

Importing Libraries

```
In [1]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
```

```
In [2]: # ignore warnings
warnings.filterwarnings('ignore')
```

1. Reading data

```
In [3]: appl_data = pd.read_csv('application_data.csv')
prev_data = pd.read_csv('previous_application.csv')
```

Inspect the DataFrame

```
In [4]: appl_data.head()
```

Out [4]:

	SK_ID_CURR	TARGET	NAME_CONTRACT_TYPE	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY	CNT_CHILDREN	AMT_INCOME
0	100002	1	Cash loans	M	N	N	Y	0
1	100003	0	Cash loans	F	N	N	N	0
2	100004	0	Revolving loans	M	Y	Y	Y	0
3	100005	0	Cash loans	F	N	Y	Y	0
4	100007	0	Cash loans	M	N	Y	Y	0

5 rows x 122 columns

```
In [5]: appl_data.shape
```

```
Out [5]: (307511, 122)
```

```
In [6]: appl_data.info(verbose=True, null_counts=True)
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 307511 entries, 0 to 307510  
Data columns (total 122 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	134133 non-null	float64
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	float64
19	DAYS_REGISTRATION	307511 non-null	float64
20	DAYS_ID_PUBLISH	307511 non-null	int64
21	OWN_CAR_AGE	154890 non-null	float64
22	FLAG_MOBIL	307511 non-null	int64
23	FLAG_EMP_PHONE	307511 non-null	int64
24	FLAG_WORK_PHONE	307511 non-null	object
25	FLAG_CONT_MOBILE	307511 non-null	int64
26	FLAG_PHONE	307511 non-null	int64
27	FLAG_EMAIL	307511 non-null	int64
28	OCCUPATION_TYPE	115161 non-null	object
29	CNT_FAM_MEMBERS	307509 non-null	float64
30	REGION_RATING_CLIENT	307511 non-null	int64
31	REGION_RATING_CLIENT_W_CITY	307511 non-null	object
32	WEEKDAY_APPR_PROCESS_START	307511 non-null	object
33	HOUR_APPR_PROCESS_START	307511 non-null	int64
34	REG_REGION_NOT_LIVE_REGION	307511 non-null	int64
35	REG_REGION_NOT_WORK_REGION	307511 non-null	int64
36	LIVE_REGION_NOT_WORK_REGION	307511 non-null	int64
37	REG_CITY_NOT_LIVE_CITY	307511 non-null	int64
38	REG_CITY_NOT_WORK_CITY	307511 non-null	float64
39	LIVE_CITY_NOT_WORK_CITY	307511 non-null	int64
40	ORGANIZATION_TYPE	307511 non-null	object
41	EXT_SOURCE_1	134133 non-null	float64
42	EXT_SOURCE_2	306851 non-null	float64
43	EXT_SOURCE_3	246346 non-null	float64
44	APARTMENTS_AVG	151530 non-null	float64
45	BASEMENTAREA_AVG	127568 non-null	float64
46	YEARS_BEGINEXPLOITATION_AVG	157504 non-null	float64
47	YEARS_BUILT_AVG	103023 non-null	float64
48	COMMONAREA_AVG	131120 non-null	float64
49	ELEVATORS_AVG	143620 non-null	float64
50	ENTRANCES_AVG	152683 non-null	float64
51	FLOORSINMAX_MODE	154890 non-null	float64
52	FLOORSINMIN_MODE	98869 non-null	float64
53	LANDAREA_AVG	124921 non-null	float64
54	LIVINGAPARTMENTS_AVG	97312 non-null	float64
55	LIVINGAREA_AVG	153161 non-null	float64
56	NONLIVINGAPARTMENTS_AVG	93997 non-null	float64
57	NONLIVINGAREA_AVG	137829 non-null	float64
58	APARTMENTS_MEDI	151530 non-null	float64
59	BASEMENTAREA_MEDI	127568 non-null	float64
60	YEARS_BEGINEXPLOITATION_MEDI	157504 non-null	float64
61	YEARS_BUILT_MEDI	103023 non-null	float64
62	COMMONAREA_MEDI	131120 non-null	float64
63	ELEVATORS_MEDI	143620 non-null	float64
64	ENTRANCES_MEDI	152683 non-null	float64
65	FLOORSINMAX_MEDI	154890 non-null	float64
66	FLOORSINMIN_MEDI	98869 non-null	float64
67	LANDAREA_MEDI	124921 non-null	float64
68	LIVINGAPARTMENTS_MEDI	97312 non-null	float64
69	LIVINGAREA_MEDI	153161 non-null	float64
70	NONLIVINGAPARTMENTS_MEDI	93997 non-null	float64
71	NONLIVINGAREA_MEDI	137829 non-null	float64
72	APARTMENTS_MODE	151530 non-null	float64
73	BASEMENTAREA_MEDI	127568 non-null	float64
74	YEARS_BEGINEXPLOITATION_MEDI	157504 non-null	float64
75	YEARS_BUILT_MEDI	103023 non-null	float64
76	COMMONAREA_MEDI	131120 non-null	float64
77	ELEVATORS_MEDI	143620 non-null	float64
78	ENTRANCES_MEDI	152683 non-null	float64
79	FLOORSINMAX_MEDI	154890 non-null	float64
80	FLOORSINMIN_MEDI	98869 non-null	float64
81	LANDAREA_MEDI	124921 non-null	float64
82	LIVINGAPARTMENTS_MEDI	97312 non-null	float64
83	LIVINGAREA_MEDI	153161 non-null	float64
84	NONLIVINGAPARTMENTS_MEDI	93997 non-null	float64
85	NONLIVINGAREA_MEDI	137829 non-null	float64
86	FONDKAPREMONT_MODE	97216 non-null	object
87	HOUSHTYPE_MODE	153214 non-null	object
88	WALLSMATERIAL_MODE	159800 non-null	float64
89	EMERGENCYSTATE_MODE	161756 non-null	object
90	EMERGENCYSTATE_MODE	161756 non-null	object
91	OBS_30_CNT_SOCIAL_CIRCLE	306490 non-null	float64
92	DEF_30_CNT_SOCIAL_CIRCLE	306490 non-null	float64
93	OBS_60_CNT_SOCIAL_CIRCLE	306490 non-null	float64
94	DEF_60_CNT_SOCIAL_CIRCLE	306490 non-null	float64
95	DAYS_LAST_PHONE_CHANGE	307511 non-null	int64
96	FLAG_DOCUMENT_2	307511 non-null	int64
97	FLAG_DOCUMENT_3	307511 non-null	int64
98	FLAG_DOCUMENT_4	307511 non-null	int64
99	FLAG_DOCUMENT_5	307511 non-null	int64
100	FLAG_DOCUMENT_6	307511 non-null	int64
101	FLAG_DOCUMENT_7	307511 non-null	int64
102	FLAG_DOCUMENT_8	307511 non-null	int64
103	FLAG_DOCUMENT_9	307511 non-null	int64
104	FLAG_DOCUMENT_10	307511 non-null	int64
105	FLAG_DOCUMENT_11	307511 non-null	int64
106	FLAG_DOCUMENT_12	307511 non-null	int64
107	FLAG_DOCUMENT_13	307511 non-null	int64
108	FLAG_DOCUMENT_14	307511 non-null	int64
109	FLAG_DOCUMENT_15	307511 non-null	int64
110	FLAG_DOCUMENT_16	307511 non-null	int64
111	FLAG_DOCUMENT_17	307511 non-null	int64
112	FLAG_DOCUMENT_18	307511 non-null	int64
113	FLAG_DOCUMENT_19	307511 non-null	int64
114	FLAG_DOCUMENT_20	307511 non-null	int64
115	FLAG_DOCUMENT_21	307511 non-null	int64
116	AMT_REQ_CREDIT_BUREAU_HOUR	265992 non-null	float64
117	AMT_REQ_CREDIT_BUREAU_DAY	265992 non-null	float64
118	AMT_REQ_CREDIT_BUREAU_WEEK	265992 non-null	float64
119	AMT_REQ_CREDIT_BUREAU_MON	265992 non-null	float64
120	AMT_REQ_CREDIT_BUREAU_QRT	265992 non-null	float64
121	AMT_REQ_CREDIT_BUREAU_YEAR	265992 non-null	float64

```
In [7]: pd.set_option('display.max_rows', appl_data.shape[0])
appl_data.describe().T
```

CNT\_FAM\_MEMBERS

307509.0

2.152865

0.910682

1.000000e+00

2.000000

2.000000

3.000000

REGION\_RATING\_CLIENT

307510.0

2.052463

0.509034

1.000000e+00

2.000000

2.000000

2.000000

REGION\_RATING\_CLIENT\_W\_CITY

307510.0

2.031521

0.502737

1.000000e+00

2.000000

2.000000

2.000000

HOUR\_APPR\_PROCESS\_START

307510.0

12.083419

3.265832

0.000000e+00

10.000000

12.000000

14.000000

REG\_REGION\_NOT\_LIVE\_REGION

307510.0

0.015144

0.121206

0.000000e+00

0.000000

0.000000

0.000000

REG\_REGION\_NOT\_WORK\_REGION

307510.0

0.005769

0.219526

0.000000e+00

0.000000

0.000000

0.000000

LIVE\_REGION\_NOT\_WORK\_REGION

307510.0

0.040659

0.197499

0.000000e+00

0.000000

0.000000

0.000000

REG\_CITY\_NOT\_LIVE\_CITY

307510.0

0.078173

0.268444

0.000000e+00

0.000000

0.000000

0.000000

REG\_CITY\_NOT\_WORK\_CITY

307510.0

0.230454

0.421124

0.000000e+00

0.000000

0.000000

0.000000

LIVE\_CITY\_NOT\_WORK\_CITY

307510.0

0.179555

0.383817

0.000000e+00

0.000000

0.000000

0.000000

EXT\_SOURCE\_1

134133.0

0.502130

0.211062

1.456813e-02

0.334057

0.505998

0.675053

EXT\_SOURCE\_2

306851.0

0.514393

0.191080

0.1873617e-08

0.392457

0.565961

0.663617

EXT\_SOURCE\_3

246346.0

0.510583

0.194844

5.2726526e-04

0.370505

0.535276

0.660051

APARTMENTS\_AVG

151545.0

0.171440

0.108240

0.000000e+00

0.057700

0.087600

0.148500

BASEMENTAREA\_AVG

127568.0

0.088442

0.084238

0.000000e+00

0.044200

0.076300

0.112200

YEARS\_BEGINEXPLOITATION\_AVG

157504.0

0.977735

0.059223

0.000000e+00

0.976700

0.981600

0.989600

YEARS\_BUILT\_AVG

103023.0

0.752471

0.113280

0.000000e+00

0.687250

0.755200

0.823200

COMMONAREA\_AVG

126486.0

0.044621

0.076036

0.000000e+00

0.007800

0.021100

0.051500

ELEVATORS\_AVG

143620.0

0.078942

0.134576

0.000000e+00

0.000000

0.000000

0.120000

ENTRANCES\_AVG

152683.0

0.149725

0.100049

0.000000e+00

0.068000

0.137900

0.206900

FLOORSINMAX\_AVG

154491.0

0.226822

0.144641

0.000000e+00

0.166700

0.166700

0.333300

FLOORSINMIN\_AVG

98869.0

0.231894

0.161380

0.000000e+00

0.083300

0.208300

0.375000

LANDAREA\_AVG

124921.0

0.066333

0.081184

0.000000e+00

0.016700

0.048100

0.058600

LIVINGAPARTMENTS\_AVG

97312.0

0.107375

0.082576

0.000000e+00

0.050400

0.075600

0.121000

LIVINGAREA\_AVG

153161.0

0.107399

0.110565

0.000000e+00

0.045300

0.074500

0.129900

NONLIVINGAPARTMENTS\_AVG

93997.0

0.080809

0.047732

0.000000e+00

0.000000

0.000000

0.003000

NONLIVINGAREA\_AVG

137829.0

0.093578

0.069523

0.000000e+00

0.000000

0.003600

0.027700

APARTMENTS\_MEDI

151545.0

0.114231

0.107936

0.000000e+00

0.052500

0.084000

0.143900

BASEMENTAREA\_MEDI

127568.0

0.087343

0.084307

0.000000e+00

0.040700

0.074600

0.112400

YEARS\_BEGINEXPLOITATION\_MEDI

157504.0

0.977065

0.064575

0.000000e+00

0.976700

0.981600

0.989600

YEARS\_BUILT\_MEDI

103023.0

0.757546

0.112066

0.000000e+00

0.691400

0.758500

0.825600

COMMONAREA\_MEDI

126486.0

0.044595

0.076147

0.000000e+00

0.007900

0.020800

0.051300

ELEVATORS\_MEDI

143620.0

0.078078

0.134661

0.000000e+00

0.000000

0.000000

0.120000

ENTRANCES\_MEDI

152683.0

0.149213

0.100366

0.000000e+00

0.068000

0.137900

0.206900

FLOORSINMAX\_MEDI

154491.0

0.225957

0.145057

0.000000e+00

0.166700

0.166700

0.333300

FLOORSINMIN\_MEDI

98869.0

0.231625

0.161934

0.000000e+00

0.083300

0.208300

0.375000

LANDAREA\_MEDI

124921.0

0.067169

0.082167

0.000000e+00

0.016700

0.048700

0.068600

LIVINGAPARTMENTS\_MEDI

97312.0

0.101954

0.093642

0.000000e+00

0.051300

0.076100

0.123100

LIVINGAREA\_MEDI

153161.0

0.108607

0.112260

0.000000e+00

0.045700

0.074900

0.130300

NONLIVINGAPARTMENTS\_MEDI

93997.0

0.080651

0.047415

0.000000e+00

0.000000

0.000000

0.003900

NONLIVINGAREA\_MEDI

137829.0

0.082836

0.070166

0.000000e+00

0.000000

0.003100

0.026600

TOTALAREA\_MEDI

159080.0

0.102547

0.107462

0.000000e+00

0.041200

0.068800

0.127800

OBS\_30\_CNT\_SOCIAL\_CIRCLE

306490.0

1.422245

2.400989

0.000000e+00

0.000000

0.000000

2.000000

DEF\_30\_CNT\_SOCIAL\_CIRCLE

306490.0

1.143421

2.446988

0.000000e+00

0.000000

0.000000

0.000000

OBS\_60\_CNT\_SOCIAL\_CIRCLE

306490.0

1.405292

2.378903

0.000000e+00

0.000000

0.000000

2.000000

DEF\_60\_CNT\_SOCIAL\_CIRCLE

306490.0

1.100049

3.862291

0.000000e+00

0.000000

0.000000

0.000000

DAYS\_LAST\_PHONE\_CHANGE

307510.0

-962.858788

826.808487

-4.292000e+03

-1570.000000

-757.000000

-274.000000

FLAG\_DOCUMENT\_1

307510.0

0.000042

0.006502

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_2

307510.0

0.017023

0.045352

0.000000e+00

0.000000

1.000000

1.000000

FLAG\_DOCUMENT\_3

307510.0

0.000081

0.009016

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_4

307510.0

0.015115

0.122010

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_5

307510.0

0.008055

0.283376

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_6

307510.0

0.000192

0.013850

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_7

307510.0

0.081376

0.273412

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_8

307510.0

0.003896

0.062295

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_9

307510.0

0.000023

0.004771

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_10

307510.0

0.003912

0.062424

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_11

307510.0

0.000007

0.002550

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_12

307510.0

0.003525

0.059268

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_13

307510.0

0.002936

0.054110

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_14

307510.0

0.001210

0.034760

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_15

307510.0

0.000928

0.099144

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_16

307510.0

0.000267

0.016327

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_17

307510.0

0.008130

0.089796

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_18

307510.0

0.000595

0.024387

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_19

307510.0

0.000507

0.022518

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_20

307510.0

0.000335

0.018299

0.000000e+00

0.000000

0.000000

0.000000

FLAG\_DOCUMENT\_21

307510.0

0.000035

0.003849

0.000000e+00

0.000000

0.000000

0.000000

AMT\_REQ\_CREDIT\_BUREAU\_HOUR

265992.0

0.060402

0.063849

0.000000e+00

0.000000

0.000000

0.000000

AMT\_REQ\_CREDIT\_BUREAU\_DAY

265992.0

0.007000

0.110757

0.000000e+00

0.000000

0.000000

0.000000

AMT\_REQ\_CREDIT\_BUREAU\_WEEK

265992.0

0.034362

0.204685

0.000000e+00

0.000000

0.000000

0.000000

AMT\_REQ\_CREDIT\_BUREAU\_MON

265992.0

0.267395

0.916002

0.000000e+00

0.000000

0.000000

0.000000

AMT\_REQ\_CREDIT\_BUREAU\_QRT

265992.0

0.265474

0.794056

0.000000e+00

0.000000

0.000000

0.000000

AMT\_REQ\_CREDIT\_BUREAU\_YEAR

265992.0

1.899974

1.869295

0.000000e+00

0.000000

1.000000

3.000000

2. Data Analysis

2.1 Identifying Missing Value columns

In [8]:

```
perc_missing = round((app1_data.isna().sum() / app1_data.shape[0], 2)

#number of columns with no missing values
perc_missing[perc_missing == 0].size
```

Out[8]: 65

In [9]:

```
#number of columns with missing values
perc_missing[perc_missing != 0].size
```

Out[9]: 57

In [10]:

```
perc_missing
```

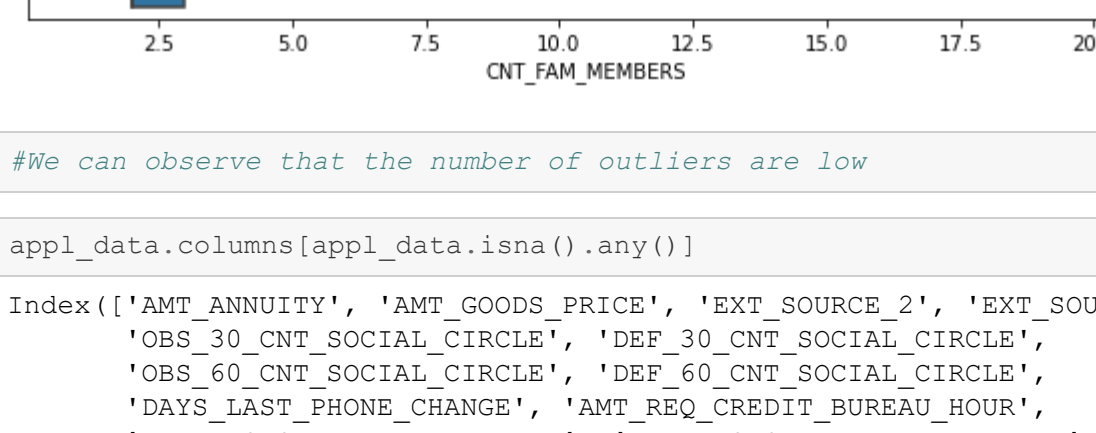
Out[10]:

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	0.00
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
OWN_CAR_AGE	0.66
FLAG_MOBIL	0.00
FLAG_EMP_PHONE	0.00
FLAG_WORK_PHONE	0.00
FLAG_CONT_MOBILE	0.00
FLAG_PHONE	0.00
FLAG_EMAIL	0.00
OCCUPATION_TYPE	0.31
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
EXT_SOURCE_1	0.56
EXT_SOURCE_2	0.00
EXT_SOURCE_3	0.20
APARTMENTS_AVG	0.51
BASEMENTAREA_AVG	0.51
YEARS_BEGINEXPLOITATION_AVG	0.49
YEARS_BUILT_AVG	0.66
COMMONAREA_AVG	0.55
ELEVATORS_AVG	0.53
ENTRANCES_AVG	0.50
FLOORSINMAX_AVG	0.50
FLOORSINMIN_AVG	0.68
LANDAREA_AVG	0.59
LIVINGAPARTMENTS_AVG	0.68
LIVINGAREA_AVG	0.51
NONLIVINGAPARTMENTS_AVG	0.69
NONLIVINGAREA_AVG	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70
ELEVATORS_MEDI	0.53
ENTRANCES_MEDI	0.50
FLOORSINMAX_MEDI	0.50
FLOORSINMIN_MEDI	0.68
LANDAREA_MEDI	0.59
LIVINGAPARTMENTS_MEDI	0.68
LIVINGAREA_MEDI	0.50
NONLIVINGAPARTMENTS_MEDI	0.69
NONLIVINGAREA_MEDI	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70
ELEVATORS_MEDI	0.53
ENTRANCES_MEDI	0.50
FLOORSINMAX_MEDI	0.50
FLOORSINMIN_MEDI	0.68
LANDAREA_MEDI	0.59
LIVINGAPARTMENTS_MEDI	0.68
LIVINGAREA_MEDI	0.50
NONLIVINGAPARTMENTS_MEDI	0.69
NONLIVINGAREA_MEDI	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70
ELEVATORS_MEDI	0.53
ENTRANCES_MEDI	0.50
FLOORSINMAX_MEDI	0.50
FLOORSINMIN_MEDI	0.68
LANDAREA_MEDI	0.59
LIVINGAPARTMENTS_MEDI	0.68
LIVINGAREA_MEDI	0.50
NONLIVINGAPARTMENTS_MEDI	0.69
NONLIVINGAREA_MEDI	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70
ELEVATORS_MEDI	0.53
ENTRANCES_MEDI	0.50
FLOORSINMAX_MEDI	0.50
FLOORSINMIN_MEDI	0.68
LANDAREA_MEDI	0.59
LIVINGAPARTMENTS_MEDI	0.68
LIVINGAREA_MEDI	0.50
NONLIVINGAPARTMENTS_MEDI	0.69
NONLIVINGAREA_MEDI	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70
ELEVATORS_MEDI	0.53
ENTRANCES_MEDI	0.50
FLOORSINMAX_MEDI	0.50
FLOORSINMIN_MEDI	0.68
LANDAREA_MEDI	0.59
LIVINGAPARTMENTS_MEDI	0.68
LIVINGAREA_MEDI	0.50
NONLIVINGAPARTMENTS_MEDI	0.69
NONLIVINGAREA_MEDI	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70
ELEVATORS_MEDI	0.53
ENTRANCES_MEDI	0.50
FLOORSINMAX_MEDI	0.50
FLOORSINMIN_MEDI	0.68
LANDAREA_MEDI	0.59
LIVINGAPARTMENTS_MEDI	0.68
LIVINGAREA_MEDI	0.50
NONLIVINGAPARTMENTS_MEDI	0.69
NONLIVINGAREA_MEDI	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70
ELEVATORS_MEDI	0.53
ENTRANCES_MEDI	0.50
FLOORSINMAX_MEDI	0.50
FLOORSINMIN_MEDI	0.68
LANDAREA_MEDI	0.59
LIVINGAPARTMENTS_MEDI	0.68
LIVINGAREA_MEDI	0.50
NONLIVINGAPARTMENTS_MEDI	0.69
NONLIVINGAREA_MEDI	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70
ELEVATORS_MEDI	0.53
ENTRANCES_MEDI	0.50
FLOORSINMAX_MEDI	0.50
FLOORSINMIN_MEDI	0.68
LANDAREA_MEDI	0.59
LIVINGAPARTMENTS_MEDI	0.68
LIVINGAREA_MEDI	0.50
NONLIVINGAPARTMENTS_MEDI	0.69
NONLIVINGAREA_MEDI	0.55
APARTMENTS_MEDI	0.51
BASEMENTAREA_MEDI	0.59
YEARS_BEGINEXPLOITATION_MEDI	0.49
YEARS_BUILT_MEDI	0.66
COMMONAREA_MEDI	0.70



[33]: plt.figure(figsize=(10, 1)) sns.boxplot(app1\_data['CNT\_FAM\_MEMBERS']).set\_title('CNT\_FAM\_MEMBERS') plt.show()

CNT\_FAM\_MEMBERS



In [36]: #We can observe that the number of outliers are low

In [37]: app1\_data.columns[app1\_data.isna().any()]

Out[37]: Index(['AMT\_ANNUITY', 'AMT\_GOODS\_PRICE', 'EXT\_SOURCE\_2', 'EXT\_SOURCE\_3', 'OBS\_30\_CNT\_SOCIAL\_CIRCLE', 'OBS\_60\_CNT\_SOCIAL\_CIRCLE', 'DEF\_30\_CNT\_SOCIAL\_CIRCLE', 'DEF\_60\_CNT\_SOCIAL\_CIRCLE', 'DAYS\_LAST\_PHONE\_CHANGE', '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\_REQ\_CREDIT\_BUREAU\_YEAR', dtype='object')]

Column: AMT\_ANNUITY

In [38]: no\_rows = app1\_data['AMT\_ANNUITY'].isna().sum() no\_rows

Out[38]: 12

We can note that the number of values here is very less.

Annuity being a major contributing factor based the business context from column\_description.csv data, imputing this value could bias the result. Considering we have large number of rows in comparison with this, we can drop these rows

In [39]: app1\_data = app1\_data.loc[~app1\_data['AMT\_ANNUITY'].isna()]

In [40]: no\_rows\_dropped = no\_rows no\_rows\_dropped

Out[40]: 12

In [41]: #imputed, so no null values anymore. app1\_data['AMT\_ANNUITY'].isna().sum()

Out[41]: 0

Column: AMT\_GOODS\_PRICE

is dependent on loan type being consumer goods. The loan type column relevant to this is NAME\_CONTRACT\_TYPE. If the NAME\_CONTRACT\_TYPE is Consumer Loan, then this is the price of the consumer item being purchased by the client using this loan.

In [42]: app1\_data['NAME\_CONTRACT\_TYPE'].value\_counts()

Out[42]: Cash loans 278220 Revolving loans 12979 Name: NAME\_CONTRACT\_TYPE, dtype: int64

We observe that we have no records with us which has the relevant NAME\_CONTRACT\_TYPE This implies this value cannot be imputed

In [43]: app1\_data['AMT\_GOODS\_PRICE'].isna().sum()

Out[43]: 278

In [44]: app1\_data.loc[app1\_data['AMT\_GOODS\_PRICE'].isna()]['NAME\_CONTRACT\_TYPE'].value\_counts()

Out[44]: Revolving loans 278 Name: NAME\_CONTRACT\_TYPE, dtype: int64

We observe that all the records with AMT\_GOODS\_PRICE as NaN are **Revolving loans** and hence they do not have an AMT\_GOODS\_PRICE associated with it. Thus, these values can be safely **imputed with 0** based on the business context.

In [45]: app1\_data.loc[app1\_data['AMT\_GOODS\_PRICE'].isna()]['AMT\_GOODS\_PRICE'] = 0.0

In [46]: app1\_data['AMT\_GOODS\_PRICE'].isna().sum()

Out[46]: 0

In [47]: app1\_data.columns[app1\_data.isna().any()]

Out[47]: Index(['EXT\_SOURCE\_2', 'EXT\_SOURCE\_3', 'OBS\_30\_CNT\_SOCIAL\_CIRCLE', 'OBS\_60\_CNT\_SOCIAL\_CIRCLE', 'DEF\_30\_CNT\_SOCIAL\_CIRCLE', 'DEF\_60\_CNT\_SOCIAL\_CIRCLE', 'DAYS\_LAST\_PHONE\_CHANGE', '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', dtype='object')]

Columns: EXT\_SOURCE\_2

Based on the information from the column\_description, EXT\_SOURCE\_2 is a "Normalized Score" that are sourced from an external data source. Considering this, EXT\_SOURCE\_2 could potentially be significant for business and hence imputing them might not be advisable.

In [48]: no\_rows = app1\_data['EXT\_SOURCE\_2'].isna().sum() no\_rows

Out[48]: 660

In [49]: app1\_data.shape

Out[49]: (307499, 53)

In [50]: #no\_rows\_dropped has the number of rows dropped already.

In [51]: no\_rows + no\_rows\_dropped

Out[51]: 672

Consider, the no\_rows\_dropped, together with the no\_rows with EXT\_SOURCE\_2 is NaN, is significantly less than the total no. of records, we can safely drop this values to avoid skewing the result.

In [52]: app1\_data = app1\_data.loc[~app1\_data['EXT\_SOURCE\_2'].isna()]

In [53]: app1\_data['EXT\_SOURCE\_2'].isna().sum()

Out[53]: 0

In [54]: no\_rows\_dropped = no\_rows\_dropped + no\_rows no\_rows\_dropped

Out[54]: 672

Columns: EXT\_SOURCE\_3

Based on the information from the column\_description, EXT\_SOURCE\_3 is a "Normalized Score" that are sourced from an external data source. Considering this, EXT\_SOURCE\_3 could potentially be significant for business and hence imputing them might not be advisable.

We have to analyse the data before decision

In [55]: no\_rows = app1\_data['EXT\_SOURCE\_3'].isna().sum() no\_rows

Out[55]: 60734

In [56]: round(no\_rows/app1\_data.shape[0], 3)

Out[56]: 0.198

This is a considerable percentage. Being a 'Score' from an external source, with no further business information, we must consider dropping these rows.

In [57]: mean = app1\_data['EXT\_SOURCE\_3'].mean() mean

Out[57]: 0.5107773158925729

In [58]: mode = app1\_data['EXT\_SOURCE\_3'].mode() mode

Out[58]: 0 0.7463 dtype: float64

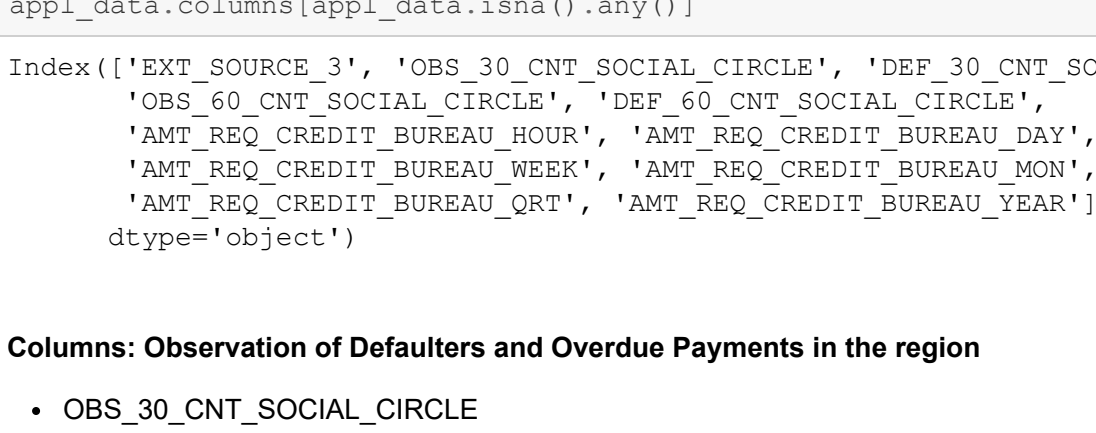
In [59]: median = round(app1\_data['EXT\_SOURCE\_3'].median(), 6) median

Out[59]: 0.535276

In [60]: es3\_data = app1\_data.loc[~app1\_data['EXT\_SOURCE\_3'].isna()]

In [61]: plt.figure(figsize=(10, 1)) sns.boxplot(app1\_data['EXT\_SOURCE\_3']).set\_title('EXT\_SOURCE\_3') plt.show()

EXT\_SOURCE\_3



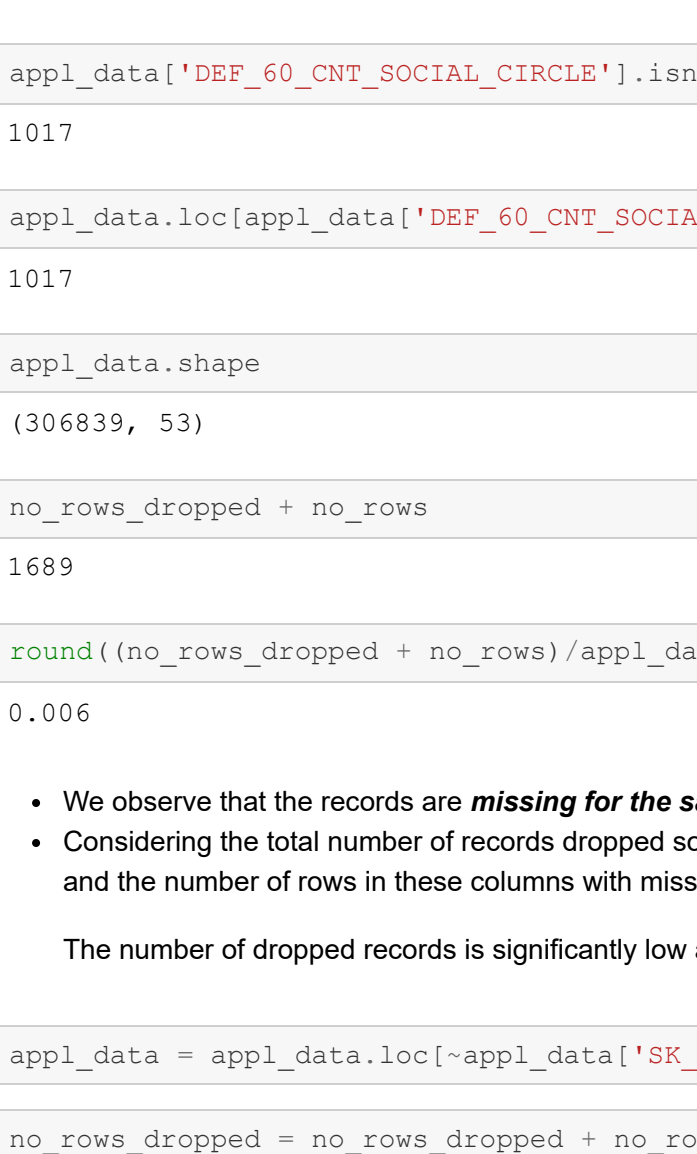
In [62]: es3\_data['EXT\_SOURCE\_3'].describe()

Out[62]: count 246105.000000 mean 0.510777 std 0.134854 min 0.000527 25% 0.370650 50% 0.535276 75% 0.689057 max 0.896010 Name: EXT\_SOURCE\_3, dtype: float64

The above implies the available score are fairly distributed without any outliers. We can also observe the percentage of missing values is higher than the standard deviation of the rest of the values.

In [63]: from scipy.stats import norm sns.distplot(app1\_data['EXT\_SOURCE\_3'], fit=norm, kde=False).set\_title('EXT\_SOURCE\_3') plt.show()

EXT\_SOURCE\_3



Considering these results, we will refrain from imputing this value When involving this variable for analysis, depending on the co-relation, we can analyse without the missing value rows for the related variables as well.

In [64]: app1\_data.columns[app1\_data.isna().any()]

Out[64]: Index(['EXT\_SOURCE\_3', 'OBS\_30\_CNT\_SOCIAL\_CIRCLE', 'OBS\_60\_CNT\_SOCIAL\_CIRCLE', 'DEF\_30\_CNT\_SOCIAL\_CIRCLE', 'DEF\_60\_CNT\_SOCIAL\_CIRCLE', 'DAYS\_LAST\_PHONE\_CHANGE', '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', dtype='object')]

Columns: Observation of Defaulters and Overdue Payments in the region

- OBS\_30\_CNT\_SOCIAL\_CIRCLE
- DEF\_30\_CNT\_SOCIAL\_CIRCLE
- OBS\_60\_CNT\_SOCIAL\_CIRCLE
- DEF\_60\_CNT\_SOCIAL\_CIRCLE

In [65]: no\_rows = app1\_data['OBS\_30\_CNT\_SOCIAL\_CIRCLE'].isna().sum() no\_rows

Out[65]: 1017

In [66]: SK\_ID\_CURRs = app1\_data.loc[app1\_data['OBS\_30\_CNT\_SOCIAL\_CIRCLE'].isna()]['SK\_ID\_CURR'] SK\_ID\_CURRs

Out[66]: 68 100080 394 100457 397 100460 457 100527 1042 101209 ... 305526 453990 305623 454093 305641 454116 305642 454117 307402 456135 Name: SK\_ID\_CURR, Length: 1017, dtype: int64

In [67]: app1\_data['DEF\_30\_CNT\_SOCIAL\_CIRCLE'].isna().sum()

Out[67]: 1017

In [68]: app1\_data.loc[app1\_data['DEF\_30\_CNT\_SOCIAL\_CIRCLE'].isna()]['SK\_ID\_CURR'].isin(SK\_ID\_CURRs).sum()

Out[68]: 1017

In [69]: app1\_data['OBS\_60\_CNT\_SOCIAL\_CIRCLE'].isna().sum()

Out[69]: 1017

In [70]: app1\_data.loc[app1\_data['OBS\_60\_CNT\_SOCIAL\_CIRCLE'].isna()]['SK\_ID\_CURR'].isin(SK\_ID\_CURRs).sum()

Out[70]: 1017

In [71]: app1\_data['DEF\_60\_CNT\_SOCIAL\_CIRCLE'].isna().sum()

Out[71]: 1017

In [72]: app1\_data.loc[app1\_data['DEF\_60\_CNT\_SOCIAL\_CIRCLE'].isna()]['SK\_ID\_CURR'].isin(SK\_ID\_CURRs).sum()

Out[72]: 1017

In [73]: app1\_data.shape

Out[73]: (306839, 53)

In [74]: no\_rows\_dropped + no\_rows

Out[74]: 1689

In [75]: round((no\_rows\_dropped + no\_rows)/app1\_data.shape[0],3)

Out[75]: 0.006

- We observe that the records are missing for the same set of applicants
- Considering the total number of records dropped earlier and the number of rows in these columns with missing values, The number of dropped records is significantly low and can be dropped.

In [76]: app1\_data = app1\_data.loc[~app1\_data['SK\_ID\_CURR'].isin(SK\_ID\_CURRs)]

In [77]: no\_rows\_dropped = no\_rows\_dropped + no\_rows no\_rows\_dropped

Out[77]: 1689

In [78]: app1\_data['OBS\_30\_CNT\_SOCIAL\_CIRCLE'].isna().sum()

Out[78]: 0

In [79]: app1\_data.columns[app1\_data.isna().any()]

Out[79]: Index(['EXT\_SOURCE\_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', dtype='object')]

Columns: Number of enquiries to Credit Bureau about the client

- 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 [80]: f, axes = plt.subplots(2, 3) plt.figure(figsize=(1, 1)) sns.histplot(data=app1\_data, x='AMT\_REQ\_CREDIT\_BUREAU\_HOUR', ax=axes[0][0]) s.set(klabel=None)

s = sns.histplot(data=app1\_data, x='AMT\_REQ\_CREDIT\_BUREAU\_DAY', ax=axes[0][1]) s.set(klabel=None)

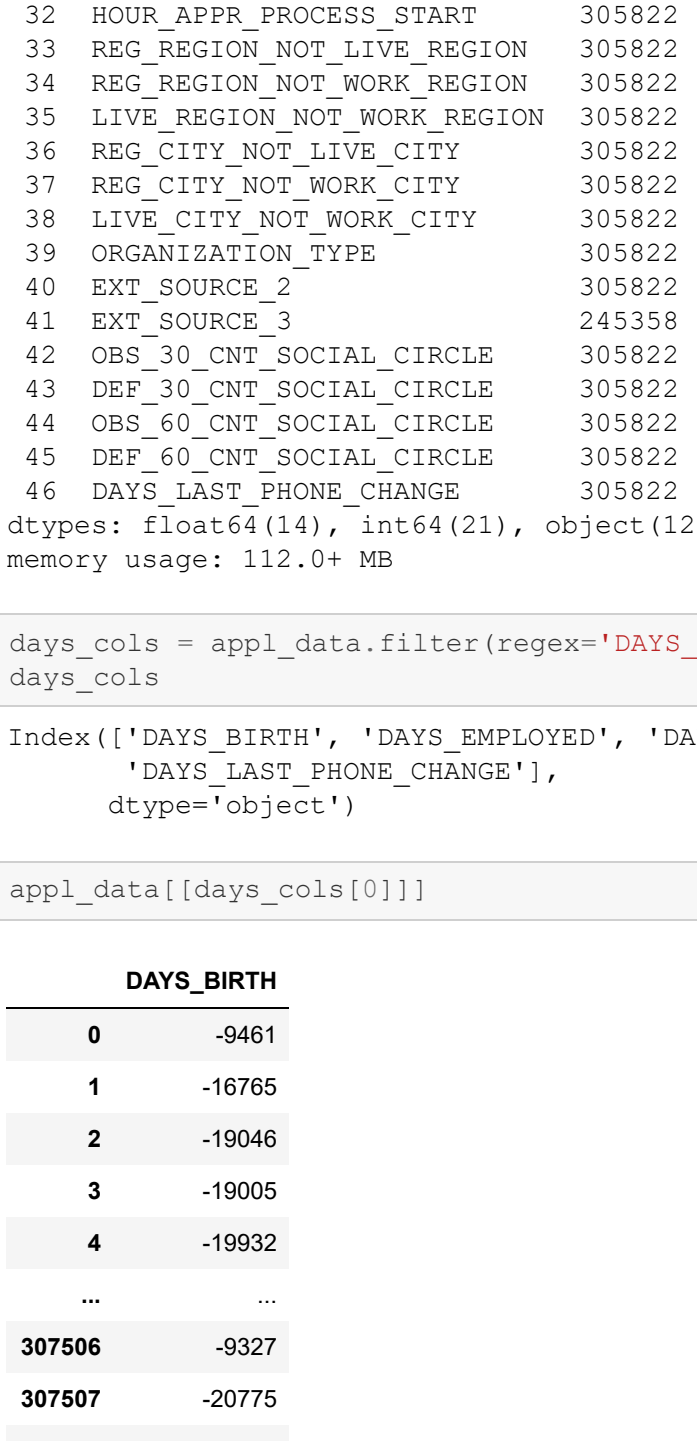
s = sns.histplot(data=app1\_data, x='AMT\_REQ\_CREDIT\_BUREAU\_WEEK', ax=axes[0][2]) s.set(klabel=None)

s = sns.histplot(data=app1\_data, x='AMT\_REQ\_CREDIT\_BUREAU\_MON', ax=axes[1][0]) s.set(klabel=None)

s = sns.histplot(data=app1\_data, x='AMT\_REQ\_CREDIT\_BUREAU\_QRT', ax=axes[1][1]) s.set(klabel=None)

s = sns.histplot(data=app1\_data, x='AMT\_REQ\_CREDIT\_BUREAU\_YEAR', ax=axes[1][2]) s.set(klabel=None)

plt.show()



<Figure size 72x432 with 0 Axes>

We observe that all these columns are mostly indicate that 0 requests are made to the credit bureau. Rest of the values are very less significant. The result is biased for the business reason that these values are this making any form of imputation unreliable. Hence we can consider dropping these columns.

In [81]: drop\_cols = app1\_data.filter(regex='AMT\_REQ\_CREDIT\_BUREAU', axis=1).columns

In [82]: #Dropping the above columns app1\_data.drop(drop\_cols,axis=1,inplace=True)

In [83]: app1\_data.columns[app1\_data.isna().any()]

Out[83]: Index(['EXT\_SOURCE\_3'], dtype='object')

## 2.4 Data Preparation for Numerical Data Analysis

In [84]: app1\_data.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 305822 entries, 0 to 307510
Data columns (total 47 columns):
 #   Column                Non-Null Count  Dtype  
---  --
 0   SK_ID_CURR            305822 non-null  int64  
 1   TARGET                305822 non-null  object  
 2   NAME_CONTRACT_TYPE    305822 non-null  object  
 3   CODE_GENDER           305822 non-null  object  
 4   FLAG_OWN_CAR          305822 non-null  object  
 5   FLAG_OWN_REALTY       305822 non-null  object  
 6   REG_REGION_NOT_LIVE_REGION  305822 non-null  object  
 7   AMT_INCOME_TOTAL     305822 non-null  float64 
 8   AMT_ANNUITY           305822 non-null  float64 
 9   AMT_GOODS_PRICE       305822 non-null  float64 
10  AMT_REQ_CREDIT_BUREAU_HOUR  305822 non-null  float64 
11  NAME_TYPE_SUITE       305822 non-null  object  
12  NAME_INCOME_TYPE      305822 non-null  object  
13  NAME_EDUCATION_TYPE   305822 non-null  object  
14  NAME_FAMILY_STATUS    305822 non-null  object  
15  NAME_HOUSING_TYPE     305822 non-null  object  
16  REGION_POPULATION_RELATIVE  305822 non-null  float64 
17  DAYS_BIRTH            305822 non-null  int64  
18  DAYS_EMPLOYED         305822 non-null  int64  
19  DAYS_REGISTRATION     305822 non-null  float64 
20  DAYS_ID_PUBLISH       305822 non-null  int64  
21  FLAG_MOBIL            305822 non-null  int64  
22  FLAG_EMP_PHONE        305822 non-null  int64  
23  FLAG_WORK_PHONE       305822 non-null  int64  
24  FLAG_CONT_MOBILE      305822 non-null  int64  
25  FLAG_PHONE            305822 non-null  int64  
26  FLAG_EMAIL            305822 non-null  int64  
27  OCCUPATION_TYPE       305822 non-null  object  
28  CNT_FAM_MEMBERS       305822 non-null  float64 
29  REGION_RATING_CLIENT  305822 non-null  int64  
30  REGION_RATING_CLIENT_W_CITY  305822 non-null  int64  
31  WEEKDAY_APPR_PROCESS_START  305822 non-null  object  
32  HOUR_APPR_PROCESS_START  305822 non-null  int64  
33  REG_REGION_NOT_LIVE_REGION  305822 non-null  object  
34  REG_REGION_NOT_WORK_REGION  305822 non-null  int64  
35  LIVE_REGION_NOT_WORK_REGION  305822 non-null  int64  
36  REG_CITY_NOT_LIVE_CITY  305822 non-null  int64  
37  REG_CITY_NOT_WORK_CITY  305822 non-null  int64  
38  LIVE_CITY_NOT_WORK_CITY  305822 non-null  int64  
39  ORGANIZATION_TYPE     305822 non-null  object  
40  EXT_SOURCE_2          245358 non-null  float64 
41  EXT_SOURCE_3          245358 non-null  float64 
42  OBS_30_CNT_SOCIAL_CIRCLE  305822 non-null  int64  
43  OBS_60_CNT_SOCIAL_CIRCLE  305822 non-null  int64  
44  OBS_60_CNT_SOCIAL_CIRCLE  305822 non-null  float64 
45  DEF_30_CNT_SOCIAL_CIRCLE  305822 non-null  float64 
46  DAYS_LAST_PHONE_CHANGE  305822 non-null  float64 
dtypes: float64(11), int64(21), object(12)
memory usage: 112.0+ MB(64(21), object(12))
```

In [85]: days\_cols = app1\_data.filter(regex='DAYS\_', axis=1).columns days\_cols

Out[85]: Index(['DAYS\_BIRTH', 'DAYS\_EMPLOYED', 'DAYS\_REGISTRATION', 'DAYS\_ID\_PUBLISH', 'DAYS\_LAST\_PHONE\_CHANGE'], dtype='object')

In [86]: app1\_data[days\_cols]

Out[86]:

	DAYS_BIRTH
0	-9461
1	-16765
2	-19006
3	-19005
4	-19932
...	...
307506	-9327
307507	-20775
307508	-14966
307509	-11961
307510	-16856

305822 rows x 1 columns

In [87]: #DAYS\_BIRTH is the Client's age in days at the time of application #convert DAYS\_BIRTH to AGE #convert DAYS\_BIRTH to AGE app1\_data['AGE'] = -round(app1\_data['DAYS\_BIRTH']/365,0) app1\_data['AGE']

Out[87]:

0	26.0
1	46.0
2	52.0
3	52.0
4	55.0
...	...
307506	25.0
307507	57.0
307508	41.0
307509	33.0
307510	46.0

Name: AGE, Length: 305822, dtype: float64

In [88]: #DAYS\_EMPLOYED #convert days before the application the person started current employment #convert DAYS\_EMPLOYED to YEARS\_EMPLOYED app1\_data['YEARS\_EMPLOYED'] = abs(round(app1\_data['DAYS\_EMPLOYED']/365, 0))

#DAYS\_REGISTRATION #how many days before the application did client change his registration #convert DAYS\_REGISTRATION to YEARS\_REGISTRATION app1\_data['YEARS\_REGISTRATION'] = abs(round(app1\_data['DAYS\_REGISTRATION']/365, 0))

#DAYS\_ID\_PUBLISH #how many days before the application did client change the identity document with which he applied for the loan #convert DAYS\_ID\_PUBLISH to YEARS\_ID\_PUBLISH app1\_data['YEARS\_ID\_PUBLISH'] = abs(round(app1\_data['DAYS\_ID\_PUBLISH']/365, 0))

#DAYS\_LAST\_PHONE\_CHANGE #convert DAYS\_LAST\_PHONE\_CHANGE to YEARS\_LAST\_PHONE\_CHANGE #convert DAYS\_LAST\_PHONE\_CHANGE app1\_data['YEARS\_LAST\_PHONE\_CHANGE'] = abs(round(app1\_data['DAYS\_LAST\_PHONE\_CHANGE']/365, 0))

In [89]: app1\_data[['AGE', 'YEARS\_EMPLOYED', 'YEARS\_REGISTRATION', 'YEARS\_ID\_PUBLISH', 'YEARS\_LAST\_PHONE\_CHANGE']]

Out[89]:

	AGE	YEARS_EMPLOYED	YEARS_REGISTRATION	YEARS_ID_PUBLISH	YEARS_LAST_PHONE_CHANGE
0	26.0	2.0	10.0	6.0	3.0
1	46.0	3.0	3.0	1.0	2.0
2	52.0	1.0	12.0	7.0	2.0
3	52.0	8.0	27.0	7.0	2.0
4	55.0	8.0	12.0	9.0	3.0
...	...	...	...	...	...
307506	26.0	1.0	23.0	5.0	1.0
307507	57.0	1001.0	12.0	11.0	0.0
307508	41.0	22.0	18.0	14.0	5.0
307509	33.0	13.0	7.0	3.0	1.0
307510	46.0	3.0	14.0	1.0	2.0

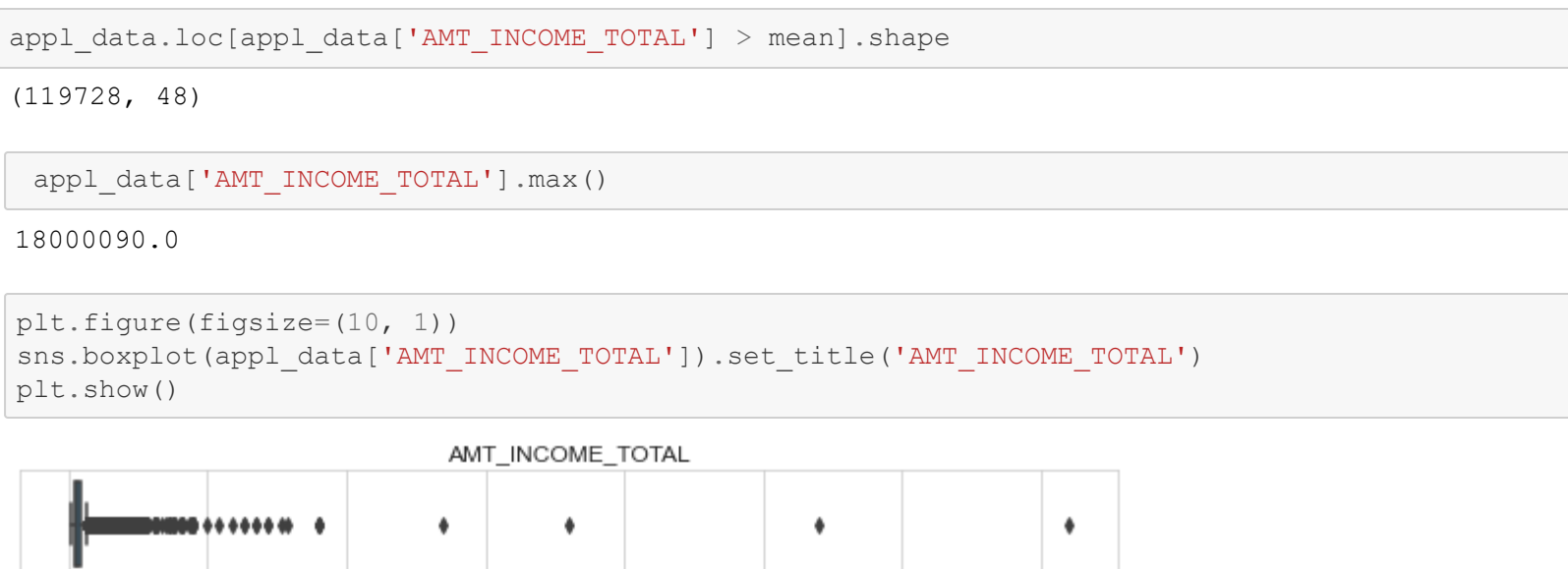
305822 rows x 5 columns

## 2.5 Binning required columns

In [92]: #Binning for Annuity and AMT\_INCOME\_TOTAL #convert AMT\_INCOME\_TOTAL to ANNUITY #convert AMT\_INCOME\_TOTAL to ANNUITY app1\_data['ANNUITY\_CLASS'] = pd.cut(app1\_data['AMT\_ANNUITY'],\ bins=(0,50000,70000,100000,260000),\ labels=['low','medium','high','v.high'])

plt.figure(figsize=(16, 1.5)) sns.set\_style('whitegrid') sns.boxplot(data=app1\_data, x='AMT\_ANNUITY', y='Annuity\_Class').set\_title('AMT\_ANNUITY') plt.show()

AMT\_ANNUITY



In [93]: app1\_data['Annuity\_Class'].value\_counts()

Out[93]: low 284653 medium 18596 high 2571 v.high 502 Name: Annuity\_Class, dtype: int64

In [94]: app1\_data.shape

Out[94]: (305822, 48)

In [95]: app1\_data['AMT\_INCOME\_TOTAL'].describe()

Out[95]:

AMT_INCOME_TOTAL	
min	0.0
25%	0.0
50%	0.0
75%	0.0
max	12.0

In [96]: #The above record is having a very high income value and is biasing the other records. #Removing the outlier app1\_data = app1\_data.sort\_values(by='AMT\_INCOME\_TOTAL',ascending=False)[1:]

In [98]: mean = app1\_data['AMT\_INCOME\_TOTAL'].mean() mean

Out[98]: 165297.0670391667

In [99]: app1\_data.loc[app1\_data['AMT\_INCOME\_TOTAL'] > mean].shape

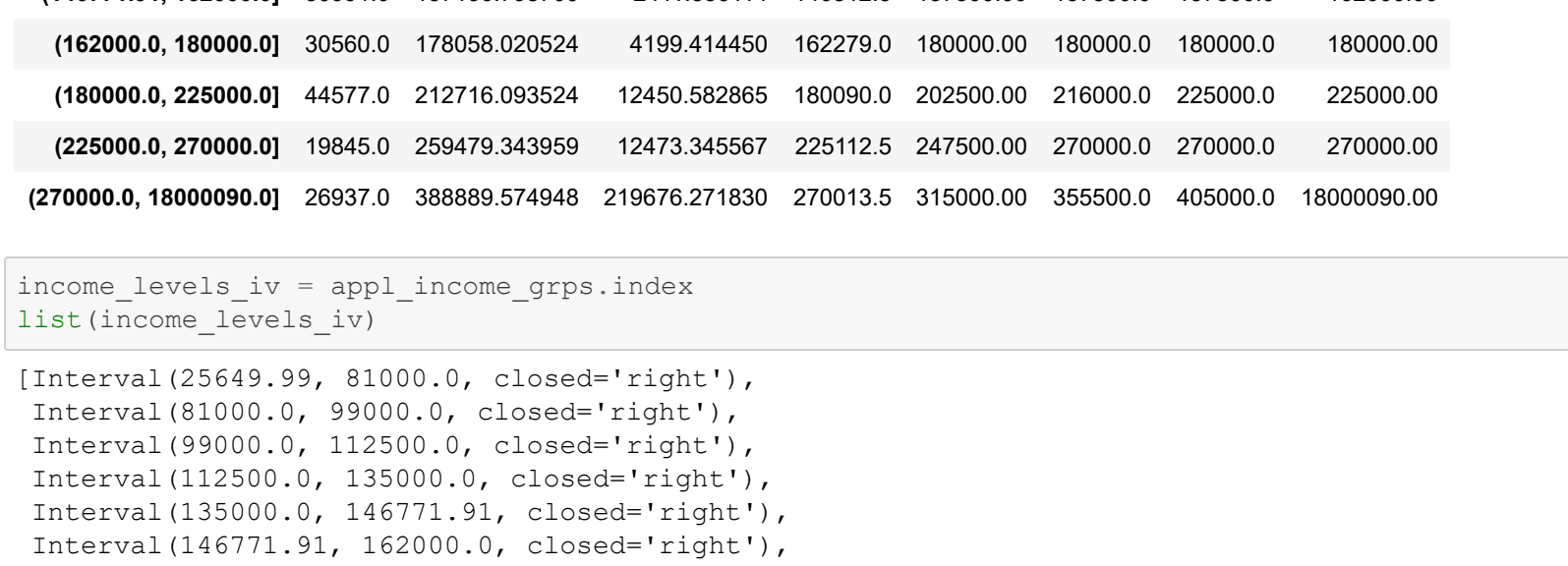
Out[99]: (119728, 48)

In [100]: app1\_data['AMT\_INCOME\_TOTAL'].max()

Out[100]: 18000090.0

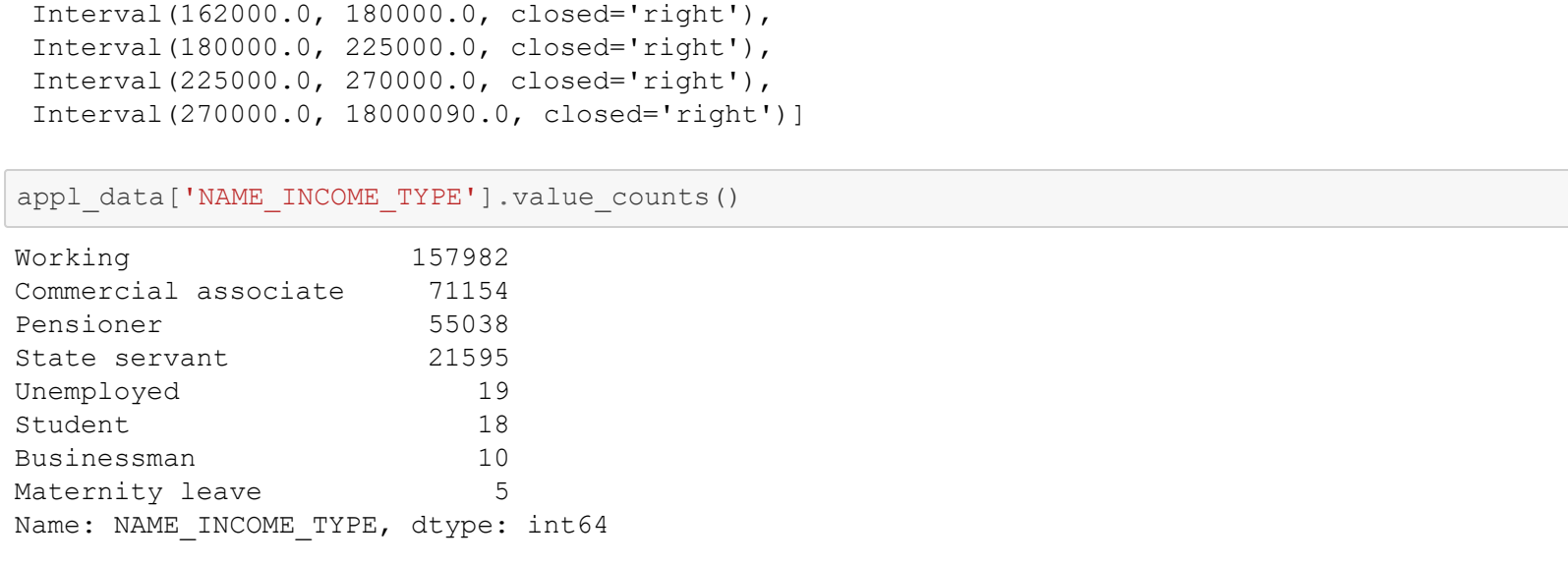
In [101]: plt.figure(figsize=(10, 1)) sns.boxplot(app1\_data['AMT\_INCOME\_TOTAL']).set\_title('AMT\_INCOME\_TOTAL') plt.show()

AMT\_INCOME\_TOTAL



In [102]: plt.figure(figsize=(10, 1)) sns.distplot(app1\_data['AMT\_INCOME\_TOTAL'])

AMT\_INCOME\_TOTAL



In [103]: app1\_data['Income\_Class'] = pd.cut(app1\_data['AMT\_INCOME\_TOTAL'],\ bins=(0,30000,60000,120000,240000,480000,960000),\ labels=['low','medium','high','v.high'])

In [104]: app1\_income\_grps = app1\_data.groupby(by='Income\_Class')['AMT\_INCOME\_TOTAL'].describe()

Out[104]:

Income_Class	count	mean	std	min	25%	50%	75%	max
(25649.99, 81000.0]	33211.0	60246.127214	12096.275007	25650.0	58500.00	67500.00	76500.00	81000.00
(81000.0, 99000.0]	30101.0	91088.792065	3401.030365	81160.0	90000.00	90000.00	90000.00	99000.00
(99000.0, 112500.0]	36694.0	111324.451613	2844.519397	99031.5	112500.00	112500.00	112500.00	112500.00
(112500.0, 135000.0]	48589.0	134847.429542	5860.487584	112500.00	135000.00	135000.00	135000.00	135000.00
(135000.0, 146771.91]	4307.0	142941.286222	2050.488281	135296.5	143696.25	144000.00	144000.00	146771.91
(146771.91, 162000.0]	30960.0	157136.798700	2411.856114	146812.5	157500.00	157500.00	157500.00	162000.00
(162000.0, 180000.0]	30580.0	170568.020524	4199.444500	162279.0	180000.00	180000.00	180000.00	180000.00
(180000.0, 225000.0]	44577.0	212716.093524	12450.582865	180000.00	202500.00	216000.00	225000.00	225000.00
(225000.0, 270000.0]	18845.0	259479.343959	12475.345567	225112.5	247500.00	270000.00	270000.00	270000.00
(270000.0, 18000090.0]	28937.0	388889.574948	219678.271830	270013.5	315000.00	355500.00	405000.00	18000090.00

In [105]: income\_levels\_iv = app1\_income\_grps.index list(income\_levels\_iv)

Out[105]: [Interval(25649.99, 81000.0, closed='right'), Interval(81000.0, 99000.0, closed='right'), Interval(99000.0, 112500.0, closed='right'), Interval(112500.0, 135000.0, closed='right'), Interval(135000.0, 146771.91, closed='right'), Interval(146771.91, 162000.0, closed='right'), Interval(162000.0, 180000.0, closed='right'), Interval(180000.0, 225000.0, closed='right'), Interval(225000.0, 270000.0, closed='right'), Interval(270000.0, 18000090.0, closed='right')]

In [106]: app1\_data['NAME\_INCOME\_TYPE'].value\_counts()

Out[106]: Working 157982 Commercial associate 71154 Pensioner 55038 State servant 21595 Unemployed 19 Student 18 Businessman 10 Maternity leave 5 Name: NAME\_INCOME\_TYPE, dtype: int64

In [107]: app1\_data['NAME\_INCOME\_TYPE'] = app1\_data['NAME\_INCOME\_TYPE'].replace({'Maternity leave': 'Businessman', 'Student': 'Unemployed', 'Pensioner': 'Other', 'State servant': 'Other', 'Unemployed': 'Unemployed'})

Out[108]: app1\_data['NAME\_INCOME\_TYPE']

Out[108]:

203693	Commercial associate
246858	Commercial associate
77768	Working
131127	Working
287463	Working
...	...
240137	Other
246104	Other
186643	Other
20127	Other
1678	Working

Name: NAME\_INCOME\_TYPE, Length: 305821, dtype: object

In [109]: #app1\_data.info()

## 2.6 Imbalance Ratio

- We observe the imbalance ratio in the target variable and if it is high, then we shall split the datasets for brevity.

In [110]: normal\_appl = app1\_data['TARGET'].value\_counts().sum() total\_appl = app1\_data['TARGET'].value\_counts()[0] defaulter\_appl = app1\_data['TARGET'].value\_counts()[1]

normal\_percentage = normal\_appl/total\_appl print('Non-Defaulter Applicants: ', normal\_appl, 'Percentage', round(normal\_percentage, 2))

defaulter\_percentage = defaulter\_appl/total\_appl print('Defaulters: ', defaulter\_appl, 'Percentage', round(defaulter\_percentage, 2))

Non-Defaulter Applicants: 281085 Percentage 0.92 Defaulters: 24736 Percentage 0.08

Note:

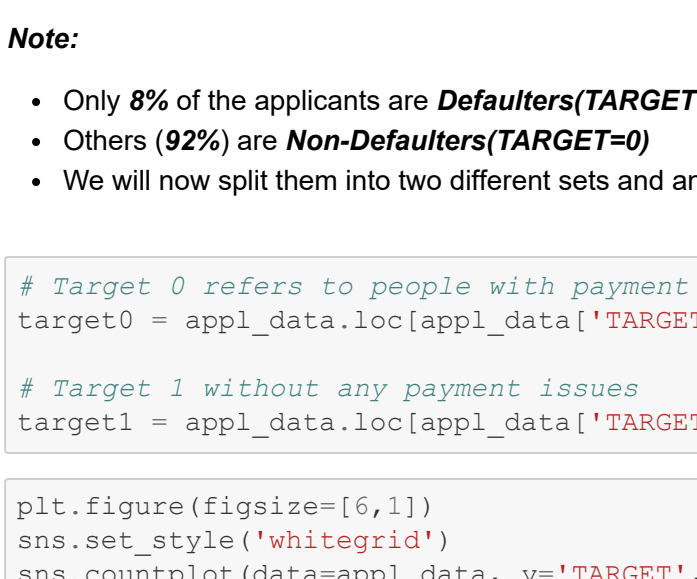
- Only 8% of the applicants are Defaulters(TARGET=1).
- Others (92%) are Non-Defaulters(TARGET=0).
- We will now split them into two different sets and analyze accordingly.

In [111]: #Target=0 refers to people with payment difficulties target0 = app1\_data.loc[app1\_data['TARGET'] == 0]

# Target=1 without any payment issues target1 = app1\_data.loc[app1\_data['TARGET'] == 1]

In [112]: plt.figure(figsize=(6, 1)) sns.set\_style('whitegrid') sns.histplot(data=app1\_data, y='TARGET', color='green').set\_title('Imbalance Ratio', fontsize=16) plt.show()

Imbalance Ratio



## 2.7 Univariate Analysis

Here, we look at variables directly influencing the target.

We use the term Target for the target variable.

We split the data into two datasets based on the Imbalance Ratio between the two types of applicants(Clients)

- Target 1: Defaulter Applicants (Clients) with payment difficulties: he/she had late payment more than X days on at least one of the first Y installments of the loan in our sample.
- Target 0: Non-Defaulter Applicants (All other cases) with no difficulties paying due on time

\*As described in the columns\_description.csv document

## 2.7.1 Reusable function for univariate analysis



```
In [113]: def Univariate(x):
plt.figure(figsize=(12,6))
sns.set_style('whitegrid')
sns.countplot(data=appl_data, x=x, hue='TARGET')
plt.xticks(rotation=90)
plt.title(f'Distribution of {x}', fontsize=25)
plt.xlabel('log')
plt.ylabel('N', fontsize=20)
plt.show()
```

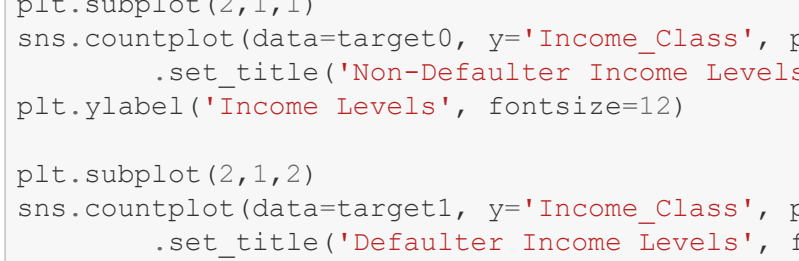
## 2.7.2 Gender vs Target

```
In [114]: plt.figure(figsize=(15,11))
sns.set_style('whitegrid')
plt.ylabel('Gender', fontsize=10)
gender_order = target0['CODE_GENDER'].value_counts().index

plt.subplot(1,2,1)
sns.countplot(data=target0, y='CODE_GENDER', order=gender_order, palette=['cyan', 'Indigo'])
plt.set_title('Non-Defaulter (By Gender)', fontsize=15)
plt.ylabel('Gender', fontsize=12)
plt.xscale('log')

plt.subplot(1,2,2)
sns.countplot(data=target1, y='CODE_GENDER', order=gender_order, palette=['cyan', 'Indigo'])
plt.set_title('Defaulter (By Gender)', fontsize=15)
plt.ylabel('Gender', fontsize=12)
plt.xscale('log')

plt.xlabel('No. of Applicants', fontsize=12)
plt.show()
```



### Observations

We generally notice that

- **Female** applicants are higher than the **Male** applicants and directly influences the numbers
- **Gender** has **No** impact on the **Target**, other than the generic imbalance across the Defaulters and the Non-Defaulters.
- We will **Ignore Gender** for all further analysis

## 2.7.3 Income vs Target

```
In [115]: plt.figure(figsize=(15,10))
sns.set_style('whitegrid')

plt.subplot(2,1,1)
sns.countplot(data=target0, y='Income_Class', palette=['cyan', 'Indigo'])
plt.set_title('Non-Defaulter Income Levels', fontsize=15)
plt.ylabel('Income Levels', fontsize=12)
plt.xscale('log')

plt.subplot(2,1,2)
sns.countplot(data=target1, y='Income_Class', palette=['cyan', 'Indigo'])
plt.set_title('Defaulter Income Levels', fontsize=15)
plt.ylabel('Income Levels', fontsize=12)
plt.xscale('log')

plt.xlabel('No. of Applicants', fontsize=12)
plt.show()
```



```
In [116]: income_levels_v[4]
Out[116]: Interval(135000.0, 146771.91, closed='right')

In [117]: appl_data['Income_Class'].value_counts()
Out[117]:
(135000.0, 135000.0]    48598
(180000.0, 225000.0]    44577
(99000.0, 112500.0]     36694
(25649.99, 81000.0]     33211
(146771.91, 162000.0]    30991
(162000.0, 180000.0]    30560
(81000.0, 95000.0]       30101
(270000.0, 1800000.0]    26937
(225000.0, 270000.0]    19845
(135000.0, 146771.91]   4307
Name: Income_Class, dtype: int64
```

### Observations

- The (135000.0, 146771.91) income level applicants seem to be 'Very Less(Only 4307)' in number.
- **Regardless\*\* Income Levels** different applicants face payment difficulties.
- No of defaulters and non-defaulters seems to be proportionally **same**
- Considering the above observations, it seems like **along with the income**, there are other factors in play when it comes to the **Target**. We might require bi-variate or multi-variate analysis

## 2.7.4 Annuity vs Target

```
In [118]: plt.figure(figsize=(10,7))
sns.set_style('whitegrid')

plt.subplot(2,1,1)
sns.countplot(data=target0, y='Annuity_Class', palette=['cyan', 'Indigo'])
plt.set_title('Non-Defaulter Annuity Levels', fontsize=15)
plt.ylabel('Annuity Levels', fontsize=12)
plt.xscale('log')

plt.subplot(2,1,2)
sns.countplot(data=target1, y='Annuity_Class', palette=['cyan', 'Indigo'])
plt.set_title('Defaulter Annuity Levels', fontsize=15)
plt.ylabel('Annuity Levels', fontsize=12)
plt.xscale('log')

plt.xlabel('No. of Applicants', fontsize=12)
plt.show()
```



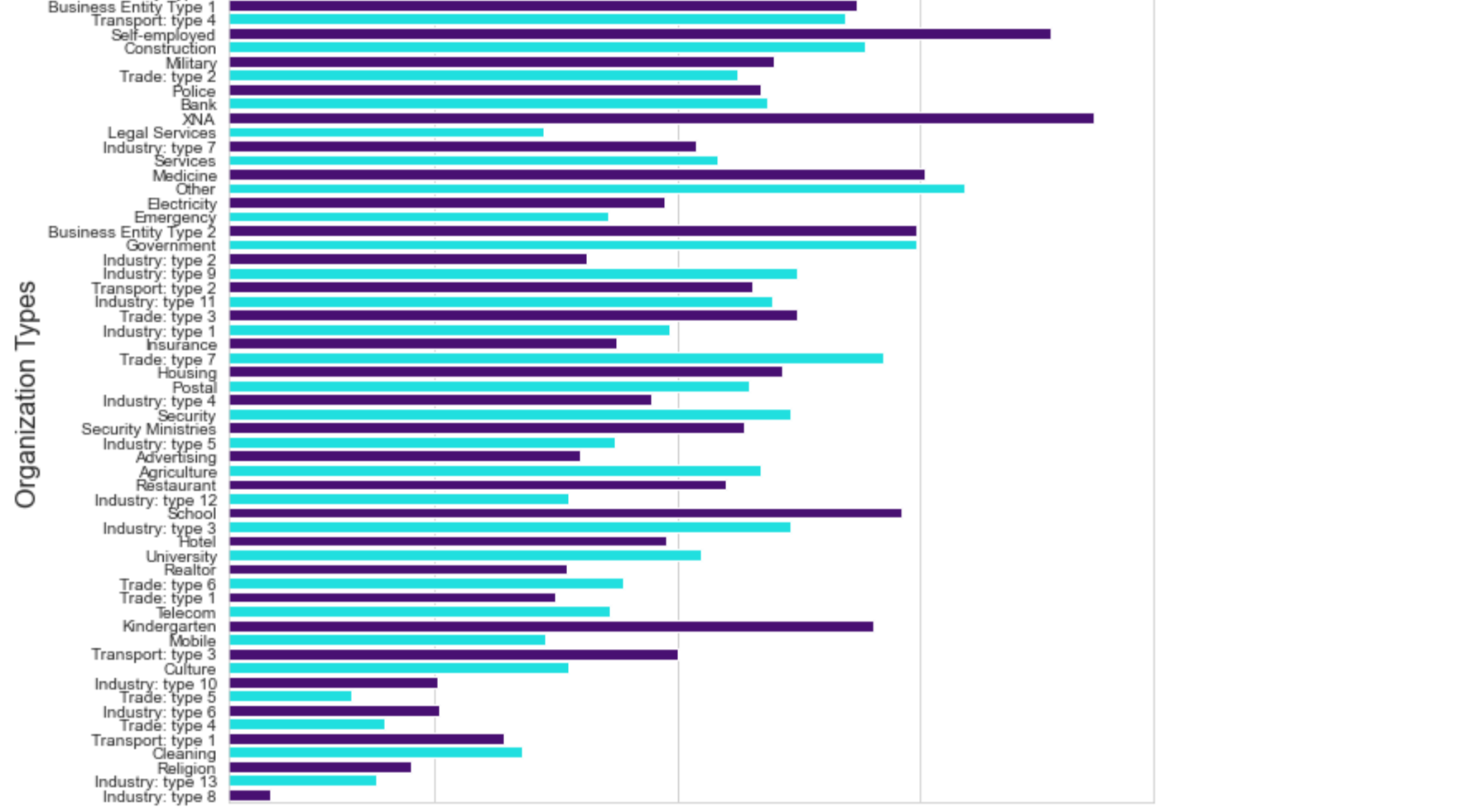
## 2.7.5 Family Size vs Target

```
In [119]: plt.figure(figsize=(15,10))
sns.set_style('whitegrid')

plt.subplot(2,1,1)
sns.countplot(data=target0, y='CNT_FAM_MEMBERS', palette=['cyan', 'Indigo'])
plt.set_title('Non-Defaulter Family Size', fontsize=15)
plt.ylabel('No. of Family Members', fontsize=12)
plt.xscale('log')

plt.subplot(2,1,2)
sns.countplot(data=target1, y='CNT_FAM_MEMBERS', palette=['cyan', 'Indigo'])
plt.set_title('Defaulter Family Size', fontsize=15)
plt.ylabel('No. of Family Members', fontsize=12)
plt.xscale('log')

plt.xlabel('No. of Applicants', fontsize=12)
plt.show()
```



```
In [120]: target0['CNT_FAM_MEMBERS'].value_counts()
Out[120]:
2.0    145548
1.0     61760
3.0     47741
4.0     22443
5.0      3135
6.0       352
7.0        75
8.0         14
9.0          6
10.0         2
11.0         2
12.0         2
13.0         2
14.0         2
15.0         1
16.0         1
17.0         1
18.0         1
19.0         1
20.0         1
Name: CNT_FAM_MEMBERS, dtype: int64

In [121]: target1['CNT_FAM_MEMBERS'].value_counts()
Out[121]:
2.0    11971
1.0     4553
3.0     4587
4.0     2129
5.0      326
6.0       75
7.0        6
8.0         6
9.0         6
10.0         1
11.0         1
12.0         1
13.0         1
14.0         1
15.0         1
16.0         1
17.0         1
18.0         1
19.0         1
20.0         1
Name: CNT_FAM_MEMBERS, dtype: int64
```

### Observations

- Interestingly, All important social variables so far seem not to individually effecting on the Target.
- Let's run them for rest of the columns.
- Considering the above observations, it seems like along with the income, there are other factors in play when it comes to the Target. We might require bi-variate or multi-variate analysis

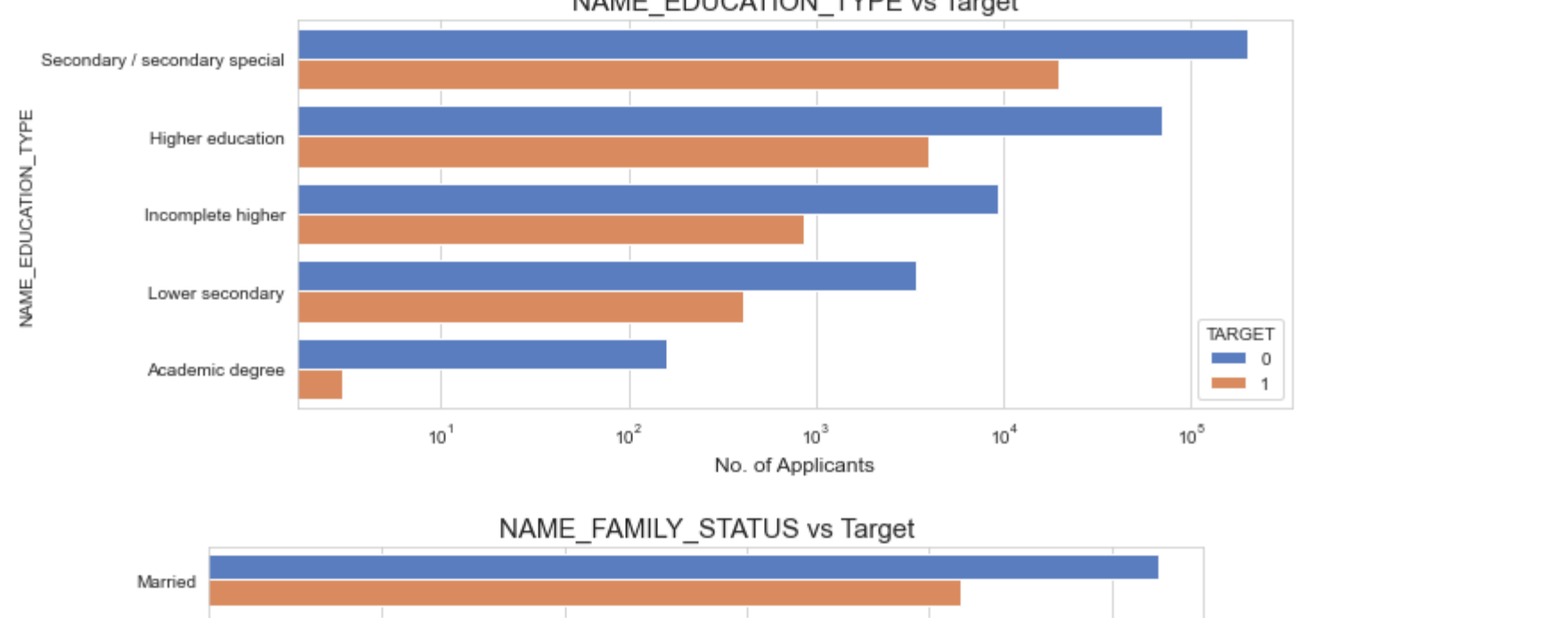
## 2.7.6 Loan Type vs Target

```
In [122]: plt.figure(figsize=(15,11))
sns.set_style('whitegrid')

plt.subplot(1,2,1)
sns.countplot(data=target0, y='NAME_CONTRACT_TYPE', palette=['cyan', 'Indigo'])
plt.set_title('Non-Defaulter Contract Types', fontsize=14)
plt.ylabel('Contract Types', fontsize=12)
plt.xlabel('No. of Applicants', fontsize=12)
plt.xscale('log')

plt.subplot(1,2,2)
sns.countplot(data=target1, y='NAME_CONTRACT_TYPE', palette=['cyan', 'Indigo'])
plt.set_title('Defaulter Contract Types', fontsize=14)
plt.ylabel('Contract Types', fontsize=12)
plt.xlabel('No. of Applicants', fontsize=12)
plt.xscale('log')

plt.show()
```



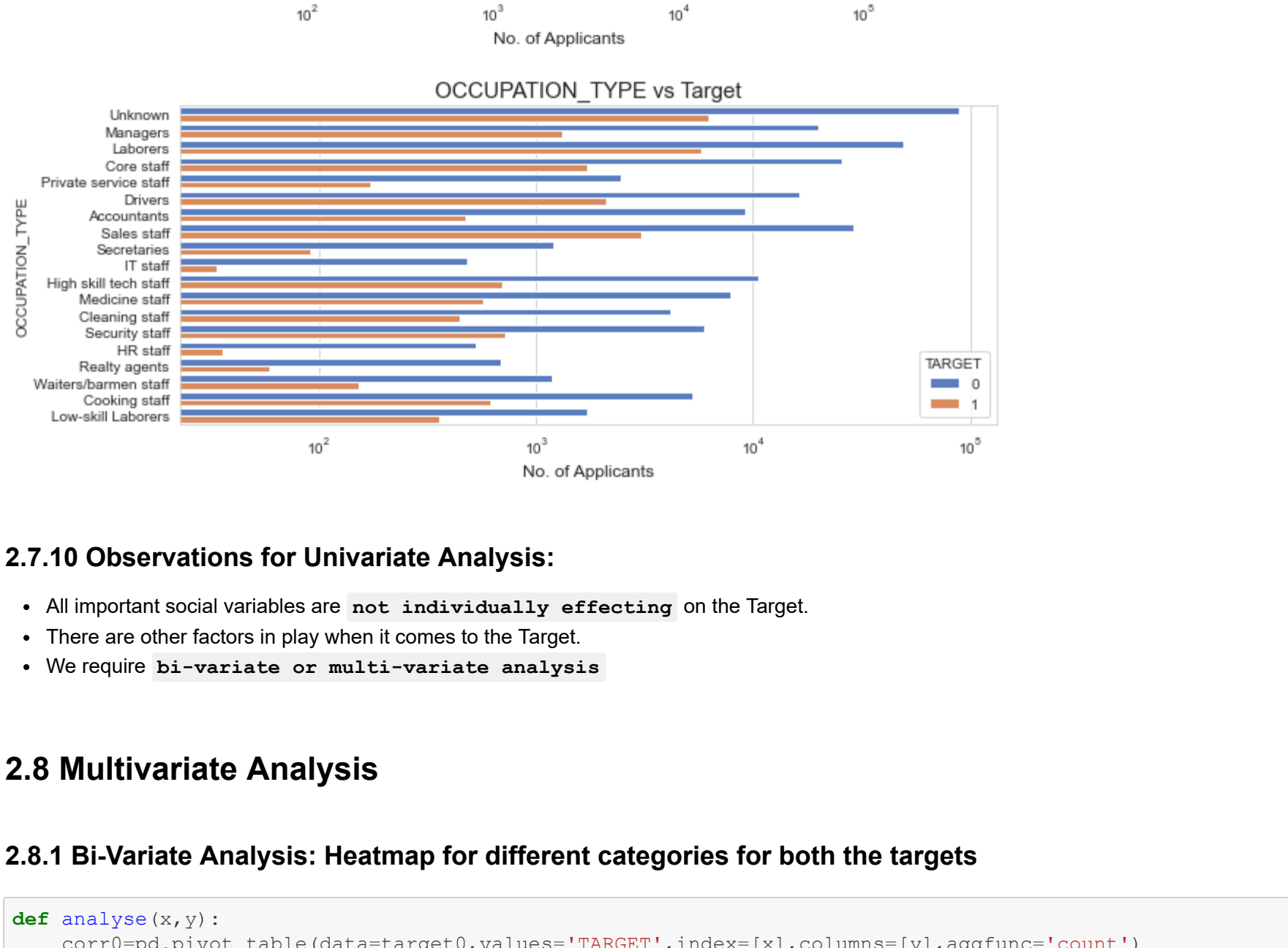
## 2.7.7 Organization Type vs Target

```
In [123]: plt.figure(figsize=(10,20))
sns.set_style('whitegrid')

plt.subplot(2,1,1)
sns.countplot(data=target0, y='ORGANIZATION_TYPE', palette=['cyan', 'Indigo'])
plt.set_title('Non-Defaulter Organization Types (By Gender)', fontsize=16)
plt.ylabel('Organization Types', fontsize=12)
plt.xlabel('log', fontsize=12)
plt.xscale('log')

plt.subplot(2,1,2)
sns.countplot(data=target1, y='ORGANIZATION_TYPE', palette=['cyan', 'Indigo'])
plt.set_title('Defaulter Organization Types (By Gender)', fontsize=16)
plt.ylabel('Organization Types', fontsize=12)
plt.xlabel('log', fontsize=12)
plt.xscale('log')

plt.show()
```



### Observations so far:

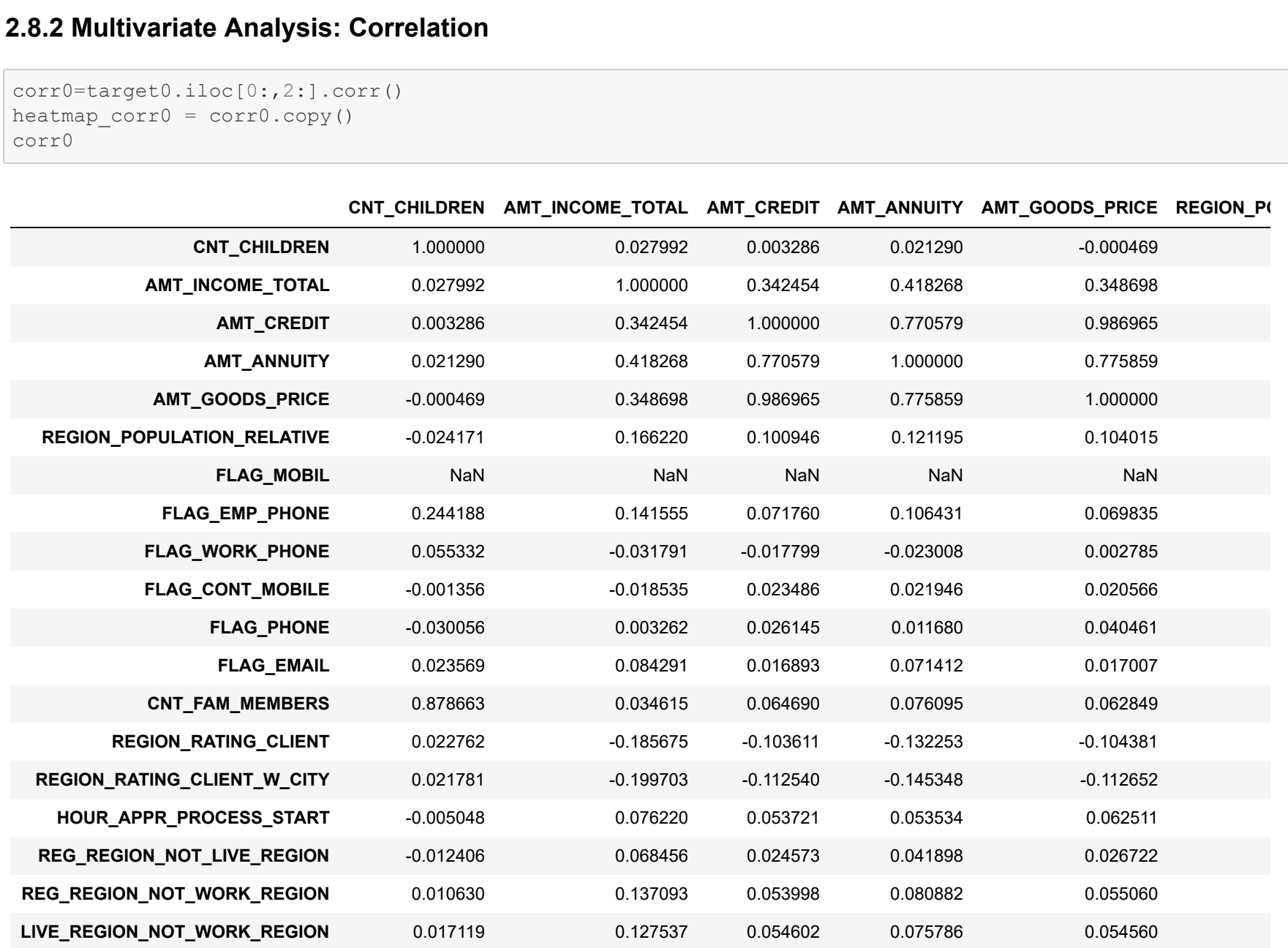
- Interestingly, All important social variables so far seem not to individually effecting on the Target.
- Let's run them for rest of the columns.
- Considering the above observations, it seems like along with the income, there are other factors in play when it comes to the Target. We might require bi-variate or multi-variate analysis

## 2.7.9 Running Univariate Analysis for all to assert the above

```
In [124]: def Univariate(x):
plt.figure(figsize=(10,4))
sns.set_style('whitegrid')
sns.countplot(data=appl_data, y=x, hue='TARGET', palette='muted')
plt.title(f'Distribution of {x}', fontsize=15)
plt.xlabel('No. of Applicants', fontsize=12)
plt.ylabel('log', fontsize=12)
plt.xscale('log')
plt.show()
```

```
In [125]: cols=['INCOME_TYPE', 'NAME_EDUCATION_TYPE', 'NAME_FAMILY_STATUS', 'NAME_HOUSING_TYPE', 'NAME_TYPE_SUITE', 'OCCUPATION_TYPE']
```

```
In [126]: for x in cols:
Univariate(x)
```



## 2.7.10 Observations for Univariate Analysis:

- All important social variables are not individually effecting on the Target.
- There are other factors in play when it comes to the Target.
- We require **bi-variate or multi-variate analysis**

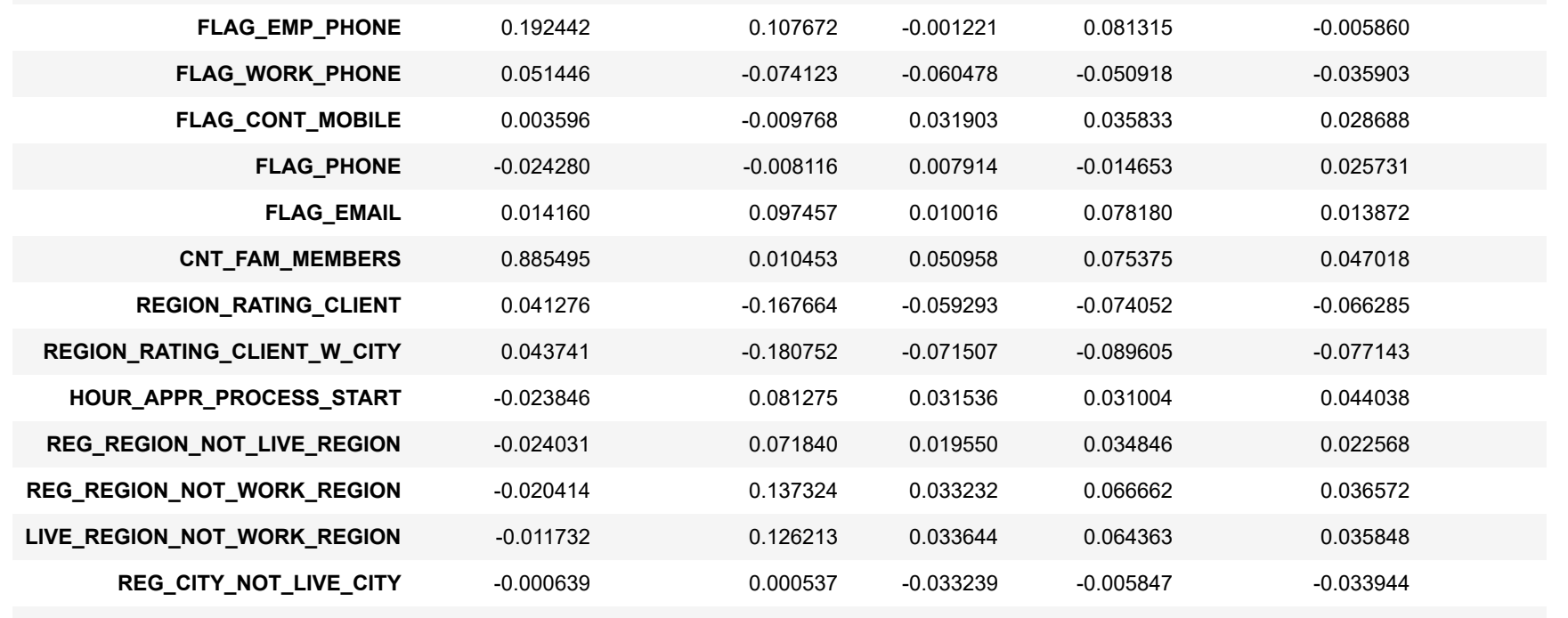
## 2.8 Multivariate Analysis

### 2.8.1 Bi-Variate Analysis: Heatmap for different categories the target

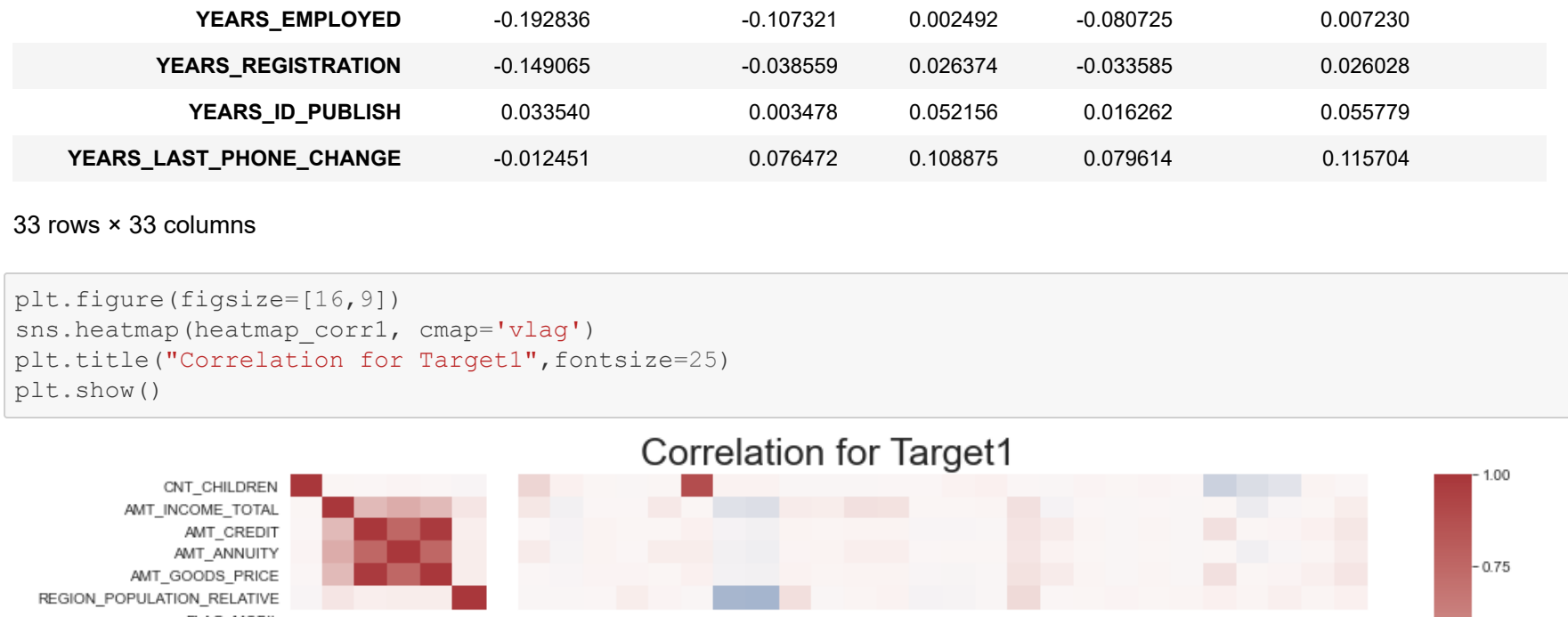
```
In [127]: def analyze(x,y):
corr=corr_pivot_table(data=target0, values='TARGET', index=[x], columns=[y], aggfunc='count')
corr=pd.pivot_table(data=target0, values='TARGET', index=[x], columns=[y], aggfunc='count')

plt.figure(figsize=(16,4))
plt.subplot(1,2,1)
sns.heatmap(corr0, cmap='Greens')
plt.xlabel(f'Target0 = {x}')
ax = plt.subplot(1,2,2)
sns.heatmap(corr1, cmap='Blues')
plt.xlabel(f'Target1 = {x}')
plt.ylabel(yticklabels([]))
plt.show()
```

```
In [128]: analyze('CNT_FAM_MEMBERS', 'Annuity_Class')
```



```
In [129]: analyze('INCOME_TYPE', 'Annuity_Class')
```



```
In [130]: analyze('Income_Class', 'Annuity_Class')
```



### Observations for Multivariate Analysis:

- No significance variations were observed for specific variables.
- This calls for multi-variate analysis
- Cross-referencing with Previous Application data could potentially yield more results

## 2.8.2 Multivariate Analysis: Correlation

```
In [131]: corr=corr_target.iloc[0,:2].corr()
heatmap_corr0 = corr0.copy()
```

```
Out[131]:
```

	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_PT
CNT_CHILDREN	1.000000	0.027992	0.003286	0.021290	-0.004469	
AMT_INCOME_TOTAL	0.027992	1.000000	0.342454	0.418268	0.348698	
AMT_CREDIT	-0.002209	0.342454	1.000000	0.757079	0.986865	
AMT_ANNUITY	0.021290	0.418268	0.757079	1.000000	0.775659	
AMT_GOODS_PRICE	-0.004469	0.348698	0.986865	0.775659	1.000000	
REGION_POPULATION_RELATIVE	-0.024171	0.166220	0.109946	0.121195	0.104015	
FLAG_EMP_PHONE	NaN	NaN	NaN	NaN	NaN	
FLAG_WORK_PHONE	0.244188	0.141555	0.071780	0.106431	0.069353	
FLAG_CONTACT_MOBILE	0.053326	-0.037391	-0.017799	-0.023008	0.027885	
FLAG_PHONE	-0.001356	-0.185355	0.023496	0.021946	0.020566	
FLAG_EMAIL	-0.002569	0.005262	0.026145	0.011680	0.040461	
FLAG_PHONE	0.023569	0.084291	0.016893	0.071412	0.017007	
CNT_FAM_MEMBERS	0.878663	0.034615	0.064890	0.076095	0.062849	
REGION_RATING_CLIENT	0.041276	-0.167864	-0.059293	-0.074052	-0.066285	
REGION_RATING_CLIENT_W_CITY	0.043741	-0.180752	-0.071507	-0.089605	-0.077143	
HOUR_APPR_PROCESS_START	-0.023846	0.081275	0.031536	0.031004	0.044038	
REG_REGION_NOT_LIVE_REGION	-0.024031	0.071840	0.019550	0.034846	0.022568	
REG_REGION_NOT_WORK_REGION	-0.024044	0.033242	0.036662	0.036662	0.036572	
LIVE_REGION_NOT_WORK_REGION	-0.017132	0.126213	0.033644	0.043633	0.035848	
REG_CITY_NOT_LIVE_CITY	-0.003639	0.000537	-0.028239	-0.005847	-0.033944	
REG_CITY_NOT_WORK_CITY	0.042572	0.009430	-0.083117	0.001438	-0.039099	
LIVE_CITY_NOT_WORK_CITY	0.053732	0.017115	-0.071052	0.003930	-0.017516	
EXT_SOURCE_2	0.012235	0.136898	0.121132	0.116656	0.130772	
EXT_SOURCE_3	-0.022444	-0.061981	0.078405	0.042357	0.079518	
OBS_30_CNT_SOCIAL_CIRCLE	0.025718	-0.006558	0.019218	0.004447	0.019872	
DEF_30_CNT_SOCIAL_CIRCLE	0.026097	-0.024705	-0.026280	-0.022747	-0.022423	
OBS_60_CNT_SOCIAL_CIRCLE	0.001541	-0.005862	0.019906	0.005479	0.020437	
DEF_60_CNT_SOCIAL_CIRCLE	-0.024783	-0.024381	-0.031081	-0.027748	-0.026598	
AGE	-0.258924	0.002472	0.135758	0.014412	0.135844	
YEARS_EMPLOYED	-0.192836	-0.107321	0.002492	-0.080725	0.007230	
YEARS_REGISTRATION	-0.149065	-0.038559	0.026374	-0.033585	0.026028	
YEARS_ID_PUBLISH	0.033540	0.003478	0.021566	0.016262	0.055779	
YEARS_LAST_PHONE_CHANGE	-0.012451	0.076472	0.108875	0.079614	0.115704	

```
In [134]: plt.figure(figsize=(16,9))
sns.heatmap(heatmap_corr0, cmap='Spectral')
plt.title('Correlation for Target0', fontsize=25)
plt.show()
```



```
In [133]: heatmap_target1.iloc[0,:2].corr()
Out[133]:
```

	CNT_CHILDREN	AMT_INCOME_TOTAL	AMT_CREDIT	AMT_ANNUITY	AMT_GOODS_PRICE	REGION_PT
CNT_CHILDREN	1.000000	0.002735	-0.002009	0.003286	-0.008419	
AMT_INCOME_TOTAL	0.002735	1.000000	0.324956	0.397928	0.327434	
AMT_CREDIT	-0.002009	0.324956	1.000000	0.757079	0.982846	
AMT_ANNUITY	0.003286	0.397928	0.757079	1.000000	0.752511	
AMT_GOODS_PRICE	-0.004469	0.327434	0.982846	0.752511	1.000000	
REGION_POPULATION_RELATIVE	-0.021240	0.117307	0.069774	0.072334	0.078964	
FLAG_EMP_PHONE	NaN	NaN	NaN	NaN	NaN	
FLAG_WORK_PHONE	0.102442	0.107872	-0.001221	0.081315	-0.005860	
FLAG_CONTACT_MOBILE	0.051446	-0.074123	-0.060478	-0.050918	-0.035903	
FLAG_PHONE	0.003596	-0.005978	0.011903	0.035853	0.028688	
FLAG_EMAIL	-0.024809	-0.008116	0.007914	-0.014653	0.025731	
FLAG_PHONE	0.014160	0.097457	0.010016	0.071880	0.013872	
CNT_FAM_MEMBERS	0.884595	0.010453	0.005058	0.075375	0.047018	
REGION_RATING_CLIENT	0.041276	-0.167864	-0.059293	-0.074052	-0.066285	
REGION_RATING_CLIENT_W_CITY	0.043741	-0.180752	-0.071507	-0.089605	-0.077143	
HOUR_APPR_PROCESS_START	-0.023846	0.081275	0.031536	0.031004	0.044038	
REG_REGION_NOT_LIVE_REGION	-0.024031	0.071840	0.019550	0.034846	0.022568	
REG_REGION_NOT_WORK_REGION	-0.024044	0.033242	0.036662	0.036662	0.036572	
LIVE_REGION_NOT_WORK_REGION	-0.017132	0.126213	0.033644	0.043633	0.035848	
REG_CITY_NOT_LIVE_CITY	-0.003639	0.000537	-0.028239	-0.005847	-0.033944	
REG_CITY_NOT_WORK_CITY	0.042572	0.009430	-0.083117	0.001438	-0.039099	
LIVE_CITY_NOT_WORK_CITY	0.053732	0.017115	-0.071052	0.003930	-0.017516	
EXT_SOURCE_2	0.012235	0.136898	0.121132	0.116656	0.130772	
EXT_SOURCE_3	-0.022444	-0.061981	0.078405	0.042357	0.079518	
OBS_30_CNT_SOCIAL_CIRCLE	0.025718	-0.006558	0.019218	0.004447	0.019872	
DEF_30_CNT_SOCIAL_CIRCLE	0.026097	-0.024705	-0.026280	-0.022747	-0.022423	
OBS_60_CNT_SOCIAL_CIRCLE	0.001541	-0.005862	0.019906	0.005479	0.020437	
DEF_60_CNT_SOCIAL_CIRCLE	-0.024783	-0.024381	-0.031081	-0.027748	-0.026598	
AGE	-0.258924	0.002472	0.135758	0.014412	0.135844	
YEARS_EMPLOYED	-0.192836	-0.107321	0.002492	-0.080725	0.007230	
YEARS_REGISTRATION	-0.149065	-0.038559	0.026374	-0.033585	0.026028	
YEARS_ID_PUBLISH	0.033540	0.003478	0.021566	0.016262	0.055779	
YEARS_LAST_PHONE_CHANGE	-0.012451	0.076472	0.108875	0.079614	0.115704	

```
In [134]: plt.figure(figsize=(16,9))
sns.heatmap(heatmap_corr0, cmap='Spectral')
plt.title('Correlation for Target1', fontsize=25)
plt.show()
```





```
1335: ## correlation between columns by observations
def correlation(x,y):
    plt.figure(figsize=(20,10))
    ax = plt.subplot(2,1,1)
    sns.scatterplot(target0[x],target0[y])
    plt.title('f'(x) VS f'(y) for Target 0',fontsize=20)
    plt.xlabel('f'(x)')
    ax.xaxis.get_label().set_fontsize(18)
    plt.ylabel('f'(y)')
    ax = plt.subplot(2,1,2)
    sns.scatterplot(target1[x],target1[y])
    plt.title('f'(x) VS f'(y) for Target 1',fontsize=20)
    plt.xlabel('f'(x)')
    ax.xaxis.get_label().set_fontsize(18)
    plt.ylabel('f'(y)')
    plt.show()

In [136]: correlation('AMT_ANNUITY','AMT_CREDIT')
```

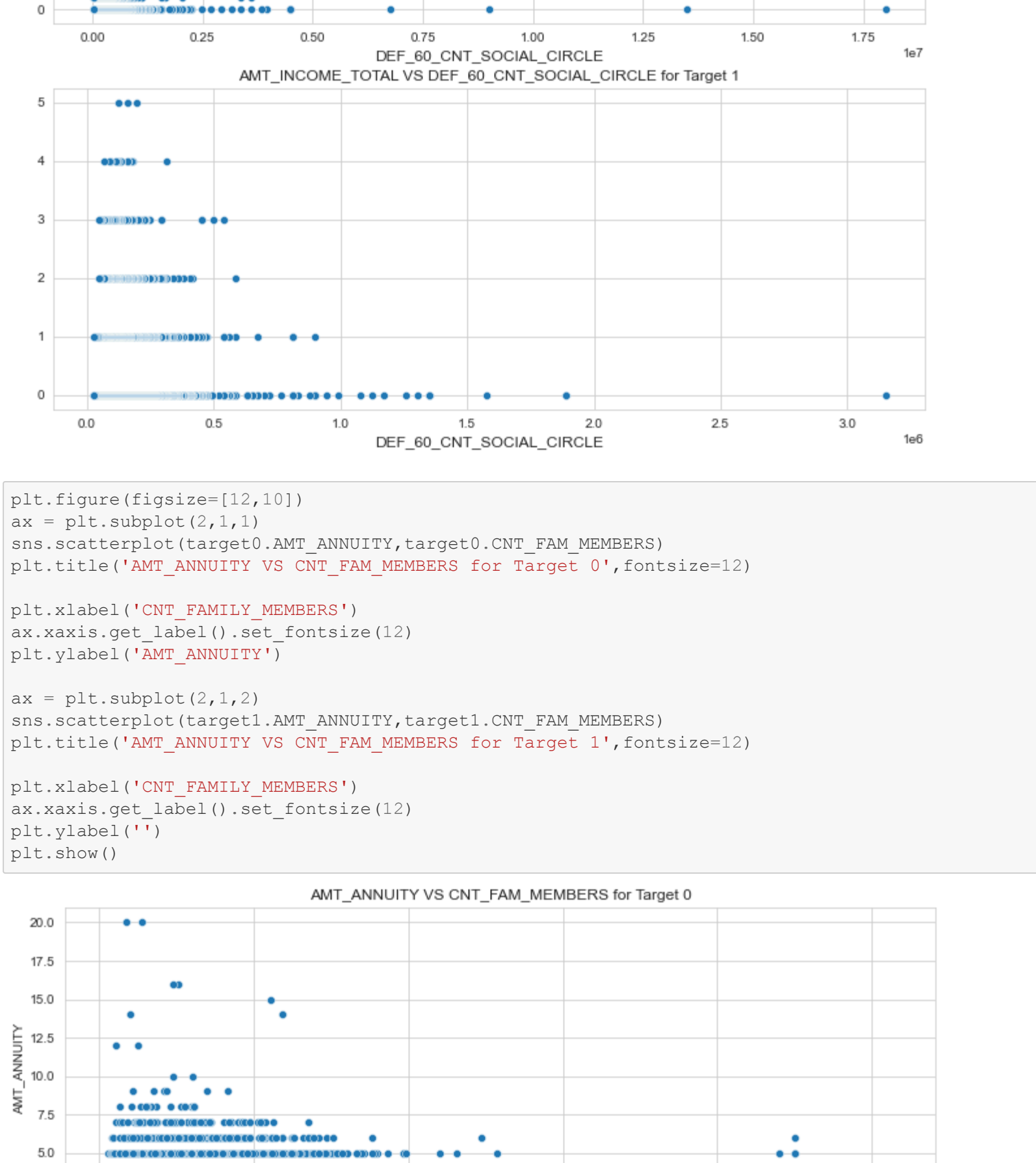


```
In [137]: correlation('AMT_ANNUITY','AMT_GOODS_PRICE')

AMT_ANNUITY VS AMT_GOODS_PRICE for Target 0
AMT_ANNUITY VS AMT_GOODS_PRICE for Target 1

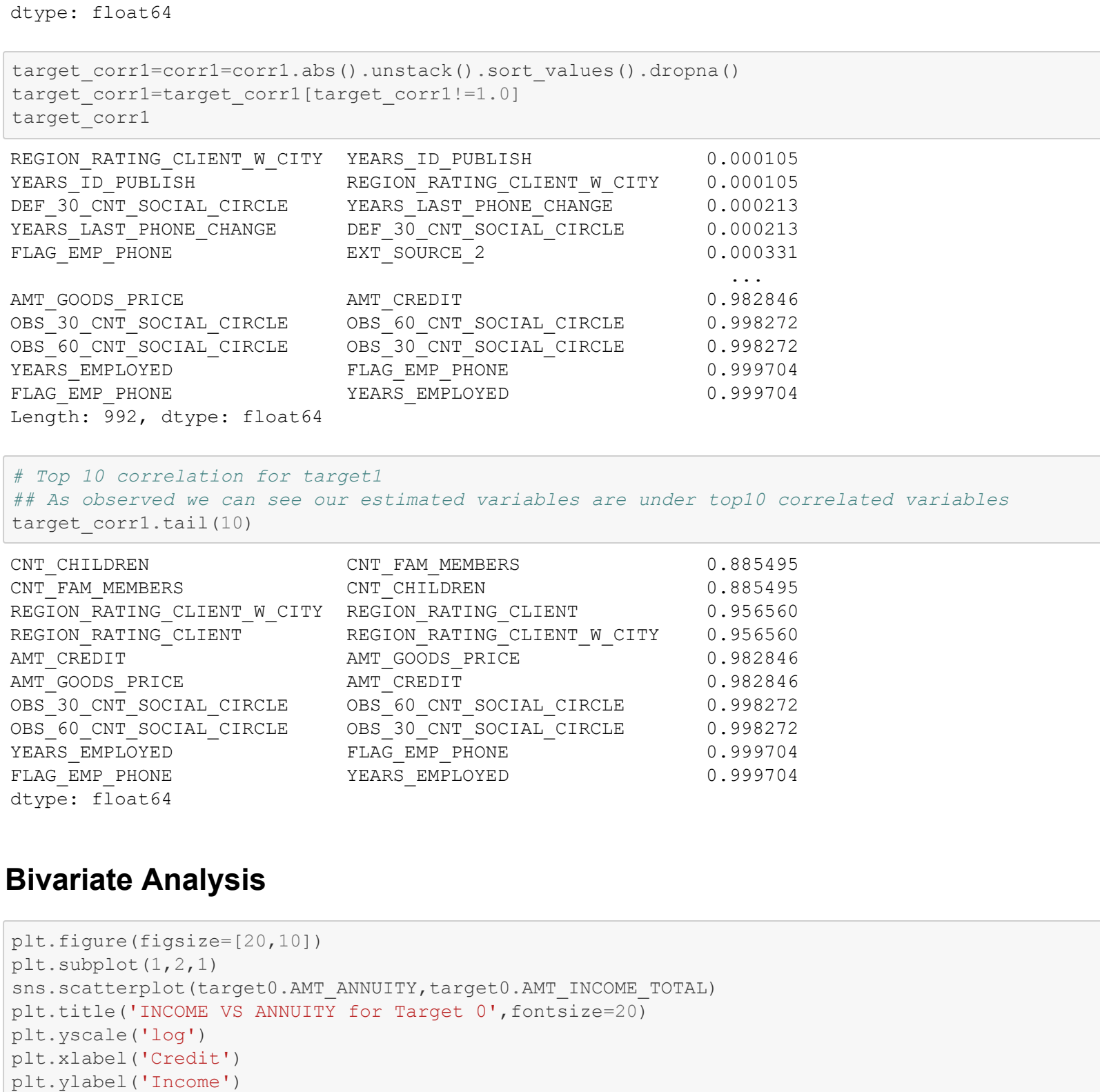
In [138]: correlation('AMT_GOODS_PRICE','AMT_CREDIT')

AMT_GOODS_PRICE VS AMT_CREDIT for Target 0
AMT_GOODS_PRICE VS AMT_CREDIT for Target 1
```



```
In [139]: plt.figure(figsize=(12,10))
ax = plt.subplot(2,1,1)
sns.scatterplot(target0.AMT_INCOME_TOTAL,target0.DEF_60_CNT_SOCIAL_CIRCLE)
plt.title('AMT_INCOME_TOTAL VS DEF_60_CNT_SOCIAL_CIRCLE for Target 0',fontsize=12)
plt.xlabel('AMT_INCOME_TOTAL')
ax = plt.subplot(2,1,2)
sns.scatterplot(target1.AMT_INCOME_TOTAL,target1.DEF_60_CNT_SOCIAL_CIRCLE)
plt.title('AMT_INCOME_TOTAL VS DEF_60_CNT_SOCIAL_CIRCLE for Target 1',fontsize=12)
plt.xlabel('AMT_INCOME_TOTAL')
ax.xaxis.get_label().set_fontsize(12)
plt.ylabel('DEF_60_CNT_SOCIAL_CIRCLE')
plt.show()

AMT_INCOME_TOTAL VS DEF_60_CNT_SOCIAL_CIRCLE for Target 0
AMT_INCOME_TOTAL VS DEF_60_CNT_SOCIAL_CIRCLE for Target 1
```



### Finding top10 correlations for target0 and target1

```
In [141]: target_corr0=target_corr0.abs().unstack().sort_values().dropna()
target_corr0
target_corr1

Out[141]: FLAG_PHONE      EXT_SOURCE_3      0.005002
EXT_SOURCE_3      FLAG_PHONE      0.000002
YEARS_LAST_PHONE_CHANGE      DEF_60_CNT_SOCIAL_CIRCLE      0.000032
DEF_60_CNT_SOCIAL_CIRCLE      YEARS_LAST_PHONE_CHANGE      0.000032
REG_CITY_NOT_WORK_CITY      DEF_30_CNT_SOCIAL_CIRCLE      0.000386
AMT_CREDIT      AMT_GOODS_PRICE      0.986965
OBS_60_CNT_SOCIAL_CIRCLE      OBS_30_CNT_SOCIAL_CIRCLE      0.998511
OBS_30_CNT_SOCIAL_CIRCLE      OBS_60_CNT_SOCIAL_CIRCLE      0.998512
FLAG_EMP_PHONE      YEARS_EMPLOYED      0.999755
YEARS_EMPLOYED      FLAG_EMP_PHONE      0.999755
Length: 992, dtype: float64

In [142]: # Top 10 correlation for target0
## As observed we can see our estimated variables are under top10 correlated variables
target_corr0.tail(10)

Out[142]: CNT_FAM_MEMBERS      CNT_CHILDREN      0.878663
CNT_CHILDREN      CNT_FAM_MEMBERS      0.878663
REGION_RATING_CLIENT      REGION_RATING_CLIENT_W_CITY      0.949976
REGION_RATING_CLIENT_W_CITY      REGION_RATING_CLIENT      0.949976
AMT_GOODS_PRICE      AMT_CREDIT      0.986965
AMT_CREDIT      AMT_GOODS_PRICE      0.986965
OBS_60_CNT_SOCIAL_CIRCLE      OBS_30_CNT_SOCIAL_CIRCLE      0.998511
OBS_30_CNT_SOCIAL_CIRCLE      OBS_60_CNT_SOCIAL_CIRCLE      0.998512
FLAG_EMP_PHONE      YEARS_EMPLOYED      0.999755
YEARS_EMPLOYED      FLAG_EMP_PHONE      0.999755
dtype: float64

In [143]: target_corr1=target_corr1.abs().unstack().sort_values().dropna()
target_corr1
target_corr0

Out[143]: REGION_RATING_CLIENT_W_CITY      YEARS_ID_PUBLISH      0.000105
YEARS_ID_PUBLISH      REGION_RATING_CLIENT_W_CITY      0.000105
DEF_30_CNT_SOCIAL_CIRCLE      YEARS_LAST_PHONE_CHANGE      0.000213
YEARS_LAST_PHONE_CHANGE      DEF_30_CNT_SOCIAL_CIRCLE      0.000213
FLAG_EMP_PHONE      EXT_SOURCE_2      0.000331
AMT_GOODS_PRICE      AMT_CREDIT      0.982846
OBS_30_CNT_SOCIAL_CIRCLE      OBS_60_CNT_SOCIAL_CIRCLE      0.998272
OBS_60_CNT_SOCIAL_CIRCLE      OBS_30_CNT_SOCIAL_CIRCLE      0.998272
YEARS_EMPLOYED      FLAG_EMP_PHONE      0.999704
FLAG_EMP_PHONE      YEARS_EMPLOYED      0.999704
Length: 992, dtype: float64

In [144]: # Top 10 correlation for target1
## As observed we can see our estimated variables are under top10 correlated variables
target_corr1.tail(10)

Out[144]: CNT_FAM_MEMBERS      CNT_FAM_MEMBERS      0.885495
CNT_CHILDREN      CNT_FAM_MEMBERS      0.885495
REGION_RATING_CLIENT      REGION_RATING_CLIENT      0.956560
REGION_RATING_CLIENT_W_CITY      REGION_RATING_CLIENT_W_CITY      0.956560
AMT_CREDIT      AMT_GOODS_PRICE      0.982846
AMT_GOODS_PRICE      AMT_CREDIT      0.982846
OBS_60_CNT_SOCIAL_CIRCLE      OBS_30_CNT_SOCIAL_CIRCLE      0.998272
OBS_30_CNT_SOCIAL_CIRCLE      OBS_60_CNT_SOCIAL_CIRCLE      0.998272
YEARS_EMPLOYED      FLAG_EMP_PHONE      0.999704
FLAG_EMP_PHONE      YEARS_EMPLOYED      0.999704
dtype: float64
```

### Bivariate Analysis

```
In [145]: plt.figure(figsize=(20,10))
plt.subplot(1,2,1)
sns.scatterplot(target0.AMT_ANNUITY,target0.AMT_INCOME_TOTAL)
plt.title('INCOME VS ANNUITY for Target 0',fontsize=20)
plt.yscale('log')
plt.xlabel('Income')
plt.ylabel('Income')
plt.subplot(1,2,2)
sns.scatterplot(target1.AMT_ANNUITY,target1.AMT_INCOME_TOTAL)
plt.title('INCOME VS ANNUITY for Target 1',fontsize=20)
plt.yscale('log')
plt.xlabel('Income')
plt.ylabel('Income')
plt.show()

INCOME VS ANNUITY for Target 0
INCOME VS ANNUITY for Target 1

In [146]: plt.figure(figsize=(20,10))
ax = plt.subplot(2,1,1)
sns.scatterplot(target0.AMT_CREDIT,target0.AMT_GOODS_PRICE)
plt.title('GOODS PRICE VS CREDIT for Target 0',fontsize=20)
plt.xlabel('Credit')
ax.xaxis.get_label().set_fontsize(18)
plt.ylabel('Goods Price')
ax = plt.subplot(2,1,2)
sns.scatterplot(target1.AMT_CREDIT,target1.AMT_GOODS_PRICE)
plt.title('GOODS PRICE VS CREDIT for Target 1',fontsize=20)
plt.xlabel('Credit')
ax.xaxis.get_label().set_fontsize(18)
plt.ylabel('Goods Price')
plt.show()

GOODS PRICE VS CREDIT for Target 0
GOODS PRICE VS CREDIT for Target 1
```

### Goods\_price and Credit are highly correlated

### Previous application data

```
In [147]: prev_data.head()

Out[147]: SK_ID_CURR  SK_ID_PREV  NAME_CONTRACT_TYPE  AMT_ANNUITY  AMT_APPLICATION  AMT_CREDIT  AMT_DOWN_PAYMENT  AMT_C
0  2030495      217877      Consumer loans      1730.430      17145.0      17145.0      0.0
1  2824225      108129      Cash loans      25688.615      617500.0      679671.0      NaN
2  2523466      122040      Cash loans      15100.735      112000.0      136444.5      NaN
3  2819243      176158      Cash loans      47041.335      450000.0      47070.0      NaN
4  1784265      202054      Cash loans      31924.395      337500.0      404035.0      NaN
5 rows x 8 columns

In [148]: pd.set_option('display.max_rows', prev_data.shape[1])
prev_data.describe().T

Out[148]:
```

		count	mean	std	min	25%	50%	75%
SK_ID_PREV	1670214.0	1.923098e+06	532957.268096	1.00000e+06	1.461857e+06	1.92310e+06	2.38420e+06	2.86820e+06
	1670214.0	2.785572e+05	102814.823849	1.00000e+05	1.893290e+05	2.78714e+05	3.675140e+05	4.68240e+05
AMT_ANNUITY	1672979.0	1.959512e+04	14782.137335	0.00000e+00	6.321780e+03	1.125000e+04	2.05842e+04	2.65842e+04
	1672979.0	1.575339e+05	292775.762387	0.00000e+00	1.827200e+04	7.104800e+04	1.805800e+05	6.180580e+05
AMT_APPLICATION	1670213.0	1.961140e+05	318574.616546	0.00000e+00	2.418950e+04	8.054100e+04	2.164185e+05	3.740000e+05
	1670213.0	6.697402e+03	20921.459410	-9.00000e-01	0.00000e+00	1.638000e+03	7.440000e+03	1.500000e+04
AMT_DOWN_PAYMENT	1734730.0	2.278473e+05	313366.557937	0.00000e+00	5.084100e+04	1.123300e+05	2.340000e+05	3.630000e+05
	1734730.0	1.248418e+01	0.334028	0.00000e+00	1.000000e+01	1.200000e+01	1.500000e+01	1.500000e+01
HOUR_APPR_PROCESS_START	1670214.0	9.964678e-01	0.056330	0.00000e+00	1.000000e+00	1.000000e+00	1.000000e+00	1.000000e+00
	1734730.0	7.963882e-02	0.107823	-1.497876e-02	0.000000e+00	5.160506e-02	1.080991e-01	1.080991e-01
RATE_DOWN_PAYMENT	774370.0	1.883586e-01	0.087671	3.478125e-02	1.607163e-01	1.891222e-01	1.932399e-01	1.932399e-01
	774370.0	1.735025e-01	0.100879	3.731501e-01	1.715644e-01	8.350951e-01	8.525370e-01	8.525370e-01
DAYS_DECISION	1670213.0	8.806797e+02	779.090967	-2.920000e+03	-1.300000e+03	-6.810000e+02	2.800000e+02	2.800000e+02
	1670213.0	1.336119e+02	1197.443459	-1.000000e+00	-1.000000e+00	3.000000e+00	8.200000e+00	8.200000e+00
SELLERPLACE_AREA	1672979.0	1.318616e+01	1.676728	0.000000e+00	6.000000e+00	1.200000e+01	2.400000e+01	2.400000e+01
	1672979.0	2.609540e+01	48.9917	-1.158354	-2.920000e+03	3.652430e+03	3.652430e+03	3.652430e+03
DAYS_FIRST_DRAWING	997149.0	1.622985e+03	1697.149	-2.870000e+03	-1.628000e+03	-8.310000e+02	-4.110000e+02	-4.110000e+02
	997149.0	1.382677e+04	72444.869708	-1.628000e+03	-1.628000e+03	-8.310000e+02	-4.110000e+02	-4.110000e+02
DAYS_LAST_DUE_1ST_VERSION	997149.0	3.370276e+04	104687.034789	-2.870000e+03	-1.342000e+03	-3.670000e+02	-2.290000e+01	-2.290000e+01
	997149.0	1.658234e+04	149647.415123	-2.870000e+03	-1.342000e+03	-4.960000e+02	-4.400000e+01	-4.400000e+01
DAYS_TERMINATION	997149.0	8.190234e+04	153303.516729	-2.870000e+03	-1.270000e+03	-4.960000e+02	-4.400000e+01	-4.400000e+01
	997149.0	1.325702e-01	0.471134	0.000000e+00	0.000000e+00	0.000000e+00	1.000000e+00	1.000000e+00

```
In [149]: perc_missing = round((prev_data.isna().sum() / prev_data.shape[0], 2)
#number of columns with missing values
perc_missing[perc_missing != 0].size

Out[149]: 23

In [150]: #number of columns with missing values
perc_missing[perc_missing != 0].size

Out[150]: 14

In [151]: perc_missing[perc_missing>0.5]

Out[151]: RATE_DOWN_PAYMENT      0.54
RATE_DOWN_PAYMENT_PRIMARY      0.54
RATE_INTEREST_PRIMARY      1.00
RATE_INTEREST_PRIVILEGED      1.00
dtype: float64

In [152]: # As thereare morethan 50% missing values in AMT_DOWN_PAYMENT,RATE_DOWN_PAYMENT
# We cannot use these cols for analysis

In [ ]:

In [153]: prev_data['SK_ID_CURR'].value_counts()

Out[153]: 187868      77
265681      73
173680      72
242412      68
206783      67
382489      1
426056      1
454726      1
381442      1
124145      1
Name: SK_ID_CURR, Length: 338857, dtype: int64

In [154]: ## Few appl numbers are repeated --> treating them as fraud

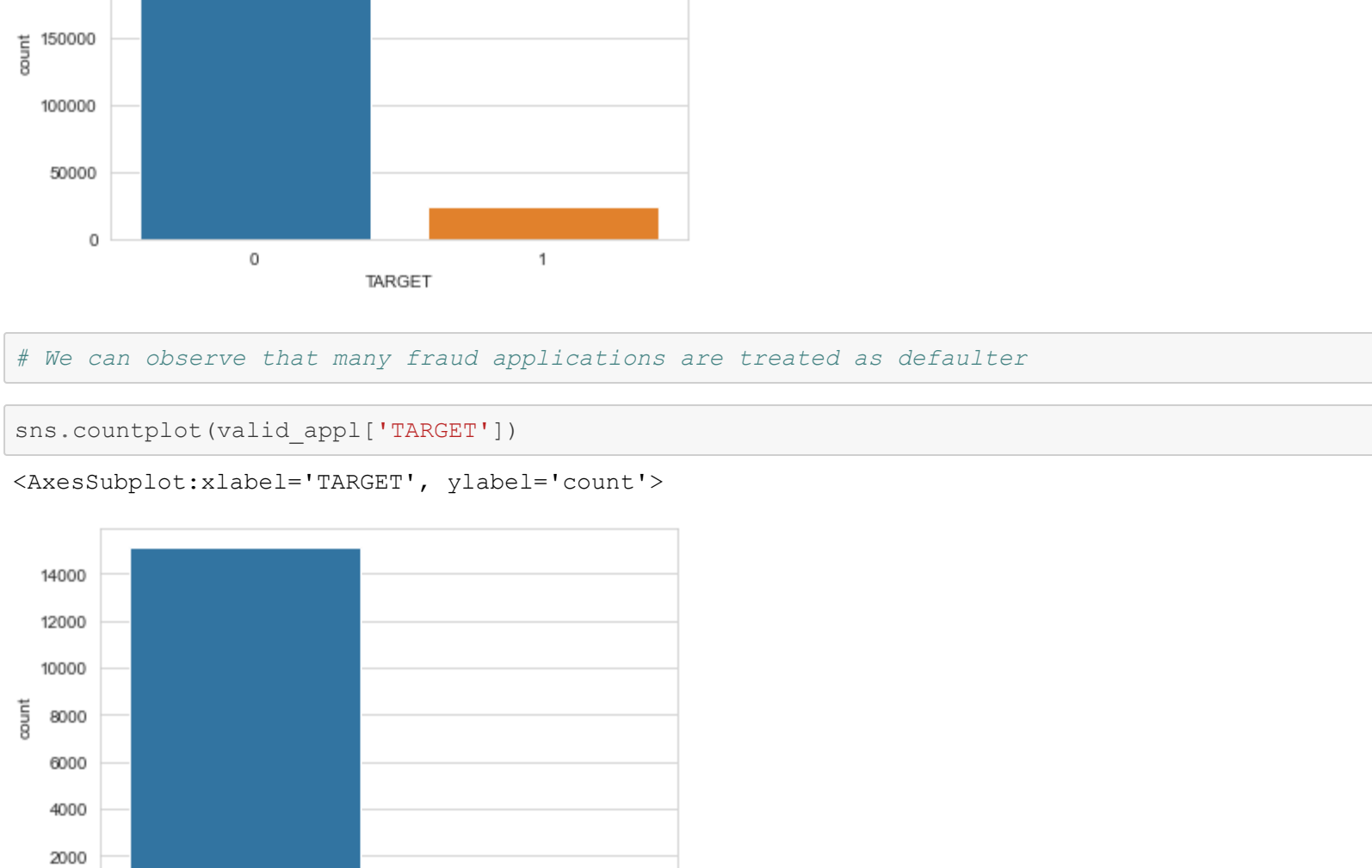
In [155]: fraud_id=[prev_data['SK_ID_CURR'].value_counts()[1:].index

In [156]: valid_appl=appl_data[~appl_data['SK_ID_CURR'].isin(fraud_id)]

In [157]: fraud_appl=appl_data[appl_data['SK_ID_CURR'].isin(fraud_id)]

In [158]: sns.countplot(fraud_appl['TARGET'])

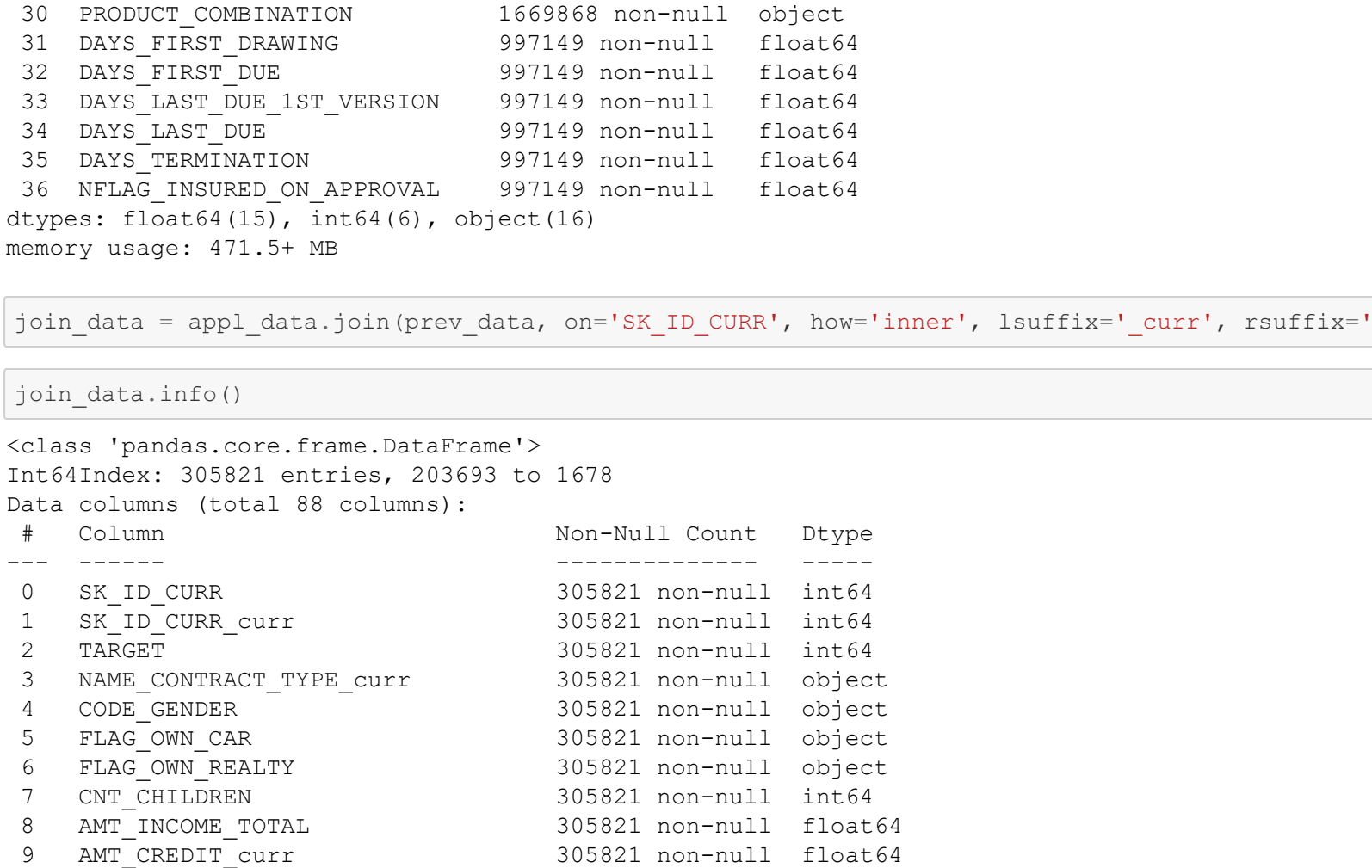
Out[158]: <AxesSubplot:label='TARGET', ylabel='count'>
```



```
In [159]: # We can observe that many fraud applications are treated as defaulter

In [160]: sns.countplot(valid_appl['TARGET'])

Out[160]: <AxesSubplot:label='TARGET', ylabel='count'>
```



```
In [161]: prev_data.info()

Out[161]: <class 'pandas.core.frame.DataFrame'>
RangeIndex: 1670214 entries, 0 to 1670213
Data columns (total 37 columns):
#   Column                Non-Null Count  Dtype
---  --
0   SK_ID_CURR             1670214 non-null   int64
1   SK_ID_PREV             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        1670214 non-null   float64
7   AMT_GOODS_PRICE         1670214 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_PRIVILEGED  5951 non-null     float64
15  NAME_CASH_LOAN_PURPOSE    1670214 non-null   object
16  NAME_CONTRACT_STATUS      1670214 non-null   object
17  DAYS_DECISION             1670214 non-null   int64
18  NAME_CLIENT_TYPE          1670214 non-null   object
19  CODE_REJECT_REASON         1670214 non-null   object
20  NAME_TYPE_SUITE           849809 non-null    object
21  NAME_CLIENT_TYPE          1670214 non-null   object
22  NAME_GOODS_CATEGORY        1670214 non-null   object
23  NAME_PORTFOLIO             1670214 non-null   object
24  NAME_PRODUCT_TYPE          1670214 non-null   object
25  CHANNEL_TYPE               1670214 non-null   object
26  SELLERPLACE_AREA           1670214 non-null   int64
27  NAME_SELLER_INDUSTRY       1670214 non-null   object
28  CNT_PAYMENT                1297984 non-null   float64
29  NAME_YIELD_GROUP           1670214 non-null   object
30  PRODUCT_COMBINATION        1669688 non-null   object
31  DAYS_FIRST_DRAWING         997149 non-null    float64
32  DAYS_FIRST_DUE            997149 non-null    float64
33  DAYS_LAST_DUE_1ST_VERSION  997149 non-null    float64
34  DAYS_LAST_DUE             997149 non-null    float64
35  DAYS_TERMINATION           997149 non-null    float64
36  NFLAG_INSURED_ON_APPROVAL  997149 non-null    float64
dtypes: float64(15), int64(6), object(16)
memory usage: 471.5+ MB

In [162]: join_data = appl_data.join(prev_data, on='SK_ID_CURR', how='inner', lsuffix='_curr', rsuffix='_prev')

In [163]: join_data.info()

Out[163]: <class 'pandas.core.frame.DataFrame'>
RangeIndex: 305821 entries, 203693 to 1678
Data columns (total 88 columns):
#   Column                Non-Null Count  Dtype
---  --
0   SK_ID_CURR             305821 non-null   int64
1   SK_ID_CURR_curr        305821 non-null   int64
2   TARGET                 305821 non-null   int64
3   NAME_CONTRACT_TYPE_curr  305821 non-null   object
4   CODE_GENDER            305821 non-null   object
5   FLAG_OWN_CAR            305821 non-null   object
6   FLAG_OWN_REALTY         305821 non-null   object
7   CNT_CHILDREN            305821 non-null   int64
8   AMT_INCOME_TOTAL_curr   305821 non-null   float64
9   AMT_CREDIT_curr         305821 non-null   float64
10  AMT_ANNUITY_curr         305821 non-null   float64
11  AMT_GOODS_PRICE_curr     305821 non-null   float64
12  NAME_TYPE_SUITE_curr     305821 non-null   object
13  NAME_INCOME_TYPE         305821 non-null   object
14  NAME_EDUCATION_TYPE      305821 non-null   object
15  NAME_FAMILY_STATUS       305821 non-null   object
16  NAME_HOUSING_TYPE        305821 non-null   object
17  REGION_POPULATION_RELATIVE  305821 non-null   float64
18  FLAG_MOBIL              305821 non-null   int64
19  FLAG_EMP_PHONE          305821 non-null   int64
20  FLAG_WORK_PHONE          305821 non-null   int64
21  FLAG_CONTACT_MOBILE      305821 non-null   int64
22  FLAG_PHONE              305821 non-null   int64
23  FLAG_EMAIL              305821 non-null   int64
24  OCCUPATION_TYPE          305821 non-null   object
25  CNT_FAM_MEMBERS          305821 non-null   int64
26  REGION_RATING_CLIENT      305821 non-null   int64
27  REGION_RATING_CLIENT_W_CITY  305821 non-null   int64
28  WEEKDAY_APPR_PROCESS_START_curr  305821 non-null   object
29  HOUR_APPR_PROCESS_START  305821 non-null   int64
30  REG_REGION_NOT_LIVE_REGION  305821 non-null   object
31  LIVE_REGION_NOT_WORK_REGION  305821 non-null   int64
32  REG_CITY_NOT_LIVE_CITY   305821 non-null   int64
33  LIVE_CITY_NOT_WORK_CITY  305821 non-null   int64
34  ORGANIZATION_TYPE         305821 non-null   object
35  EXT_SOURCE_2              305821 non-null   float64
36  EXT_SOURCE_3             243357 non-null   float64
37  OBS_30_CNT_SOCIAL_CIRCLE  305821 non-null   float64
38  OBS_60_CNT_SOCIAL_CIRCLE  305821 non-null   float64
39  DEF_60_CNT_SOCIAL_CIRCLE  305821 non-null   float64
40  AGE                      305821 non-null   float64
41  YEARS_EMPLOYED            305821 non-null   float64
42  REGISTRATION              305821 non-null   float64
43  YEARS_ID_PUBLISH         305821 non-null   float64
44  YEARS_LAST_PHONE_CHANGE   305821 non-null   float64
45  Annuity_Class             305821 non-null   category
46  Income_Class              305821 non-null   category
47  INCOME_TYPE               305821 non-null   object
48  SK_ID_CURR_prev          305821 non-null   int64
49  SK_ID_CURR_curr          305821 non-null   int64
50  NAME_CONTRACT_TYPE_prev   305821 non-null   object
51  AMT_APPLICATION_prev      305821 non-null   float64
52  AMT_CREDIT_prev          305821 non-null   float64
53  AMT_DOWN_PAYMENT_prev     305821 non-null   float64
54  AMT_GOODS_PRICE_prev      305821 non-null   float64
55  WEEKDAY_APPR_PROCESS_START_prev  305821 non-null   object
56  HOUR_APPR_PROCESS_START_prev  305821 non-null   int64
57  FLAG_LAST_APPL_PER_CONTRACT_prev  305821 non-null   object
58  NFLAG_LAST_APPL_IN_DAY_prev  305821 non-null   int64
59  RATE_DOWN_PAYMENT_prev    141866 non-null    float64
60  RATE_INTEREST_PRIMARY_prev  1091 non-null     float64
61  RATE_INTEREST_PRIVILEGED_prev  1091 non-null     float64
62  NAME_CASH_LOAN_PURPOSE_prev  305821 non-null   object
63  NAME_CONTRACT_STATUS_prev  305821 non-null   object
64  DAYS_DECISION_prev        305821 non-null   int64
65  NAME_PAYMENT_TYPE         305821 non-null   object
66  CODE_REJECT_REASON_prev   305821 non-null   object
67  NAME_TYPE_SUITE_prev      305821 non-null   object
68  NAME_CLIENT_TYPE_prev     305821 non-null   object
69  NAME_GOODS_CATEGORY_prev   305821 non-null   object
70  NAME_PORTFOLIO_prev        305821 non-null   object
71  NAME_PRODUCT_TYPE_prev    305821 non-null   object
72  CHANNEL_TYPE_prev          305821 non-null   object
73  SELLERPLACE_AREA_prev     305821 non-null   int64
74  NAME_SELLER_INDUSTRY_prev  305821 non-null   object
75  CNT_PAYMENT_prev          237563 non-null   float64
76  NAME_YIELD_GROUP_prev      305821 non-null   object
77  PRODUCT_COMBINATION_prev   305821 non-null   object
78  DAYS_FIRST_DRAWING_prev    182215 non-null    float64
79  DAYS_FIRST_DUE_prev        182215 non-null    float64
80  DAYS_LAST_DUE_1ST_VERSION_prev  182215 non-null    float64
81  DAYS_LAST_DUE_prev         182215 non-null    float64
82  DAYS_TERMINATION_prev      182215 non-null    float64
83  FLAG_INSURED_ON_APPROVAL_prev  182215 non-null    float64
dtypes: category(15), int64(6), object(29)
memory usage: 203.6+ MB
```

```
In [164]: join_data.loc[join_data['TARGET'] > 0]

Out[164]:
```

SK_ID_CURR	SK_ID_CURR_curr	TARGET	NAME_CONTRACT_TYPE_curr	CODE_GENDER	FLAG_OWN_CAR	FLAG_OWN_REALTY
248159	387126	1	Cash loans	F	Y	Y
248970	388040	1	Cash loans	M	Y	Y
41725	148308	1	Cash loans	M	Y	Y
167656	294352	1	Cash loans	F	N	Y
198106	327395	1	Cash loans	F	Y	N
...	...	...	...	...	...	...
86070	199880	1	Cash loans	F	Y	N
111197	229012	1	Cash loans	F	N	N
76778	188037	1	Cash loans	F	N	N
15543	121642	1	Cash loans	F	N	N
20727	124157	1	Cash loans	F	Y	Y
24736 rows x 8 columns						

```
In [165]: fill = join_data.loc[join_data['TARGET'] > 0]['SK_ID_CURR'].value_counts()
fill.sum()

Out[165]: 24736

In [166]: ## We can observe that few applications which are repeated in previous data are not present in applicat
ion data
join_data['SK_ID_CURR_prev'].value_counts()

Out[166]: 265681      21
110899      16
205683      16
173680      15
206682      15
345543      ..
368176      1
366129      1
360215      1
278507      1
Name: SK_ID_CURR_prev, Length: 180724, dtype: int64

In [167]: prev_data['SK_ID_CURR'].value_counts()

Out[167]: 187868      77
265681      73
173680      72
242412      68
206783      67
382489      1
426056      1
454726      1
381442      1
124145      1
Name: SK_ID_CURR, Length: 338857, dtype: int64

In [168]: join_data[['SK_ID_CURR', 'TARGET', 'CODE_REJECT_REASON']]

Out[168]:
```

SK_ID_CURR	TARGET	CODE_REJECT_REASON
203693	0	HC
266608	0	LIMIT
77768	199180	XAP
131127	252084	XAP
287463	432980	XAP
...	...	...
240137	378118	XAP
246104	384810	XAP
186643	316377	XAP
20727	124157	XAP
1678	101905	HC
305821 rows x 3 columns		

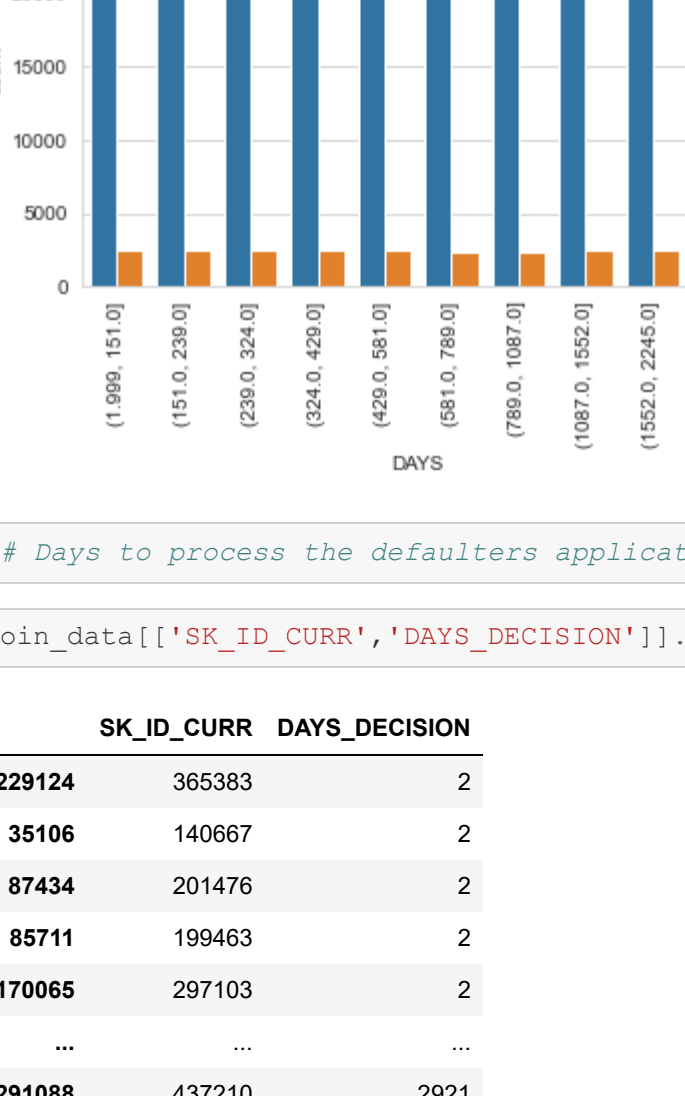
```
In [169]: join_data[['SK_ID_CURR', 'TARGET', 'CODE_REJECT_REASON']].loc[join_data['TARGET'] > 0]

Out[169]:
```

SK_ID_CURR	TARGET	CODE_REJECT_REASON	
248159	387126	1	XAP
24570	388040	1	XAP
41725	148308	1	SCO
167656	294352	1	XAP
198106	327395	1	XAP
...	...	...	...
86070	199880	1	XAP
111197	229012	1	HC
76778	188037	1	XAP
15543	121642	1	XAP
20727	124157	1	XAP
24736 rows x			



```
In [174]: sns.countplot(data=join_data[['TARGET','DAYS']],x='DAYS',hue='TARGET')
plt.xticks(rotation=90)
plt.show()
```



```
In [175]: ## Days to process the defaulters application is more compared to non defaulters
```

```
In [176]: join_data[['SK_ID_CURR','DAYS_DECISION']].loc[join_data['TARGET']==1].sort_values(by='DAYS_DECISION')
```

Out[176]:

SK_ID_CURR	DAYS_DECISION
229124	365383
35106	140667
87434	201476
85711	199463
170065	297103
...	...
291088	437210
199478	331259
293786	440369
193784	324690
223434	358790

24736 rows × 2 columns

```
In [177]: join_data[['SK_ID_CURR','DAYS_DECISION']].loc[join_data['TARGET']==0].sort_values(by='DAYS_DECISION')
```

Out[177]:

SK_ID_CURR	DAYS_DECISION
22839	128576
292252	438560
278341	422476
219630	354440
87937	202087
...	...
34160	139586
135141	256734
115640	234100
184510	313869
91988	208788

281085 rows × 2 columns

## Notes on Analysis

- Pre-liminary analysis indicate that: Some applicants are having multiple previous applications There are more relevant business information like CODE\_REJECT\_REASON which means a direct impact on the Target variable. Previous Applications dataset has more than double the number of records across both the Target segments.