



Instituto Tecnológico
de Buenos Aires

EDA - Credit Fraud

Agustin Lara

—

Analítica Predictiva

29/03/2023

AGENDA

01 Limpieza de Datos

02 Exploración

03 Corrplot

04 Conclusiones

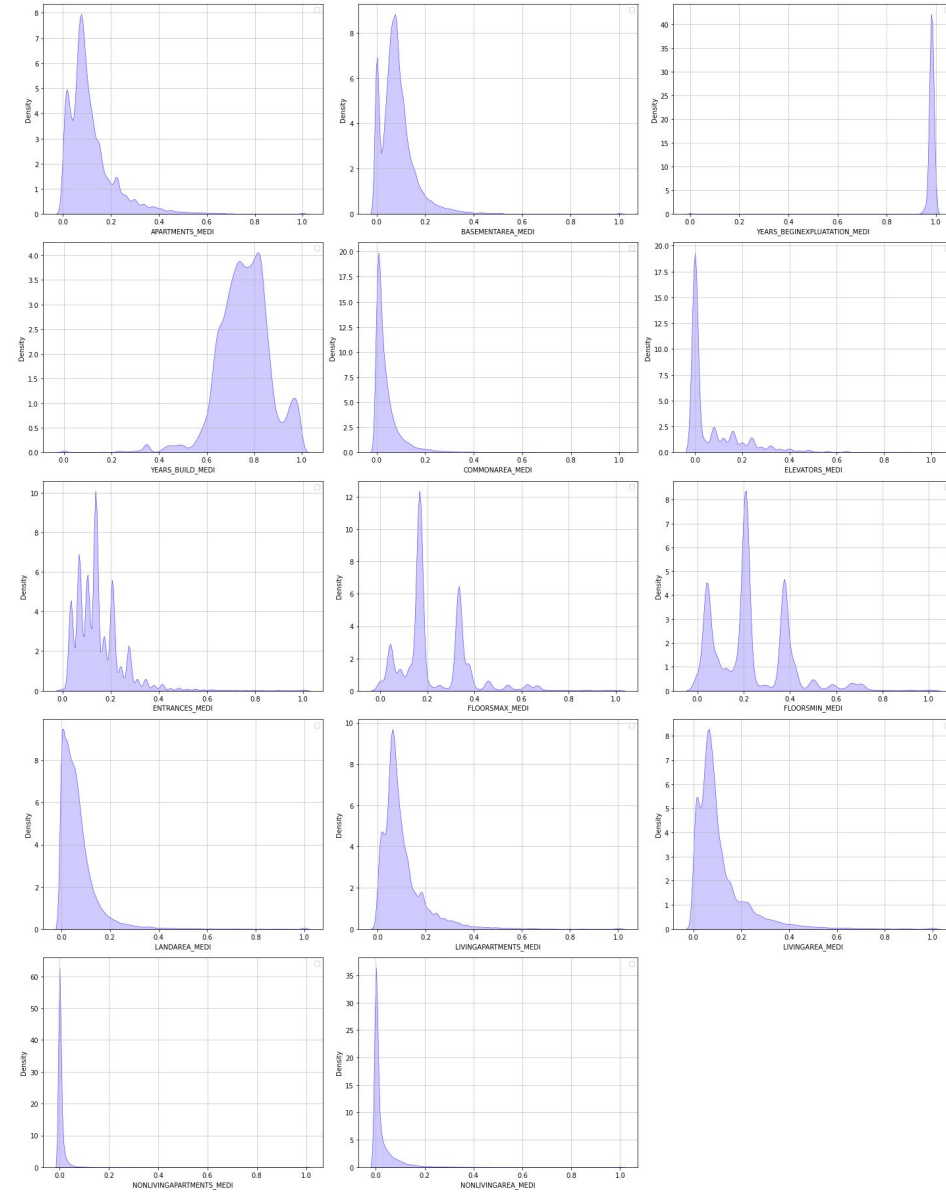
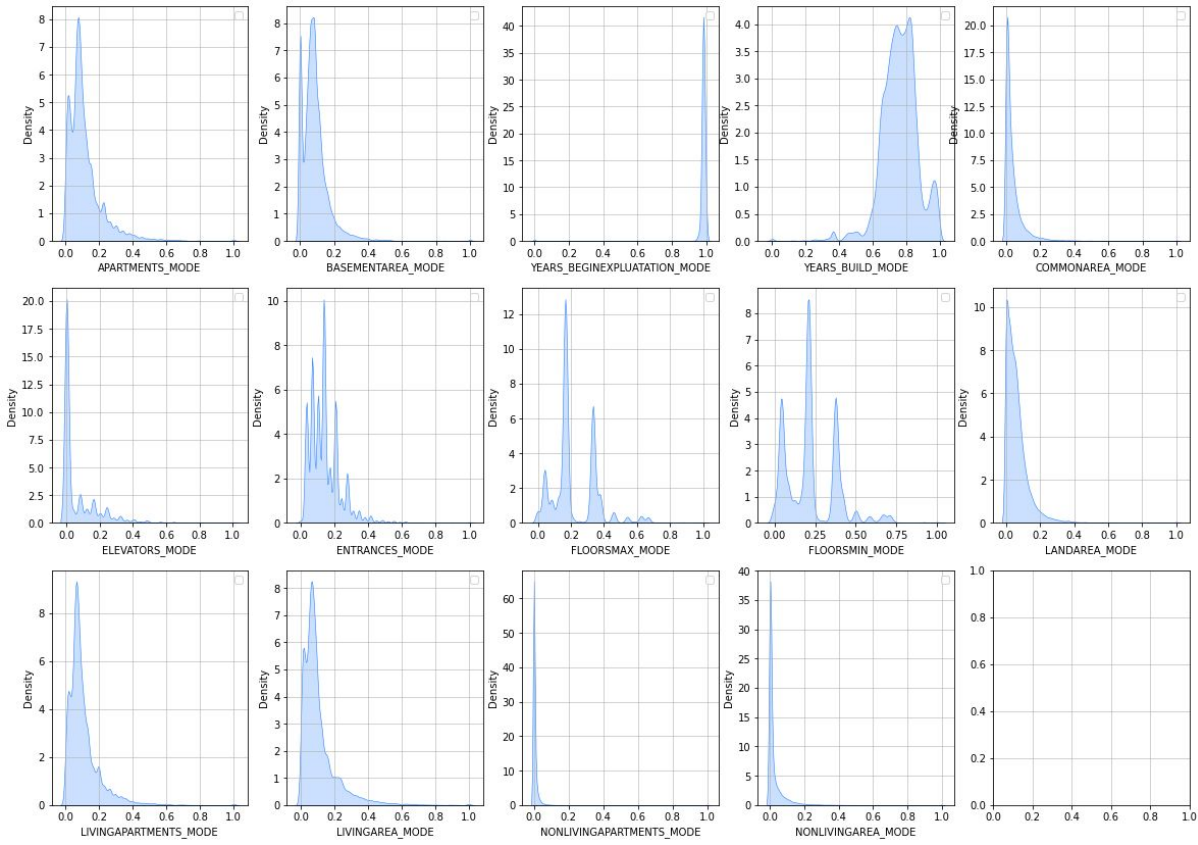
Base: “Credit Card Fraud Detection”

- 307511 registros
- 122 columnas
- dtypes: float64(65), int64(41), object(16)

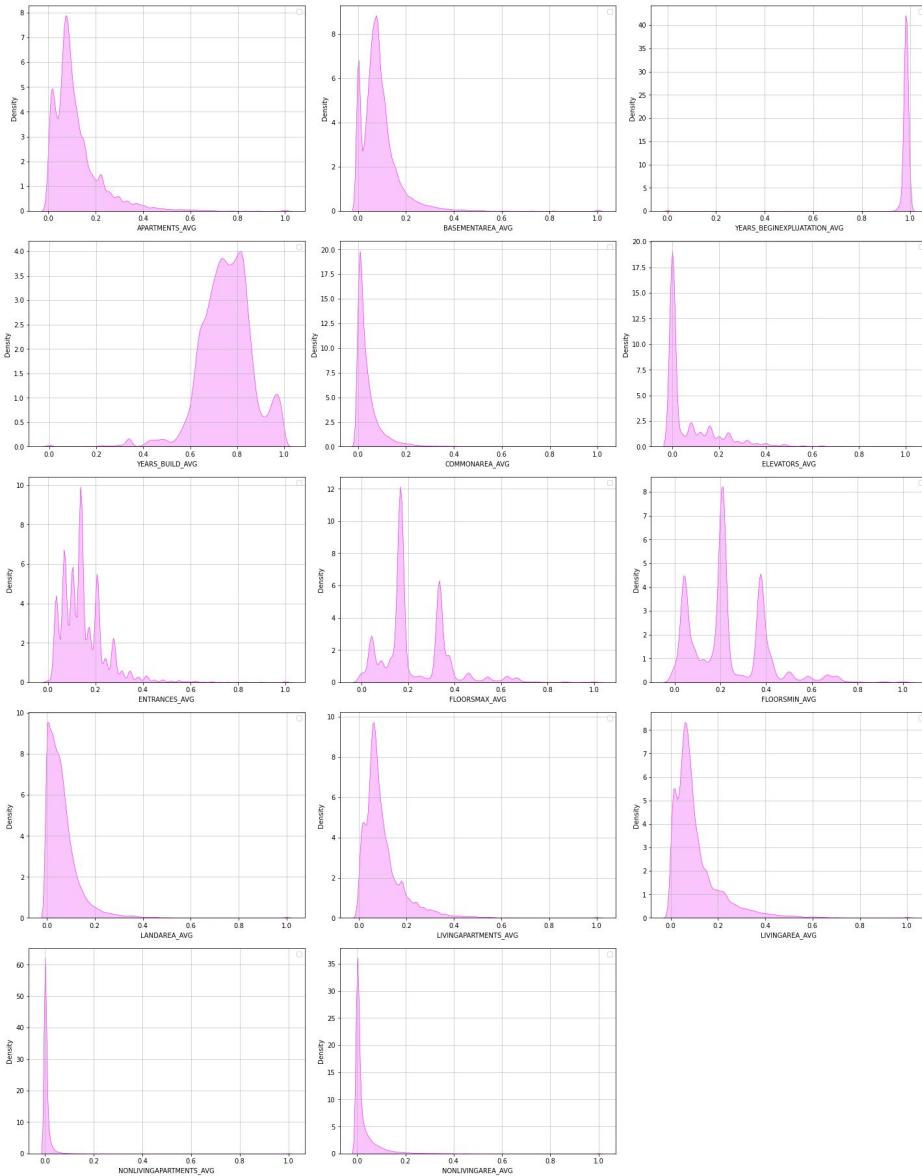
Se busca poder identificar a partir de variables categóricas y numéricas si un cliente se atrasa con los pagos de sus cuotas a la hora de tomar un préstamo



1. Variables normalizadas



2. Cambio de valores



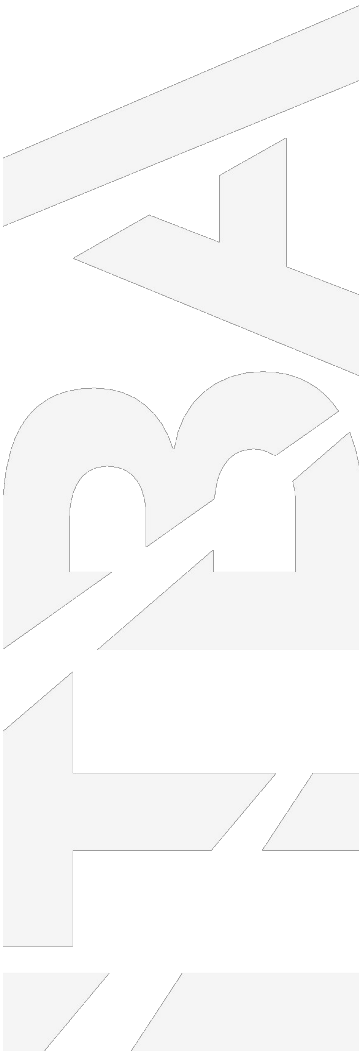
pivot - DataFrame

	FLAG OWN CAR	FLAG OWN REALTY	FLAG MOBIL	FLAG FMP PHONE	WORK PH	CONT MC	FLAG PHONE	FLAG EMAIL	DOCUMEI	DOCUMEI	DOCUMEI	DOCUMEI	DOCUMEI	DOCUMEI	DOCUMEI
0	N	N	1	0	0	0	0	0	0	0	0	0	0	0	0
1	N	N	1	0	0	0	0	0	0	1	0	0	0	0	0
2	N	N	1	0	0	0	1	0	0	1	0	0	0	0	0
3	N	N	1	0	0	1	0	0	0	0	0	0	0	0	0
4	N	N	1	0	0	1	0	0	0	0	0	0	0	0	0
5	N	N	1	0	0	1	0	0	0	0	0	0	0	0	1
6	N	N	1	0	0	1	0	0	0	0	0	0	1	0	0
7	N	N	1	0	0	1	0	0	0	0	0	0	1	0	0
8	N	N	1	0	0	1	0	0	0	0	0	0	1	0	0
9	N	N	1	0	0	1	0	0	0	0	0	0	1	0	0
10	N	N	1	0	0	1	0	0	0	0	0	1	0	0	0
11	N	N	1	0	0	1	0	0	0	1	0	0	0	0	0

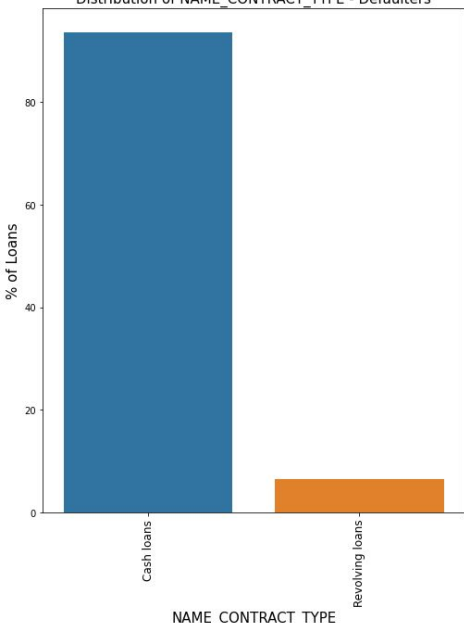
TARGET	FLAG_CONT_MOBILE	FLAG_DOCUMENT_10	FLAG_DOCUMENT_11
N	0.998128666	2.47625E-05	0.00399029
Y	0.998187311	0	0.00302114
DIFF	-5.86455E-05	2.47625E-05	0.00096914

FLAG_DOCUMENT_12	FLAG_DOCUMENT_13	FLAG_DOCUMENT_14	FLAG_DOCUMENT_15	FLAG_DOCUMENT_16
7.07499E-06	0.00	0.3088232172799500	0.12770352971141100	0.001026934
0	0.001208459	0.12084592145015100	0.4431017119838870	0.006042296
7.07499E-06	0.002520059	0.187977296	-0.315398182	-0.005015362

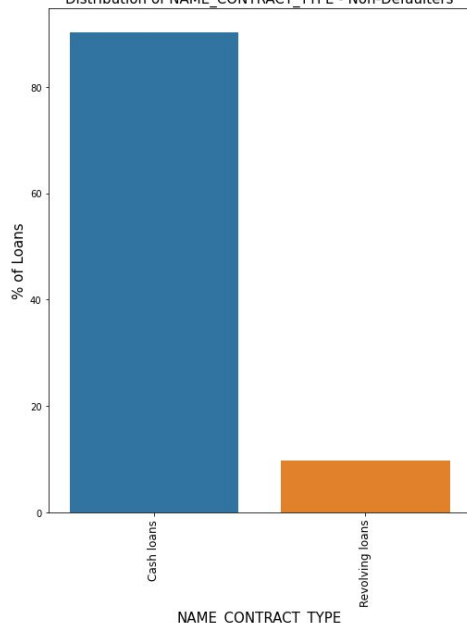
Index	TYPE IN POPULATION REL	DAYS BIRTH	DAYS EMPLOYED	DAYS REGISTRATION	DAYS ID PUBLISH	OWN CAR AGE	OCCUR
0	0.018801	-9461	-637	-3648	-2120	nan	Labore
1	0.003541	-16765	-1188	-1186	-291	nan	Core s
2	0.010032	-19046	-225	-4260	-2531	26	Labore
3	0.008019	-19005	-3039	-9833	-2437	nan	Labore
4	0.028663	-19932	-3038	-4311	-3458	nan	Core s
5	0.035792	-16941	-1588	-4970	-477	nan	Labore
6	0.035792	-13778	-3130	-1213	-619	17	Accoun
7	0.003122	-18850	-449	-4597	-2379	8	Manage
8	0.018634	-20099	365243	-7427	-3514	nan	nan
9	0.019689	-14469	-2019	-14437	-3992	nan	Labore
10	0.0228	-10197	-679	-4427	-738	nan	Core s
11	0.015221	-20417	365243	5246	2512	nan	nan



Distribution of NAME_CONTRACT_TYPE - Defaulters



Distribution of NAME_CONTRACT_TYPE - Non-Defaulters

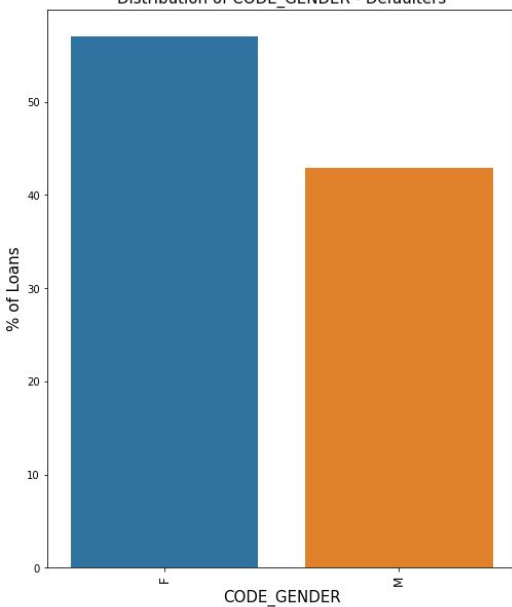


3. Comparación: Defaulters v non Defaulters

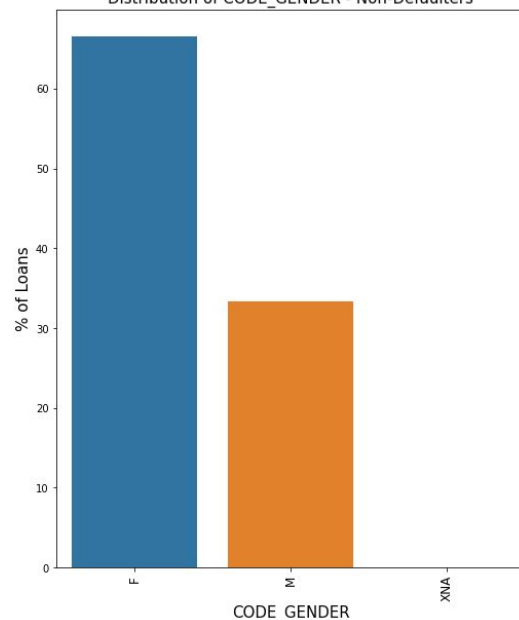
%8.07

Presenta problemas de pago

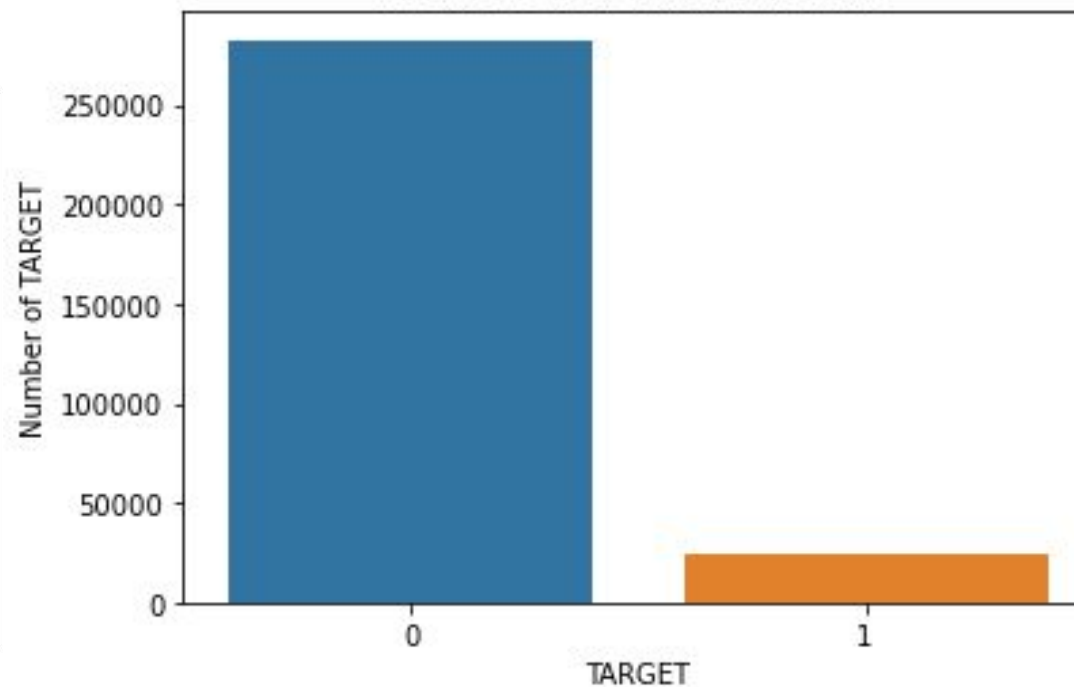
Distribution of CODE_GENDER - Defaulters

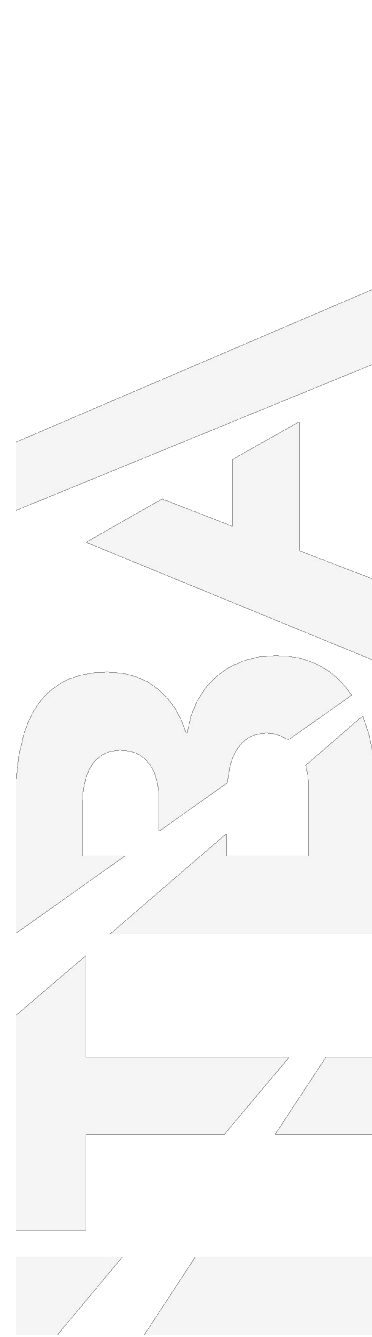
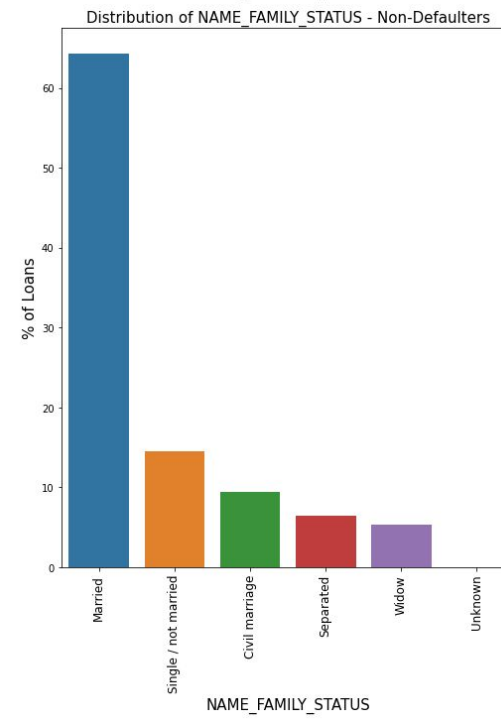
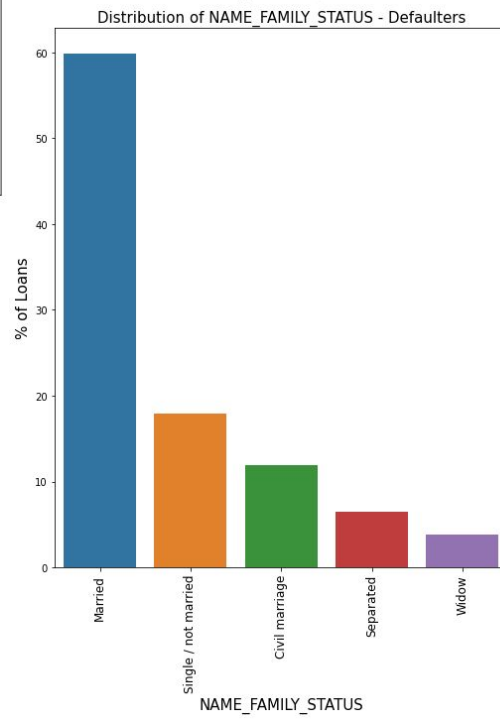
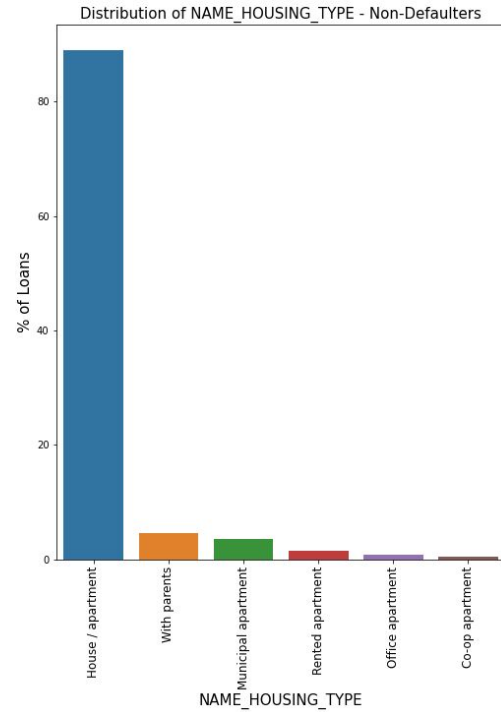
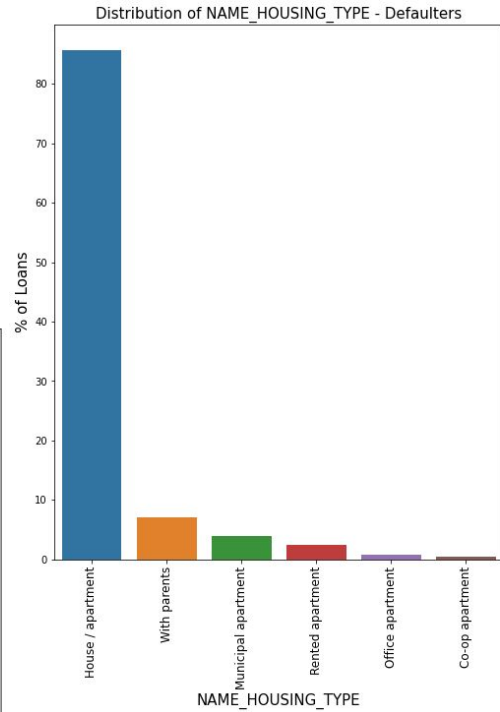
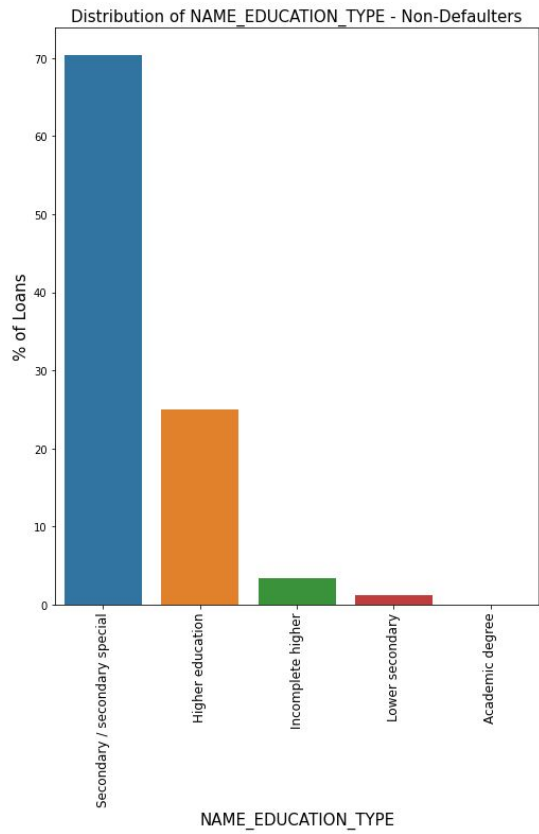
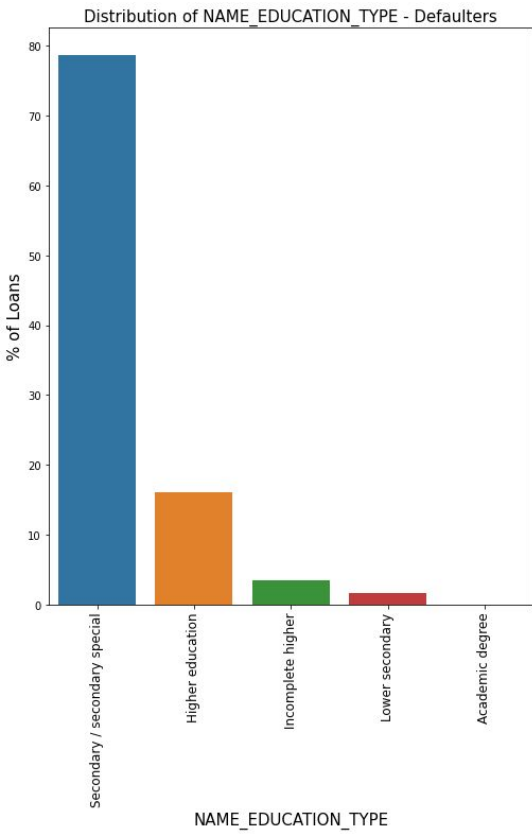


Distribution of CODE_GENDER - Non-Defaulters



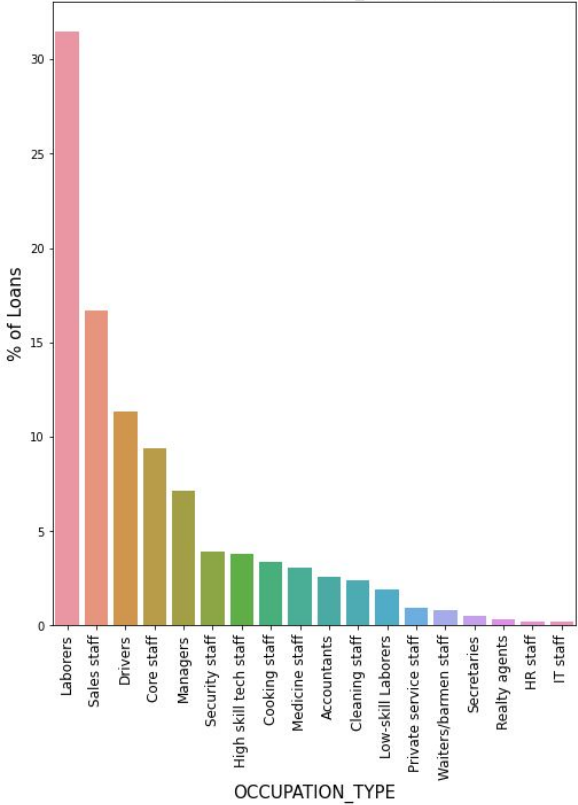
Distribution of TARGET Variable



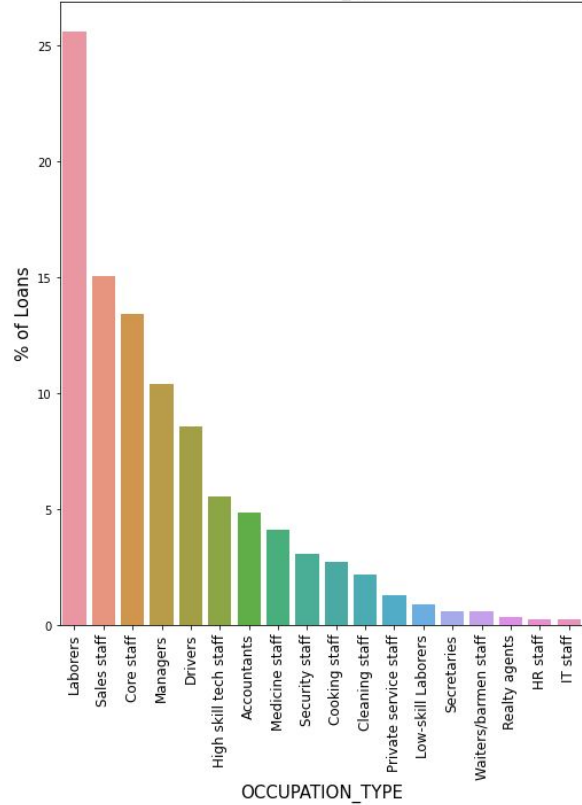




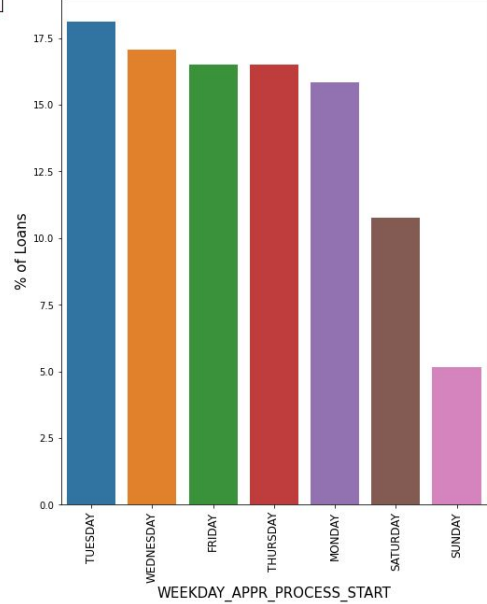
Distribution of OCCUPATION_TYPE - Defaulters



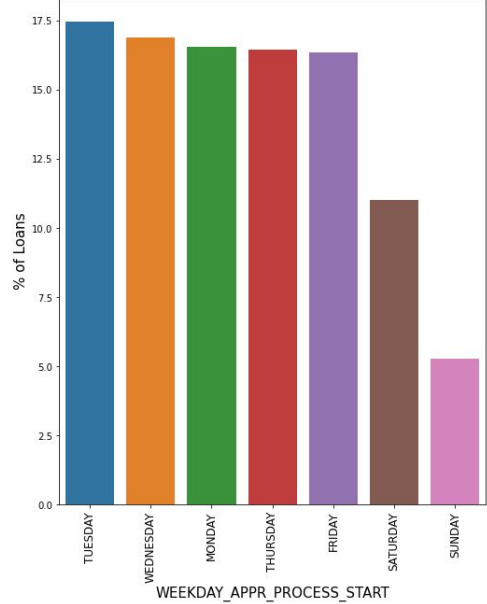
Distribution of OCCUPATION_TYPE - Non-Defaulters

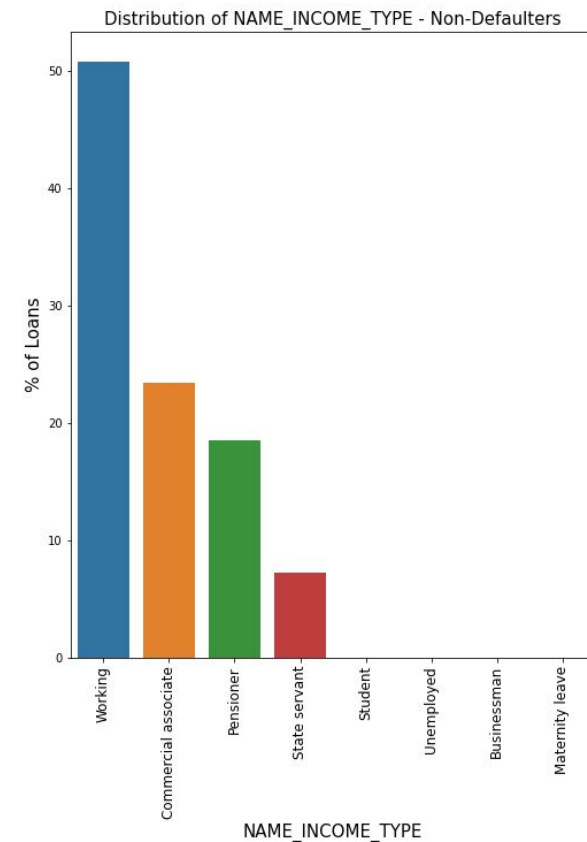
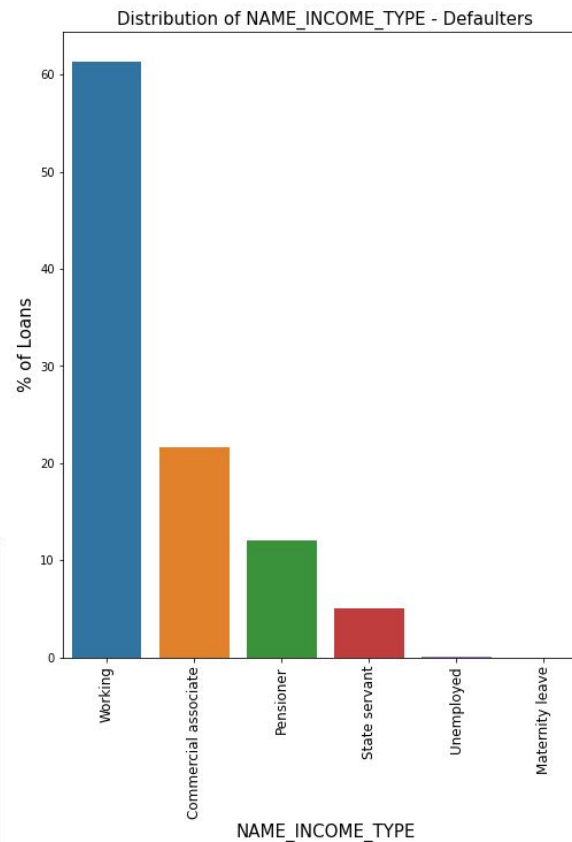
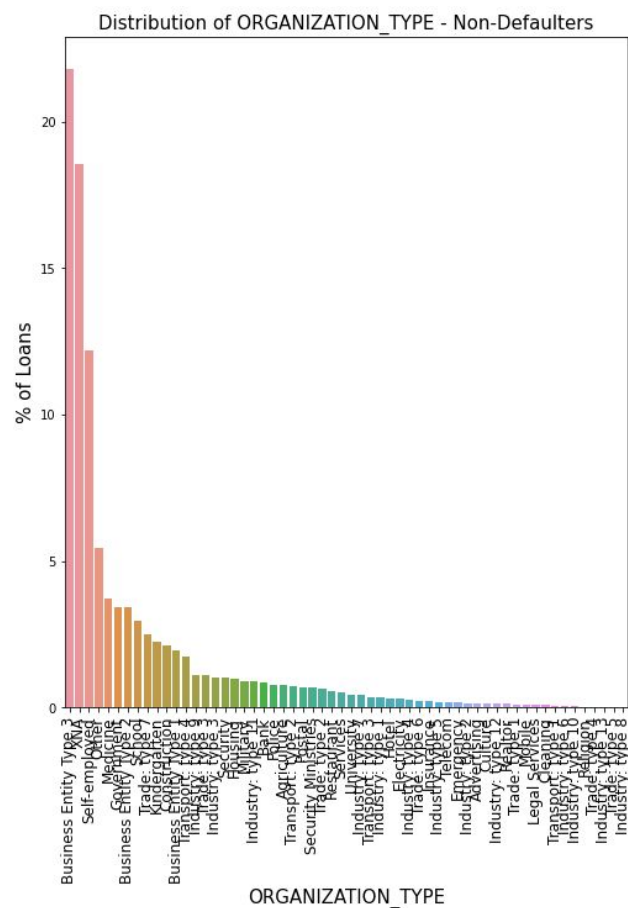
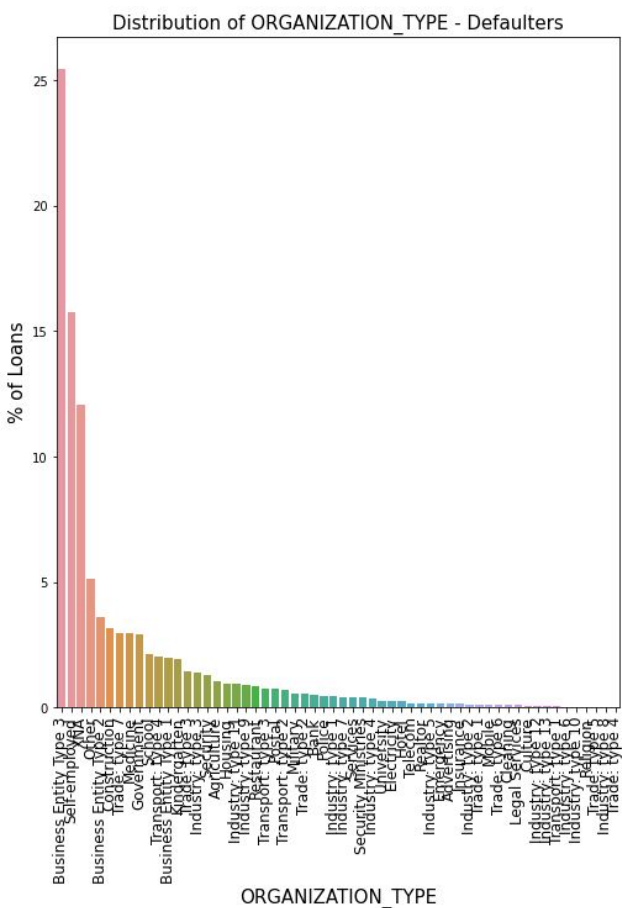


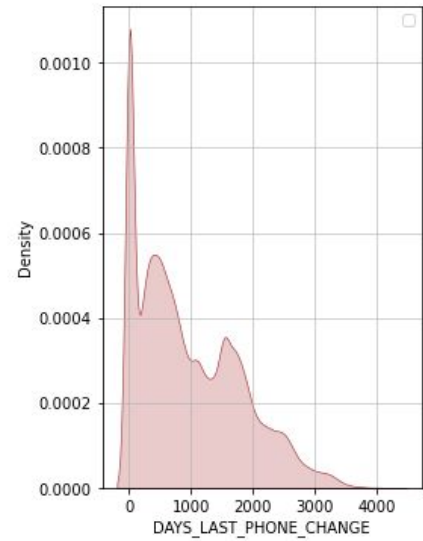
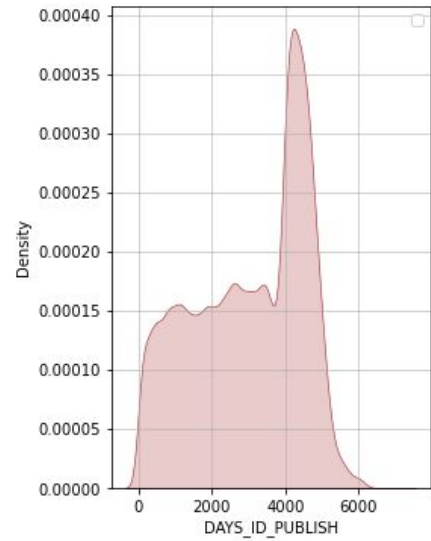
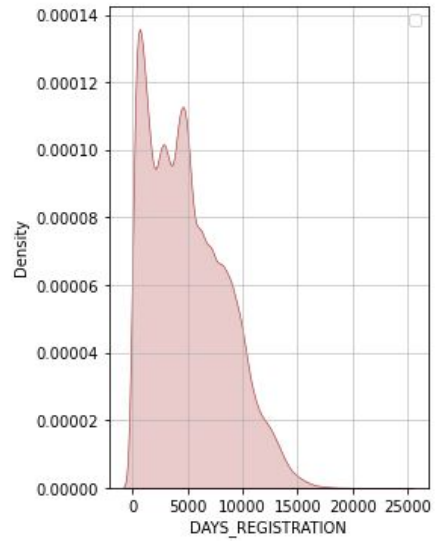
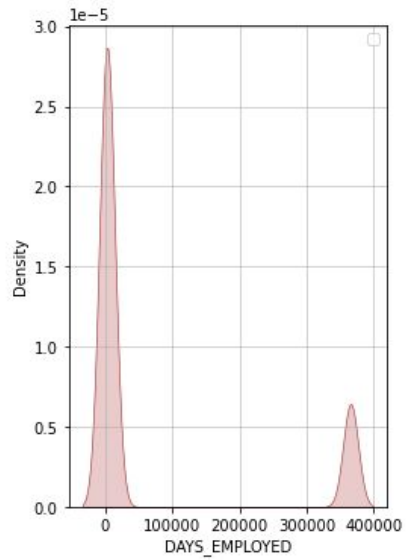
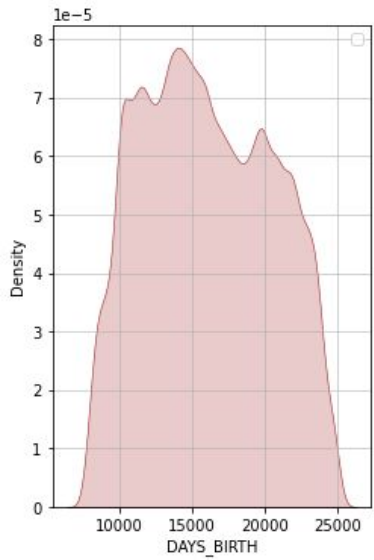
Distribution of WEEKDAY_APPR_PROCESS_START - Defaulters

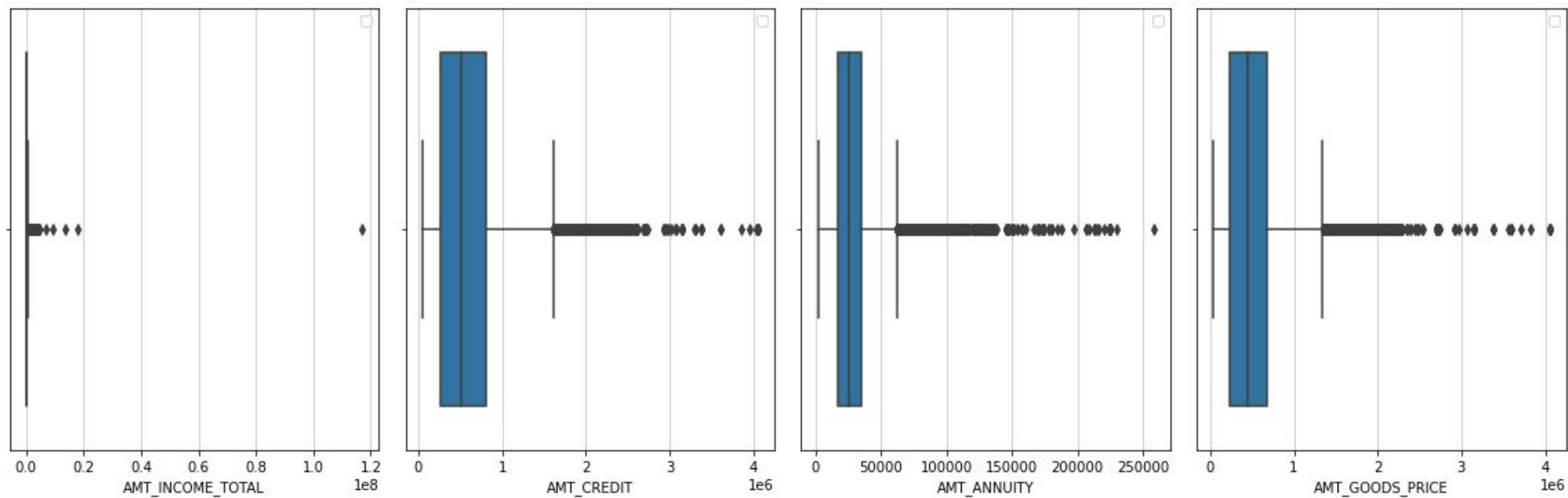


Distribution of WEEKDAY_APPR_PROCESS_START - Non-Defaulters



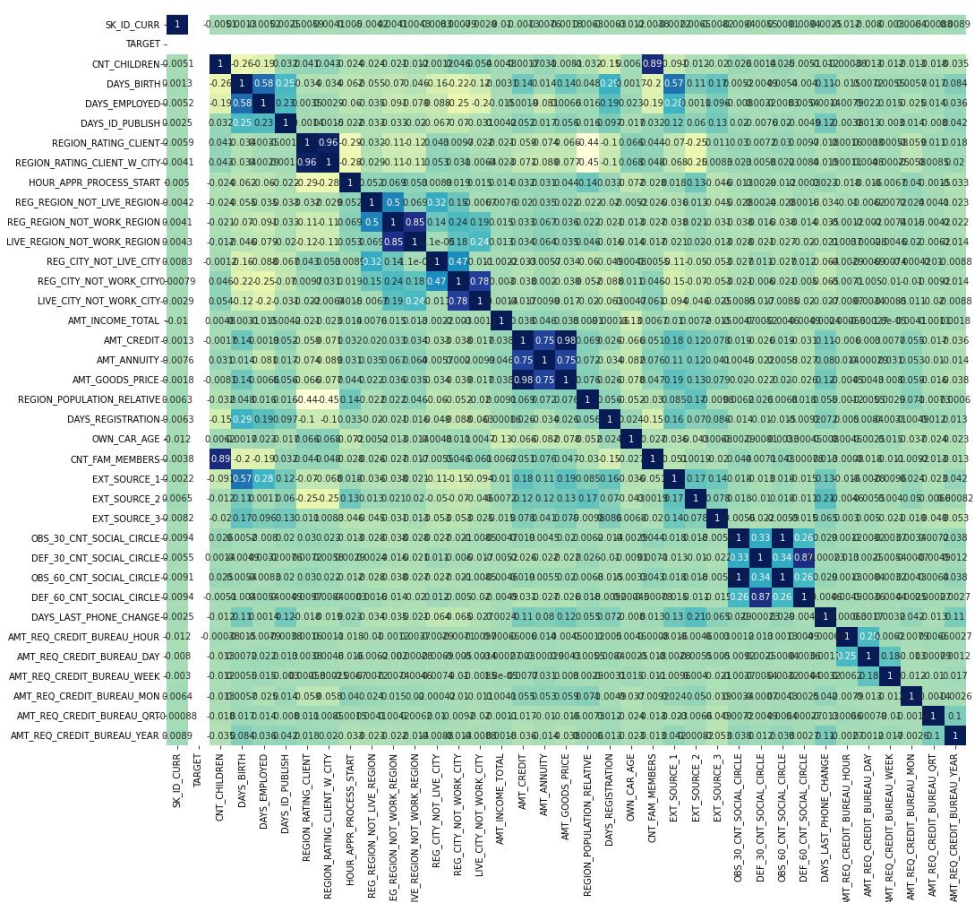




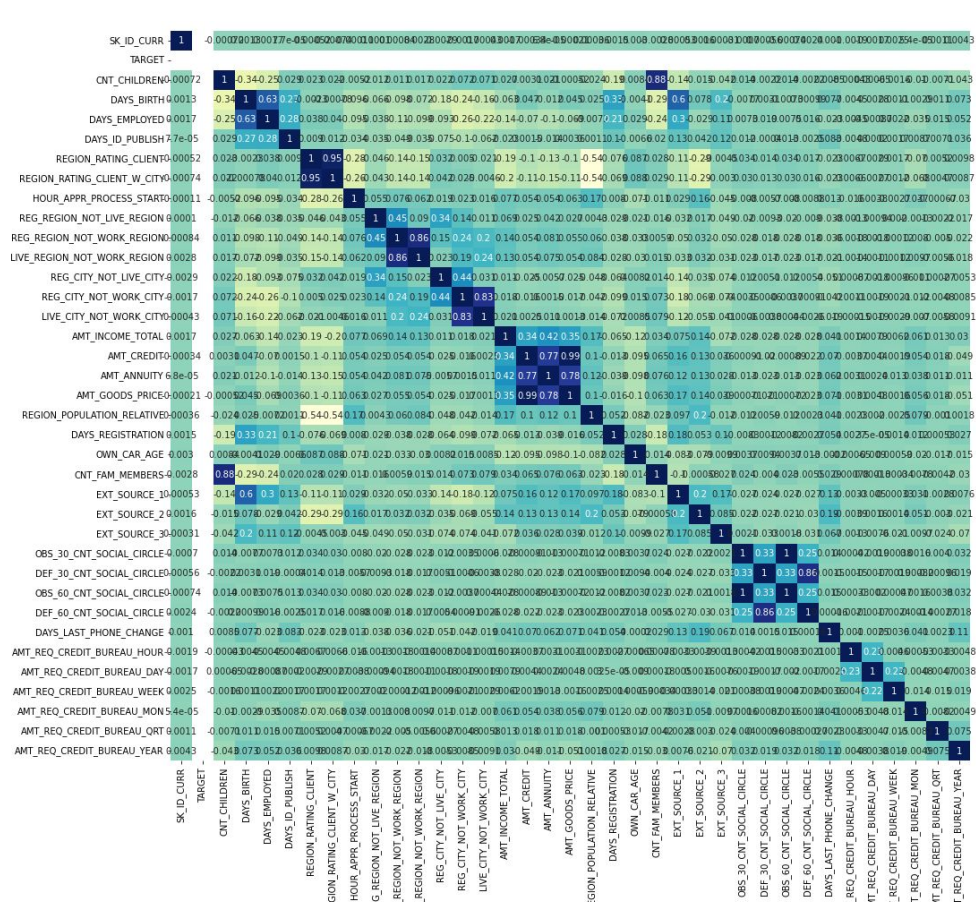




Defaulters

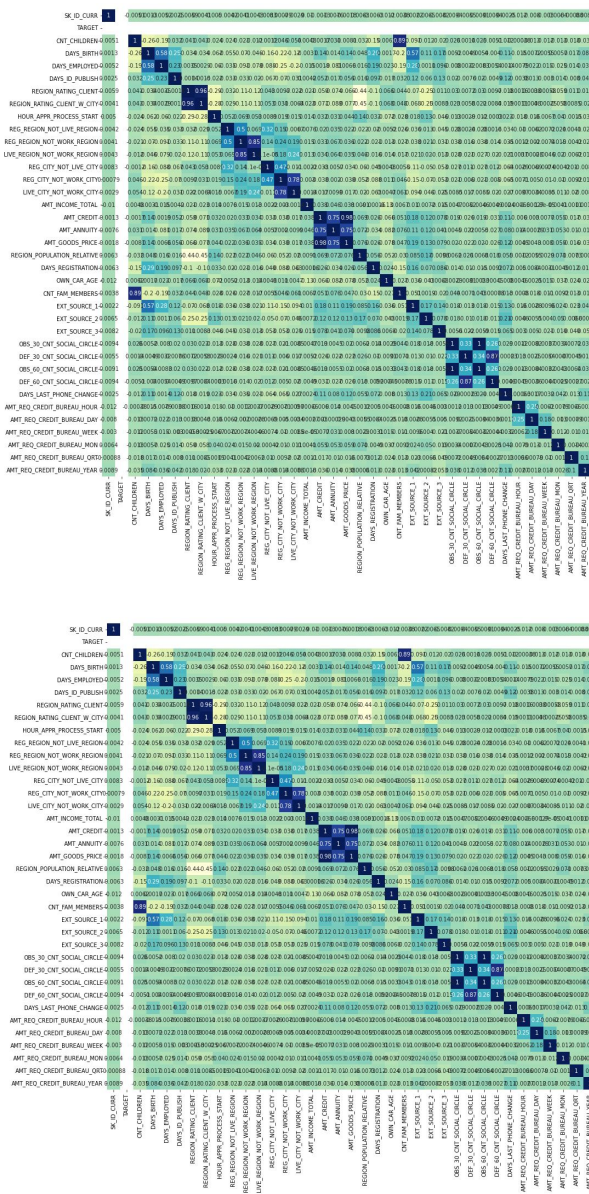


No Defaulters

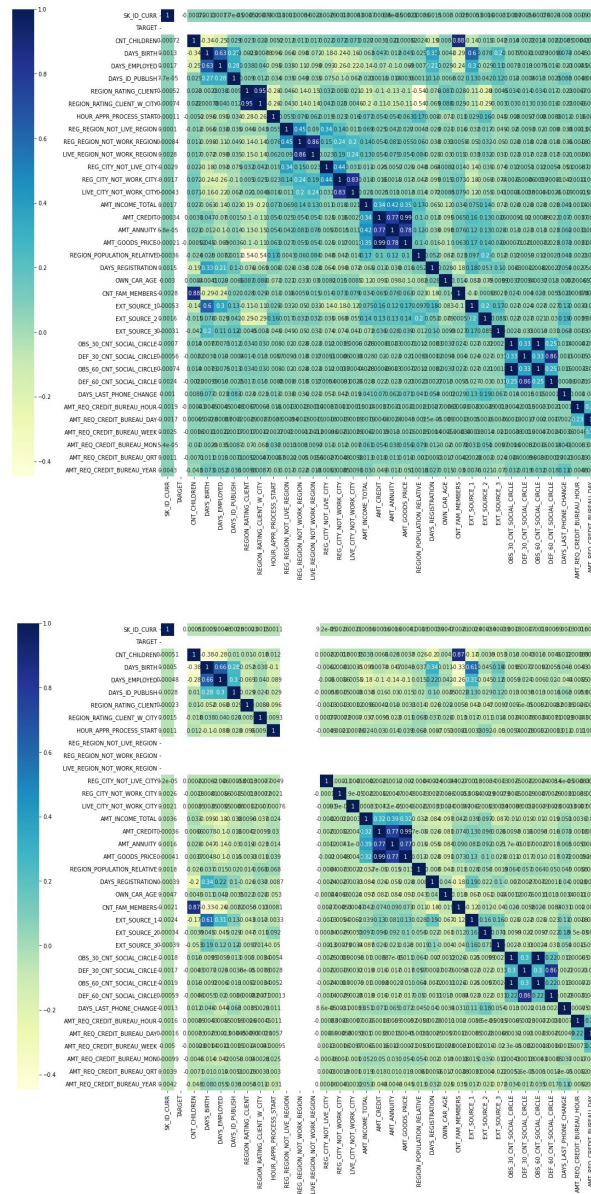




Defaulters



No Defaulters



%36.81

Personas con problemas de pago
eliminadas

CONCLUSIONES

- Es una base desbalanceada
- Mejorar limpieza de datos y outliers para no perder tantos datos
- Parece ser más sencillo identificar quien NO tendrá problemas con el pago