace = True)

In [16]: (100% data.isnull().sum()//len(data)).round(2).sort_values(ascending=False).head(20)

Out[16]:
```

	NAME_TYPE_SUITE	AMT_REQ_CREDIT_BUREAU_YEAR	AMT_REQ_CREDIT_BUREAU_MON	AMT_REQ_CREDIT_BUREAU_WEEK	AMT_REQ_CREDIT_BUREAU_DAY	AMT_REQ_CREDIT_BUREAU_QRT	DEF_60_CNT_SOCIAL_CIRCLE	OB6_60_CNT_SOCIAL_CIRCLE	DEF_30_CNT_SOCIAL_CIRCLE	OB3_30_CNT_SOCIAL_CIRCLE	EXT_SOURCE_2	AMT_GOODS_PRICE	DAYS_PUBLISH	FLAG_MOBIL	FLAG_WORK_PHONE	FLAG_EMP_PHONE	DAYS_EMPLOYED	FLAG_CONT_MOBILE	FLAG_PHONE
count	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
mean	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
std	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
min	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
25%	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
50%	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
75%	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000
max	0.420000	13.500000	13.500000	13.500000	13.500000	13.500000	0.330000	0.330000	0.330000	0.330000	0.210000	0.090000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000	0.000000

```
In [17]: data.drop(['SK_ID_CURR'],axis = 1,inplace = True)

In [ ]:
```

```
In [18]: #outliers treatment for numerical features(using IQR AND z-score)(num_feat)
data.describe()

Out[18]:
```

	TARGET	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_POPULATION_RELATIVE
count	307511.000000	307511.000000	307511.000000	307511.000000	307511.000000	307511.000000	307
mean	0.060729	0.471052	108787.919297	599025.999706	27108.973610	538316.204329	
std	0.272419	0.722121	237123.146279	402490.769966	14493.737315	369288.822246	
min	0.000000	0.000000	25650.000000	450			



In [76]: num\_feat

Out [76]:

	AMT_ANNUITY	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_GOODS_PRICE	CNT_FAM_MEMBERS	DAYS_EMPLOYED	DAYS_ID_PUI
0	24700.500000	200000.000000	406562.500000	351000.000000	1.000000		637
1	36968.500000	200000.000000	1293502.500000	1084500.000000	2.000000		1188
2	6750.000000	67500.000000	135000.000000	135000.000000	1.000000		225
3	29686.500000	135000.000000	312882.500000	297000.000000	2.000000		3039
4	21865.500000	121500.000000	513000.000000	513000.000000	1.000000		3038
...	...	...	...	...	...	...	...
307506	27558.000000	157500.000000	254700.000000	225000.000000	1.000000		236
307507	12001.500000	72000.000000	269550.000000	225000.000000	1.000000		13399
307508	29979.000000	153000.000000	677664.000000	585000.000000	1.000000		7921
307509	20205.000000	171000.000000	370107.000000	319500.000000	2.000000		4786
307510	49117.500000	157500.000000	675000.000000	675000.000000	2.000000		1262

307511 rows x 8 columns

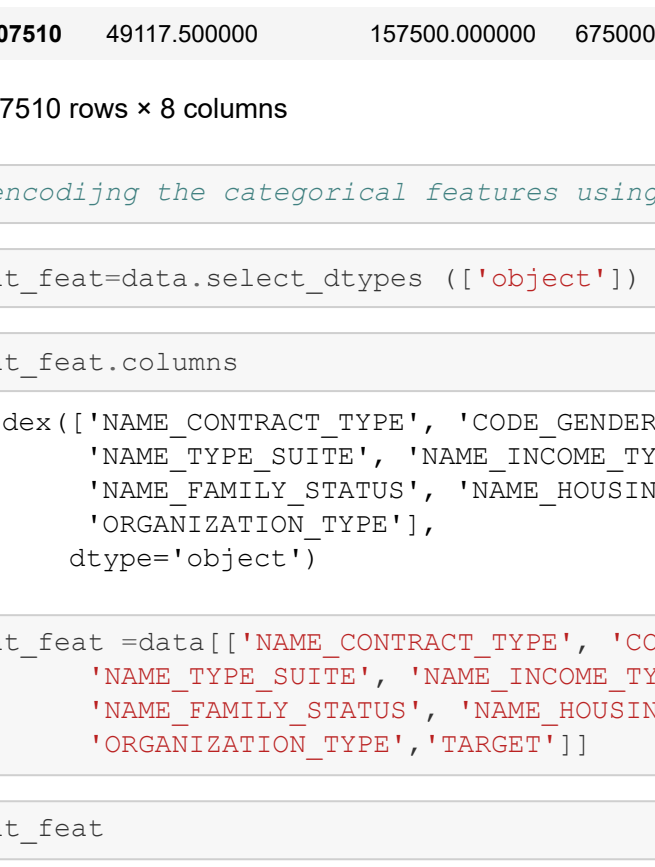
In [77]: num\_feat.describe()

Out [77]:

	AMT_ANNUITY	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_GOODS_PRICE	CNT_FAM_MEMBERS	DAYS_EMPLOYED	DAYS_ID_PUI
count	307511.000000	307511.000000	307511.000000	307511.000000	307509.000000	307511.000000	307510.0
mean	26940.161632	145831.230382	962861.869218	611901.874523	2.043765	4363.448085	2994.2
std	13689.567715	46518.877135	357882.236012	303278.121247	0.695068	4724.082848	1509.4
min	1615.500000	25650.000000	45000.000000	40500.000000	1.000000	0.000000	0.0
25%	16524.000000	112500.000000	270000.000000	238500.000000	2.000000	933.000000	1720.0
50%	24963.000000	147150.000000	513631.000000	450000.000000	2.000000	2219.000000	3254.0
75%	34596.000000	200000.000000	806850.000000	679500.000000	3.000000	5707.000000	4296.0
max	70588.000000	200000.000000	1345650.000000	1084500.000000	3.000000	13389.000000	7197.0

In [78]: sns.boxplot(num\_feat["AMT\_INCOME\_TOTAL"])

Out [78]: <matplotlib.axes.\_subplots.AxesSubplot at 0x20e94120ca0>



In [79]: num\_feat.drop(num\_feat["AMT\_INCOME\_TOTAL"],idxmax(),inplace = True)

In [80]: num\_feat

Out [80]:

	AMT_ANNUITY	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_GOODS_PRICE	CNT_FAM_MEMBERS	DAYS_EMPLOYED	DAYS_ID_PUI
1	36968.500000	200000.000000	1293502.500000	1084500.000000	2.000000		1188
2	6750.000000	67500.000000	135000.000000	135000.000000	1.000000		225
3	29686.500000	135000.000000	312882.500000	297000.000000	2.000000		3039
4	21865.500000	121500.000000	513000.000000	513000.000000	1.000000		3038
5	27517.500000	99000.000000	490495.000000	454500.000000	2.000000		1588
...	...	...	...	...	...	...	...
307506	27558.000000	157500.000000	254700.000000	225000.000000	1.000000		236
307507	12001.500000	72000.000000	269550.000000	225000.000000	1.000000		13399
307508	29979.000000	153000.000000	677664.000000	585000.000000	1.000000		7921
307509	20205.000000	171000.000000	370107.000000	319500.000000	2.000000		4786
307510	49117.500000	157500.000000	675000.000000	675000.000000	2.000000		1262

307510 rows x 8 columns

In [81]: #encoding the categorical features using (label encoding)(cat\_feat)

In [82]: cat\_feat=data.select\_dtypes(['object'])

In [83]: cat\_feat.columns

Out [83]: Index(['NAME\_CONTRACT\_TYPE', 'CODE\_GENDER', 'FLAG\_OWN\_CAR', 'FLAG\_OWN\_REALTY', 'NAME\_TYPE\_SUITE', 'NAME\_INCOME\_TYPE', 'NAME\_EDUCATION\_TYPE', 'NAME\_FAMILY\_STATUS', 'NAME\_HOUSING\_TYPE', 'WEEKDAY\_APPR\_PROCESS\_START', 'ORGANIZATION\_TYPE', 'TARGET'], dtype='object')

In [84]: cat\_feat.to\_data(['NAME\_CONTRACT\_TYPE', 'CODE\_GENDER', 'FLAG\_OWN\_CAR', 'FLAG\_OWN\_REALTY', 'NAME\_TYPE\_SUITE', 'NAME\_INCOME\_TYPE', 'NAME\_EDUCATION\_TYPE', 'NAME\_FAMILY\_STATUS', 'NAME\_HOUSING\_TYPE', 'WEEKDAY\_APPR\_PROCESS\_START', 'ORGANIZATION\_TYPE', 'TARGET'])

In [85]: cat\_feat

Out [85]:

	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	NAME_TYPE_SUITE	NAME_INCOME_TYPE	NA
0	Cash loans	M	N	N	Unacompained	Working	
1	Cash loans	F	N	N	Family	State servant	
2	Revolving loans	M	Y	Y	Unacompained	Working	
3	Cash loans	F	N	Y	Unacompained	Working	
4	Cash loans	M	N	Y	Unacompained	Working	
...	...	...	...	...	...	...	...
307506	Cash loans	M	N	N	Unacompained	Working	
307507	Cash loans	F	N	Y	Unacompained	Pensioner	
307508	Cash loans	F	N	Y	Unacompained	Working	
307509	Cash loans	F	N	Y	Unacompained	Commercial associate	
307510	Cash loans	F	N	N	Unacompained	Commercial associate	

307511 rows x 12 columns

In [86]: #cat\_feat.drop(["TARGET"],axis=1,inplace = True)

In [87]: cat\_feat["NAME\_TYPE\_SUITE"].isnull().sum()

Out [87]: 1292

In [88]: cat\_feat["NAME\_TYPE\_SUITE"].fillna(cat\_feat["NAME\_TYPE\_SUITE"].value\_counts().index[0],inplace = True)

In [89]: from sklearn.preprocessing import LabelEncoder  
lb = LabelEncoder()  
cat\_feat = cat\_feat.apply(LabelEncoder().fit\_transform)

In [90]: cat\_feat

Out [90]:

	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	NAME_TYPE_SUITE	NAME_INCOME_TYPE	NA
0	0	1	0	1	6	7	
1	0	0	0	0	1	4	
2	1	1	1	1	6	7	
3	0	0	0	1	6	7	
4	0	1	0	1	6	7	
...	...	...	...	...	...	...	...
307506	0	1	0	0	6	7	
307507	0	0	0	1	6	7	
307508	0	0	0	1	6	7	
307509	0	0	0	1	6	1	
307510	0	0	0	0	6	1	

307511 rows x 12 columns

In [91]: num\_feat

Out [91]:

	AMT_ANNUITY	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_GOODS_PRICE	CNT_FAM_MEMBERS	DAYS_EMPLOYED	DAYS_ID_PUI
1	36968.500000	200000.000000	1293502.500000	1084500.000000	2.000000		1188
2	6750.000000	67500.000000	135000.000000	135000.000000	1.000000		225
3	29686.500000	135000.000000	312882.500000	297000.000000	2.000000		3039
4	21865.500000	121500.000000	513000.000000	513000.000000	1.000000		3038
5	27517.500000	99000.000000	490495.000000	454500.000000	2.000000		1588
...	...	...	...	...	...	...	...
307506	27558.000000	157500.000000	254700.000000	225000.000000	1.000000		236
307507	12001.500000	72000.000000	269550.000000	225000.000000	1.000000		13399
307508	29979.000000	153000.000000	677664.000000	585000.000000	1.000000		7921
307509	20205.000000	171000.000000	370107.000000	319500.000000	2.000000		4786
307510	49117.500000	157500.000000	675000.000000	675000.000000	2.000000		1262

307510 rows x 8 columns

In [92]: cat\_feat

Out [92]:

	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	NAME_TYPE_SUITE	NAME_INCOME_TYPE	NA
0	0	1	0	1	6	7	
1	0	0	0	0	1	4	
2	1	1	1	1	6	7	
3	0	0	0	1	6	7	
4	0	1	0	1	6	7	
...	...	...	...	...	...	...	...
307506	0	1	0	0	6	7	
307507	0	0	0	1	6	3	
307508	0	0	0	1	6	7	
307509	0	0	0	1	6	1	
307510	0	0	0	0	6	1	

307511 rows x 12 columns

In [93]: final\_data = pd.concat([num\_feat, cat\_feat], axis=1, join='inner')

In [94]: final\_data

Out [94]:

	AMT_ANNUITY	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_GOODS_PRICE	CNT_FAM_MEMBERS	DAYS_EMPLOYED	DAYS_ID_PUI
1	36968.500000	200000.000000	1293502.500000	1084500.000000	2.000000		1188
2	6750.000000	67500.000000	135000.000000	135000.000000	1.000000		225
3	29686.500000	135000.000000	312882.500000	297000.000000	2.000000		3039
4	21865.500000	121500.000000	513000.000000	513000.000000	1.000000		3038
5	27517.500000	99000.000000	490495.000000	454500.000000	2.000000		1588
...	...	...	...	...	...	...	...
307506	27558.000000	157500.000000	254700.000000	225000.000000	1.000000		236
307507	12001.500000	72000.000000	269550.000000	225000.000000	1.000000		13399
307508	29979.000000	153000.000000	677664.000000	585000.000000	1.000000		7921
307509	20205.000000	171000.000000	370107.000000	319500.000000	2.000000		4786
307510	49117.500000	157500.000000	675000.000000	675000.000000	2.000000		1262

307510 rows x 20 columns

In [95]: final\_data.isnull().sum()

Out [95]:

AMT_ANNUITY	0
AMT_INCOME_TOTAL	0
AMT_CREDIT	0
AMT_GOODS_PRICE	0
CNT_FAM_MEMBERS	2
DAYS_EMPLOYED	0
DAYS_ID_PUBLISH	0
DAYS_REGISTRATION	0
NAME_CONTRACT_TYPE	0
CODE_GENDER	0
FLAG_OWN_CAR	0
FLAG_OWN_REALTY	0
NAME_TYPE_SUITE	0
NAME_INCOME_TYPE	0
NAME_EDUCATION_TYPE	0
NAME_FAMILY_STATUS	0
NAME_HOUSING_TYPE	0
WEEKDAY_APPR_PROCESS_START	0
ORGANIZATION_TYPE	0
TARGET	0
dtype:	int64

In [96]: final\_data["CNT\_FAM\_MEMBERS"].fillna(final\_data["CNT\_FAM\_MEMBERS"].median(),inplace=True)

In [97]: X = final\_data.drop(["TARGET"],axis=1)  
y = final\_data["TARGET"]

In [98]: from sklearn.model\_selection import train\_test\_split

In [ ]:

In [99]: X\_train, X\_test, y\_train, y\_test = train\_test\_split(X,y,train\_size=0.33, random\_state=42)

In [100]: #X\_train = final\_data

In [101]: X\_train.shape

Out [101]: (260361, 19)

In [102]: from sklearn.ensemble import RandomForestClassifier  
model= RandomForestClassifier()  
model.fit(X\_train,y\_train)  
y\_pred = model.predict(X\_test)  
from sklearn.metrics import accuracy\_score  
score=accuracy\_score(y\_test,y\_pred)  
print("Your Model Accuracy is", score)  
Your Model Accuracy is 0.9186728288611437

In [103]: from sklearn.metrics import accuracy\_score, f1\_score, precision\_score, recall\_score, roc\_auc\_score, confusion\_matrix, classification\_report  
RF\_roc=roc\_auc\_score(y\_test, y\_pred)  
RF\_acc = accuracy\_score(y\_test, y\_pred)  
RF\_prec = precision\_score(y\_test, y\_pred)  
RF\_rec = recall\_score(y\_test, y\_pred)  
RF\_f1 = f1\_score(y\_test, y\_pred)  
RF\_cls=classification\_report(y\_test,y\_pred)

In [104]: print("precision score is ", RF\_prec)  
print("roc\_auc\_score is ", RF\_roc)  
print("recall score is ", RF\_rec)  
print("f1\_score is ", RF\_f1)  
precision\_score is 0.3333333333333333  
roc\_auc\_score is 0.5000498648655191  
recall\_score is 0.0001211827435773146  
f1\_score is 0.00024422740763173837

In [105]: print('classification report is', RF\_cls)  
classification\_report is  

	precision	recall	f1-score	support
0	0.92	1.00	0.96	93227
1	0.33	0.00	0.00	8252
accuracy			0.92	101479
macro avg	0.63	0.50	0.48	101479
weighted avg	0.87	0.92	0.88	101479

In [106]: import collections  
from collections import Counter  
from imblearn.over\_sampling import RandomOverSampler  
over\_sample = RandomOverSampler(0.65)  
X\_train\_sm, y\_train\_sm = over\_sample.fit\_sample(X\_train,y\_train)  
X\_test\_sm, y\_test\_sm = over\_sample.fit\_sample(X\_test,y\_test)  
print("The number of Target befor fit (1)",format(Counter(y\_train)))  
print("The number of Target after fit (1)",format(Counter(y\_train\_sm)))  
The number of Target befor fit Counter({0: 282686, 1: 24824})  
The number of Target after fit Counter({0: 189459, 1: 161040})

In [107]: from imblearn.combine import SMOTETomek  
from imblearn.over\_sampling import SMOTE  
import collections  
from collections import Counter  
counter =Counter(y\_train)  
print('Before', counter)  
smt = SMOTE()  
X\_train\_sm, y\_train\_sm = smt.fit\_sample(X\_train,y\_train)  
X\_test\_sm, y\_test\_sm = smt.fit\_sample(X\_test,y\_test)  
counter1 = Counter(y\_train\_sm)  
print('after',counter1)  
Before Counter({0: 189459, 1: 16572})  
after Counter({0: 189459, 1: 189459})

In [108]: model= RandomForestClassifier(n\_estimators= 42,criterion='gini')  
new =model.fit(X\_train\_sm,y\_train\_sm)  
y\_pred1 = model.predict(X\_test\_sm)  
from sklearn.metrics import accuracy\_score  
score=accuracy\_score(y\_test\_sm,y\_pred1)  
print("Your Model Accuracy is", score)  
Your Model Accuracy is 0.8774764821350038

In [109]: RF\_cls1=classification\_report(y\_test\_sm,y\_pred1)

In [110]: print(RF\_cls1)  

	precision	recall	f1-score	support
0	0.82	0.97	0.89	93227
1	0.91	0.78	0.86	93227
accuracy			0.88	186454
macro avg	0.89	0.88	0.88	186454
weighted avg	0.89	0.88	0.88	186454

In [111]: print(confusion\_matrix(y\_pred1,y\_test\_sm))  
[[90784 20402]  
 [ 2443 72825]]

In [112]: from sklearn.model\_selection import StratifiedKFold  
stfold=StratifiedKFold(n\_splits=4)  
model = RandomForestClassifier()  
from sklearn.model\_selection import cross\_val\_score  
score =cross\_val\_score(model,X\_train\_sm,y\_train\_sm,cv =stfold)

In [113]: score  
array([0.82094373, 0.96777156, 0.96786623, 0.96689504])

Out [114]: RF\_roc2=roc\_auc\_score(y\_test\_sm, y\_pred1)  
RF\_acc2= accuracy\_score(y\_test\_sm, y\_pred1)  
RF\_prec2 = precision\_score(y\_test\_sm, y\_pred1)  
RF\_rec2 = recall\_score(y\_test\_sm, y\_pred1)  
RF\_f12 = f1\_score(y\_test\_sm, y\_pred1)  
RF\_cls2=classification\_report(y\_test\_sm,y\_pred1)

In [115]: print("precision score is ", RF\_prec2)  
print("roc\_auc\_score is ", RF\_roc2)  
print("recall score is ", RF\_rec2)  
print("f1\_score is ", RF\_f12)  
precision\_score is 0.9675426476058983  
roc\_auc\_score is 0.877476482135004  
recall\_score is 0.7811578190867453  
f1\_score is 0.8644173417608831

In [116]: #decision tree.  
dtree = DecisionTreeClassifier()  
from sklearn.tree import DecisionTreeClassifier

In [117]: tree = dtree.fit(X\_train\_sm,y\_train\_sm)  
y\_pred = dtree.predict(X\_test\_sm)  
print("Classification report - Vn", classification\_report(y\_test\_sm,y\_pred1))  
print(confusion\_matrix(y\_pred,y\_test\_sm))  
Classification report -  

	precision	recall	f1-score	support
0	0.81	0.85	0.83	93227
1	0.85	0.80	0.82	93227
accuracy			0.83	186454
macro avg	0.83	0.83	0.83	186454
weighted avg	0.83	0.83	0.83	186454

  
[[79593 18685]  
 [13634 74542]]

In [118]: #adaboosting  
from sklearn.ensemble import AdaBoostClassifier

In [119]: ab = AdaBoostClassifier(n\_estimators=1000, random\_state=0)

In [120]: ab.fit(X\_train\_sm,y\_train\_sm)  
y\_pred2=ab.predict(X\_test\_sm)  
print("Classification report - Vn", classification\_report(y\_test\_sm,y\_pred2))  
print(confusion\_matrix(y\_pred2,y\_test\_sm))  
Classification report -  

	precision	recall	f1-score	support
0	0.85	0.95	0.90	93227
1	0.95	0.83	0.89	93227
accuracy			0.89	186454
macro avg	0.90	0.89	0.89	186454
weighted avg	0.90	0.89	0.89	186454

  
[[98881 15552]  
 [ 4346 77675]]

In [121]: #gradient boosting  
from sklearn.ensemble import GradientBoostingClassifier

In [122]: GB = GradientBoostingClassifier()

In [123]: GB.fit(X\_train\_sm,y\_train\_sm)  
y\_pred3=GB.predict(X\_test\_sm)  
print("Classification report - Vn", classification\_report(y\_test\_sm,y\_pred3))  
print(confusion\_matrix(y\_pred3,y\_test\_sm))  
Classification report -  

	precision	recall	f1-score	support
0	0.83	0.95	0.89	93227
1	0.94	0.81		