

In [759]. import pandas as pd import matplotlib.pyplot as plt import numpy as np from numpy import unique

Creamos un dataframe en el cual se excluyen de por si columnas que no tienen nada de informacion.

In [760].

```
df = pd.read_csv('small_train.csv')
df1 = df.dropna(how='all', axis=1)
df1
```

Out[760].

		Id	Var1	Var2	Var3	Var4	Var5	Var6	Var7	Var9	Var10	...	Var221	Var222	Var223	Var224	Var225	Var226	Var227	...	Var228	Var229	Target	
0	1	NaN	NaN	NaN	NaN	NaN	NaN	1526.0	7.0	NaN	NaN	...	osk	KVEsq	jySVZNOjy	NaN	NaN	x3v	RAYp		F2FyR07tdsN7i	NaN	-1	
1	2	NaN	NaN	NaN	NaN	NaN	NaN	525.0	0.0	NaN	NaN	...	osk	2Kb5FSF	LM8689qOp	NaN	NaN	fKcE	RAYp		F2FyR07tdsN7i	NaN	-1	
2	3	NaN	NaN	NaN	NaN	NaN	NaN	526.0	7.0	NaN	NaN	...	osk	AWZaUT	NKw4yOc	jySVZNOjy	NaN	KG3k	Qw4f	02nv6sff		iB5G6K1eUxUh6	am7c	-1
3	4	NaN	NaN	NaN	NaN	NaN	NaN	0.0	NaN	NaN	NaN	...	osk	CE7uK3u	LM8689qOp	NaN	NaN	FSa2	RAYp		F2FyR07tdsN7i	NaN	-1	
4	5	NaN	NaN	NaN	NaN	NaN	NaN	1029.0	7.0	NaN	NaN	...	osk	132cve	LM8689qOp	NaN	NaN	KG3k	FSa2	RAYp		F2FyR07tdsN7i	rm96	-1
...	
49995	49996	NaN	NaN	NaN	NaN	NaN	NaN	357.0	0.0	NaN	NaN	...	osk	EROH7Cq	LM8689qOp	NaN	NaN	7FJQ	RAYp		F2FyR07tdsN7i	NaN	-1	
49996	49997	NaN	NaN	NaN	NaN	NaN	NaN	1078.0	0.0	NaN	NaN	...	osk	GfSQowC	LM8689qOp	NaN	KG3k	FSa2	RAYp		55YFYv9	am7c	-1	
49997	49998	NaN	NaN	NaN	NaN	NaN	NaN	280.0	7.0	NaN	NaN	...	osk	dH6g2t	LM8689qOp	NaN	NaN	fKcE	RAYp		TCU50_yymmGIBZ0L	NaN	-1	
49998	49999	NaN	NaN	NaN	NaN	0.0	NaN	NaN	NaN	NaN	NaN	...	osk	2FZQqC	LM8689qOp	NaN	NaN	FSa2	RAYp		F2FyR07tdsN7i	NaN	-1	
49999	50000	NaN	NaN	NaN	NaN	NaN	NaN	1694.0	7.0	NaN	NaN	...	osk	tlvC9n	LM8689qOp	NaN	NaN	x3v	RAYp		F2FyR07tdsN7i	NaN	1	

50000 rows x 214 columns

Ahora analizamos la data en sus partes numerica y categorica. Cada una separada por resultados de 1 y -1

In [761].

```
numerical_data = df1.loc[:, "Var1":"Var190"]
numerical_target_data = numerical_data.join(df1["Target"])
df_numerical_t1 = numerical_target_data.loc[numerical_target_data.loc[:, "Target"]==1]
df_numerical_t1.describe()
```

Out[761].

	Var1	Var2	Var3	Var4	Var5	Var6	Var7	Var9	Var10	Var11	...	Var181	Var182	Var183	Var184	Var186	Var187
count	22.000000	30.0	30.000000	47.0	8.500000e+01	3476.000000	3461.000000	22.000000	8.500000e+01	30.000000	...	3498.000000	4.700000e+01	1.300000e+01	30.000000	22.000000	22.000000
mean	11.760909	0.0	4371.400000	327.461157	2.093675e+05	1466.443613	7.461138	142.636364	2.28987e+05	9.066667	...	6.060348	1.139278e+06	9.811920e+04	2.533333	2.454545	41.000000
std	76.117021	0.0	23853.665323	0.0	4.477857e+05	2351.148475	6.291430	483.236720	4.662392e+05	4.570772	...	2.565105	1.740722e+06	6.299477e+05	4.980983	4.404398	133.53045
min	0.000000	0.0	0.000000	0.0	0.000000e+00	0.000000	0.000000	0.000000	0.000000e+00	8.000000	...	0.000000	0.000000e+00	0.000000e+00	0.000000	0.000000	0.000000
25%	0.000000	0.0	0.000000	0.0	0.000000e+00	609.000000	7.000000	10.000000	0.000000e+00	8.000000	...	0.000000	0.000000e+00	0.000000e+00	0.000000	0.000000	5.000000
50%	0.000000	0.0	0.000000	0.0	0.000000e+00	952.000000	7.000000	10.000000	0.000000e+00	8.000000	...	0.000000	5.09520e+05	0.000000e+00	0.000000	0.000000	6.000000
75%	8.000000	0.0	0.000000	0.0	1.967300e+05	1660.750000	7.000000	73.500000	2.439360e+05	8.000000	...	0.000000	1.377453e+06	5.652700e+04	4.000000	6.000000	20.000000
max	360.000000	0.0	130668.000000	0.0	2.592000e+06	68439.000000	35.000000	2300.000000	2.332800e+06	32.000000	...	28.000000	7.200840e+06	1.299600e+06	20.000000	18.000000	634.000000

8 rows x 175 columns

In [762].

```
numerical_data = df1.loc[:, "Var1":"Var190"]
numerical_target_data = numerical_data.join(df1["Target"])
df_numerical_t1 = numerical_target_data.loc[numerical_target_data.loc[:, "Target"]==1]
df_numerical_t1.describe()
```

Out[762].

	Var1	Var2	Var3	Var4	Var5	Var6	Var7	Var9	Var10	Var11	...	Var181	Var182	Var183	Var184	Var186	Var187
count	680.000000	1211.000000	1210.000000	1532.000000	1.402000e+03	40995.000000	41000.000000	680.000000	1.402000e+03	1211.000000	...	41483.000000	1.532000e+03	1.211000e+03	1211.000000	22.000000	22.000000
mean	11.176471	0.004129	327.461157	0.129243	2.405774e+05	1.275481	6.441259e+05	6.754488	45.088235	4.025125e+05	8.614876	...	0.611718	1.425147e+06	7.726978e+04	8.607762	3.3
std	39.098886	0.143681	2154.679029	1.294719	6.542114e+05	2711.855948	6.325972	131.182470	9.478455e+05	2.816490	...	2.489772	2.294233e+06	2.006027e+05	4.7326890	8.8	8.8
min	0.000000	0.000000	0.000000	0.000000	0.000000e+00	0.000000	0.000000	0.000000	0.000000e+00	8.000000	...	0.000000	0.000000e+00	0.000000e+00	0.000000	0.000000	0.0
25%	0.000000	0.000000	0.000000	0.000000	0.000000e+00	511.000000	0.000000	4.000000	0.000000e+00	8.000000	...	0.000000	0.000000e+00	0.000000e+00	0.000000	0.000000	0.0
50%	0.000000	0.000000	0.000000	0.000000	0.000000e+00	854.000000	7.000000	20.000000	0.000000e+00	8.000000	...	0.000000	1.125990e+05	0.000000e+00	0.000000	0.000000	0.0
75%	16.000000	0.000000	0.000000	0.000000	1.131750e+05	1410.500000	7.000000	46.000000	2.647958e+05	8.000000	...	0.000000	1.876194e+06	5.862000e+04	8.000000	1.5	1.5
max	680.000000	5.000000	42588.000000	27.000000	6.048550e+06	131761.000000	140.000000	2098.000000	1.230000e+07	40.000000	...	49.000000	1.199478e+07	3.048400e+06	1200.000000	102.0	102.0

8 rows x 175 columns

In [763].

```
categorical_data = df1.loc[:, "Var191":"Target"]
df_categorical_t1 = categorical_data.loc[categorical_data.loc[:, "Target"]==1]
df_categorical_t1.describe(include='all')
```

Out[763].

	Var191	Var192	Var193	Var194	Var195	Var196	Var197	Var198	Var199	Var200	...	Var221	Var222	Var223	Var224	Var225	Var226	Var227	...	Var228	Var229	Target
count	39	3666	3662	959	3662	3662	3675	3682	3682	2173	...	3682	3682	3423	18	2160	3662	3682	...	3682	1930	3682.0
unique	1	224	31	3	11	2	155	1281	923	2060	...	7	1261	4	1	3	23	7	...	25	4	NaN
freq	1	354	33478	11630	44417	45893	4299	4064	873	73	...	34437	4064	33729	802	9871	7658	32724	...	30502	10659	NaN
mean	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
std	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
min	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
25%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
50%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
75%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
max	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN

11 rows x 39 columns

In [764].

```
categorical_data = df1.loc[:, "Var191":"Target"]
df_categorical_tn = categorical_data.loc[categorical_data.loc[:, "Target"]==1]
df_categorical_tn.describe(include='all')
```

Out[764].

	Var191	Var192	Var193	Var194	Var195	Var196	Var197	Var198	Var199	Var200	...	Var221	Var222	Var223	Var224	Var225	Var226	Var227	...	Var228	Var229	Target
count	1044	45965	46318	11825	46318	46318	46318	46318	46314	22419	...	46318	46318	41366	802	21696	46318	46318	...	46318	19638	46318.0
unique	1	360	51	3	23	4	223	4118	4865	14487	...	7	4118	4	1	3	23	7	...	30	4	NaN
freq	1	15400	1KXc	R012	SEuy	lmaJ	LK8T	0XwJ	rhX215s	rB3_sZi	...	oSk	catzSD2	LM8689qOp	4nX2	ELof	FSa2	RAYp	...	F2FyR07tdsN7i	am7c	NaN
mean	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
std	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
min	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
25%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
50%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
75%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
max	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN

11 rows x 39 columns

Ahora analizamos si no hay alguna variable que pertenezca a solo una de estas cuatro divisiones. Esto es, quitando columnas que no tengan absolutamente ningun valor en alguna de las cuatro dataframes que hemos creado. Como ya hemos quitado columnas que no tenian ningun valor en todo el dataframe, si una columna fuera a no tener un valor en algunos de estas divisiones, significaria que se presenta solo en una instancia de compra o no compra.

Sin embargo, encontramos que no hay variables que solo pertenezcan, por ejemplo, al data frame numerico de puros compradores. Y asi para los otros tres data frames. Sabemos entonces que no hay una variable que ocurra solo cuando el cliente compra o solo cuando el cliente no compra. Podemos entonces seguir a limpiar el data frame inicial completo.

In [765].

```
just_int_n1 = df_numerical_t1.dropna(how='all', axis=1)
just_int_ntn = df_numerical_tn.dropna(how='all', axis=1)
just_int_ct1 = df_categorical_t1.dropna(how='all', axis=1)
just_int_ctn = df_categorical_tn.dropna(how='all', axis=1)
```

Out[765].

	Var191	Var192	Var193	Var194	Var195	Var196	Var197	Var198	Var199	Var200	...	Var221	Var222	Var223	Var224	Var225	Var226	Var227	...	Var228	Var229	Target
count	39	3666	3662	959	3662	3662	3675	3682	3682	2173	...	3682	3682	3423	18	2160	3662	3682	...	3682	1930	3682.0
unique	1	224	31	3	11	2	155	1281	923	2060	...	7	1261	4	1	3	23	7	...	25	4	NaN
freq	1	354	33478	11630	44417	45893	4299	4064	873	73	...	34437	4064	33729	802	9871	7658	32724	...	30502	10659	NaN
mean	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
std	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
min	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
25%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN
50%	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN</						