



Instituto Tecnológico  
de Buenos Aires

## PRESENTACION FINAL

Sofía Feilbogen - L. 61889

—

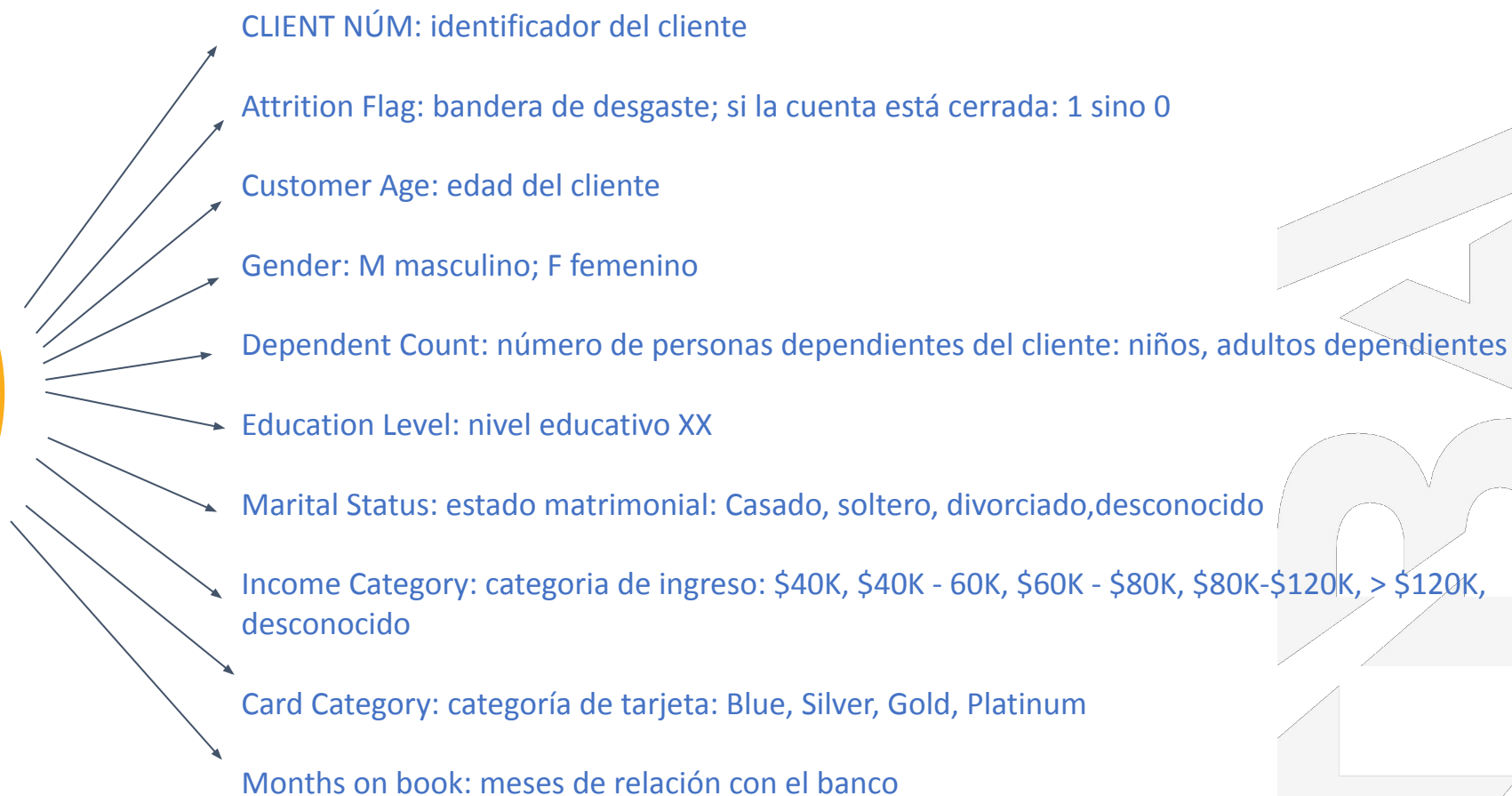
# ÍNDICE

1. Introducción
  - Datos del cliente
  - Objetivos del trabajo
  - Características del trabajo
2. Estadísticas descriptivas principales
3. Data discovery
4. Partición
5. Pipeline
6. Modelos
  - Ajustes de hiperparametros
7. Conclusión

# Introducción



## Datos del cliente

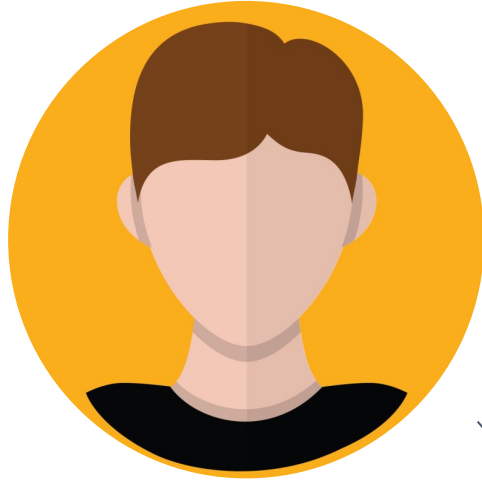


10127 registros. 23 variables.

Naive\_Bayes\_Classifier\_Attrition\_Flag\_Card\_Category\_Contacts\_Count\_12\_mon\_Dependent\_count\_Education\_Level\_Months\_Inactive\_12\_mon\_1

Naive\_Bayes\_Classifier\_Attrition\_Flag\_Card\_Category\_Contacts\_Count\_12\_mon\_Dependent\_count\_Education\_Level\_Months\_Inactive\_12\_mon\_2

## Datos del cliente



# Objetivos del trabajo



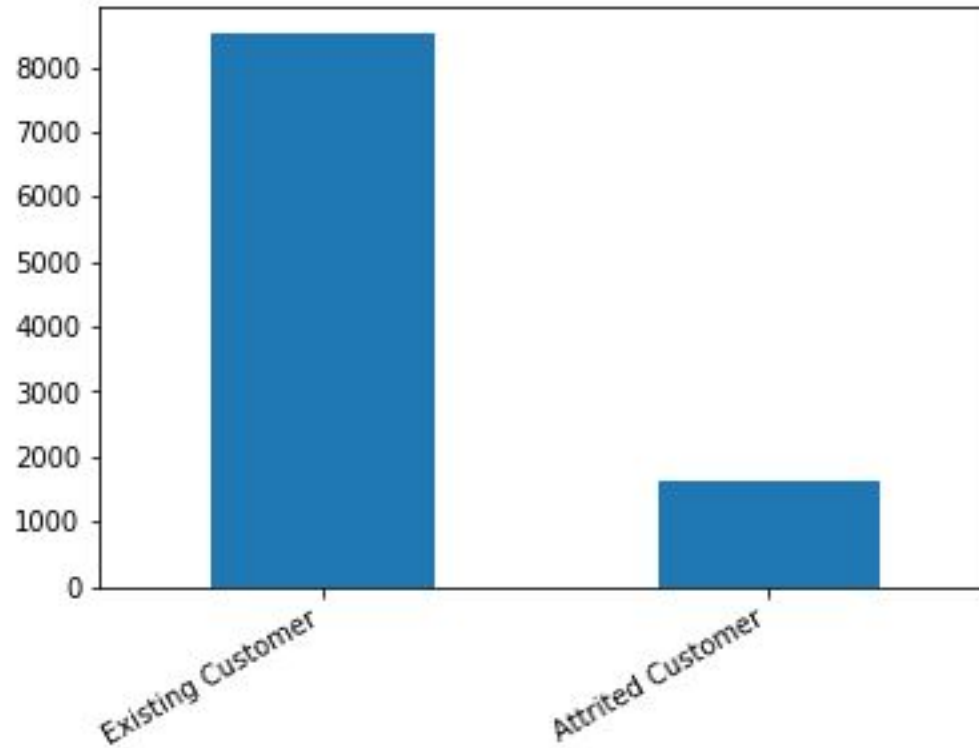
Entender qué variables afectan a un cliente para que sea perdido.

Encontrar el perfil del cliente perdido.

De esta forma entendemos sobre qué tenemos que poner foco para que un cliente existente no se pierda y así retenerlo.

De esta forma también podemos categorizar al cliente, dejar un registro de que es posible que se pierda o no en un futuro y estar más atento a los que se pierdan posiblemente.

## Características del trabajo



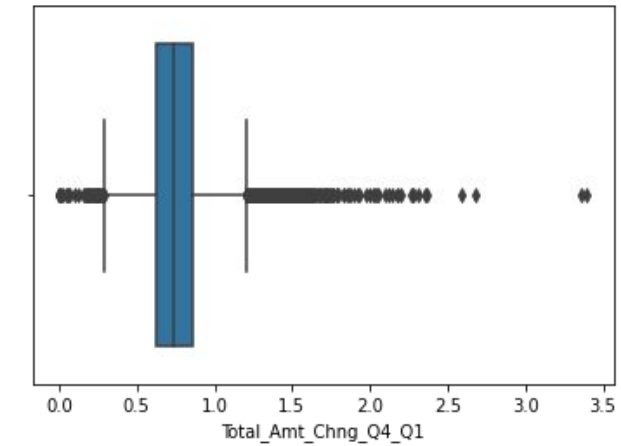
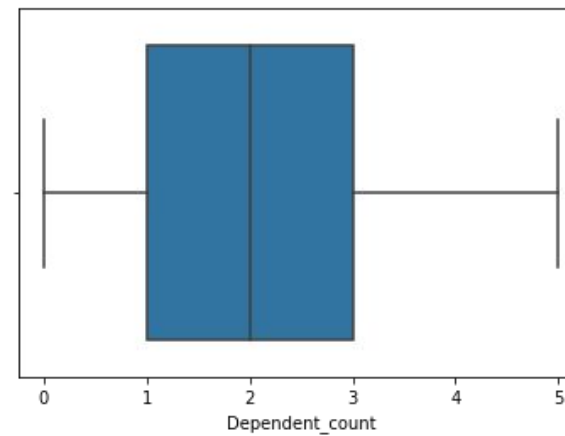
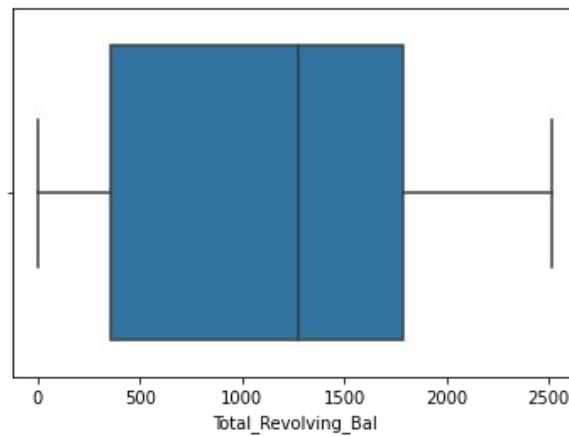
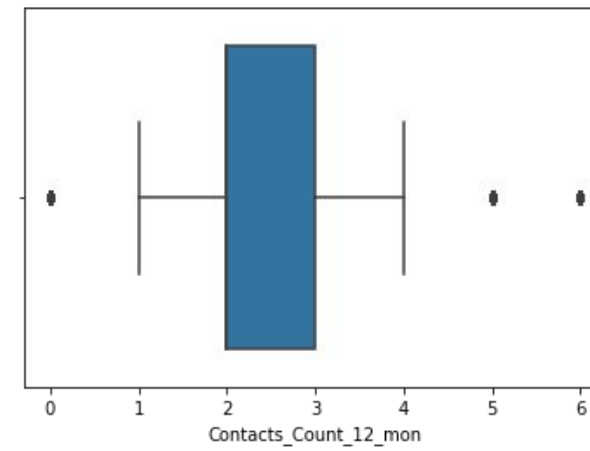
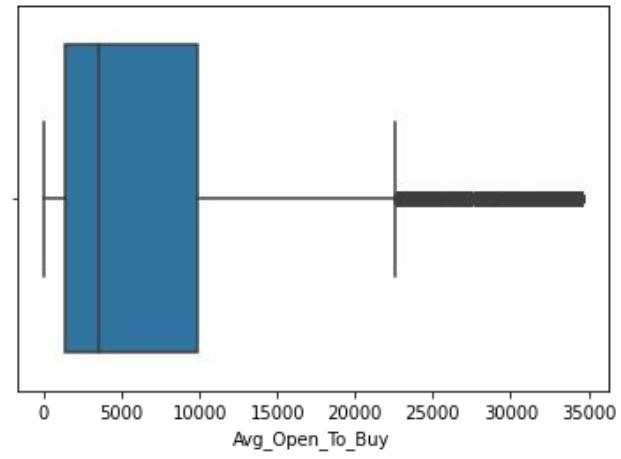
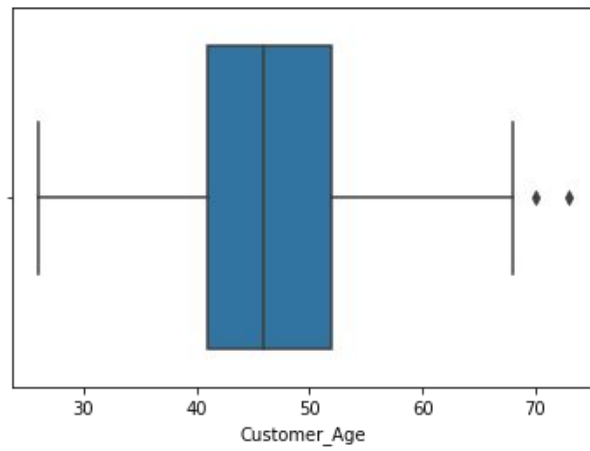
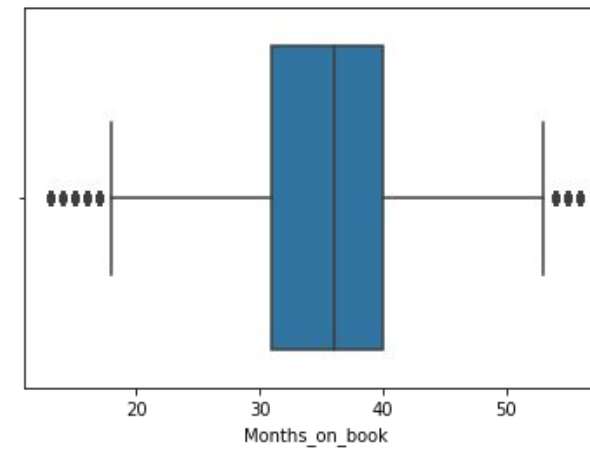
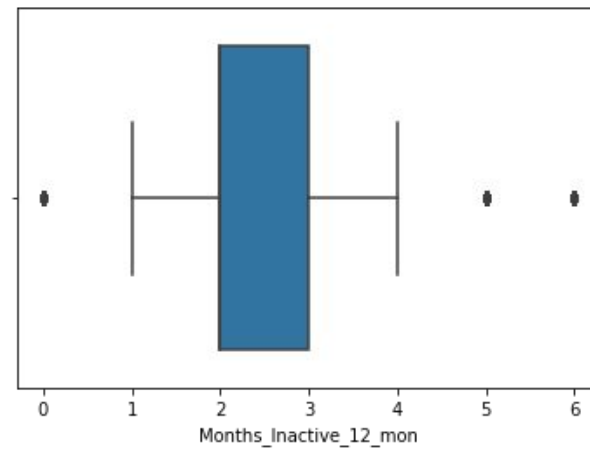
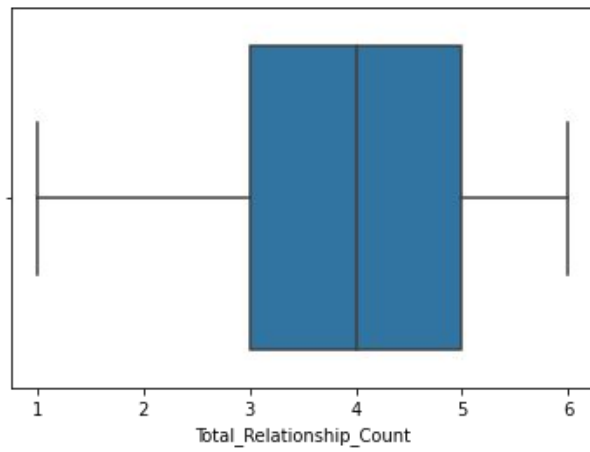
Datos asimétricos y desbalanceados

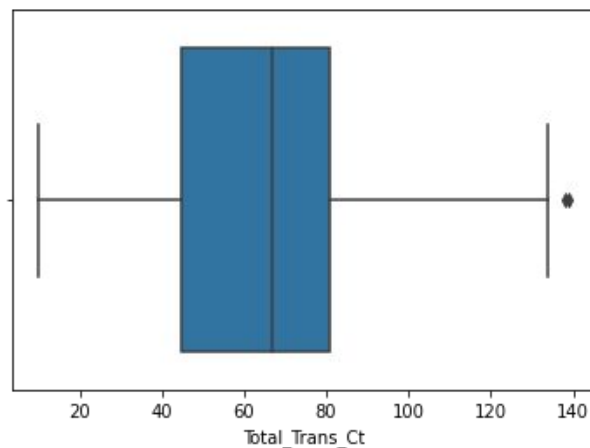
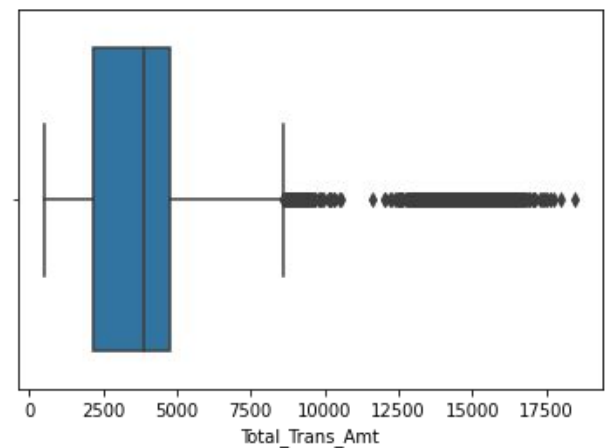
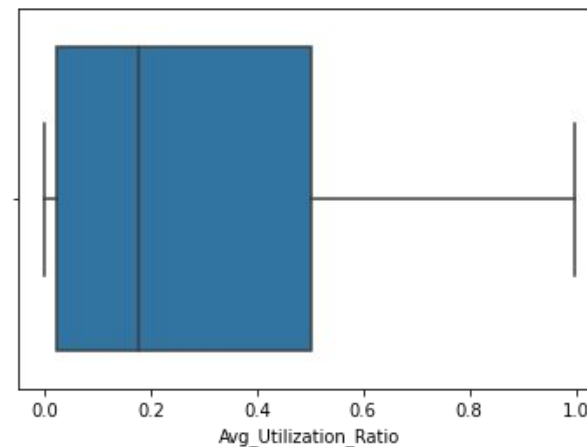
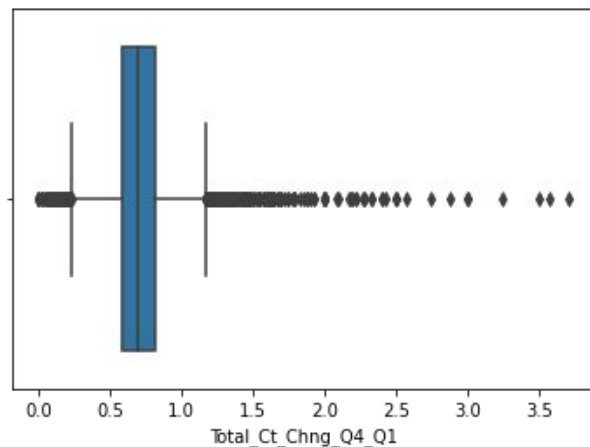
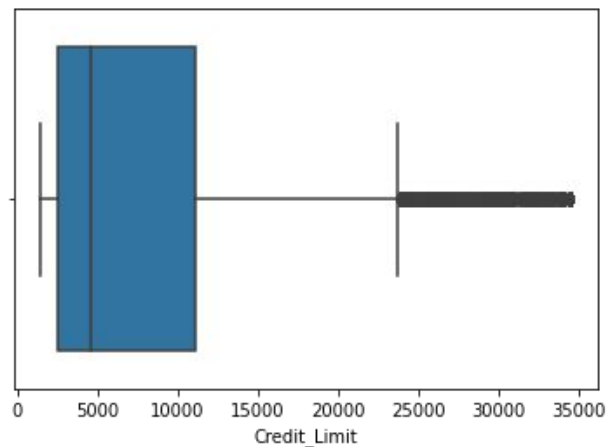


Existing Customer 8500  
Attrited Customer 1627  
20% de 10127 = 2025,4 (moderately imbalanced)

# Estadísticas descriptivas principales







- No se encuentran missings
- Se encuentran outliers en 10 variables

Education Level, Marital Status, Income Category → "Unknown" → NaN

Edad más chica: 26

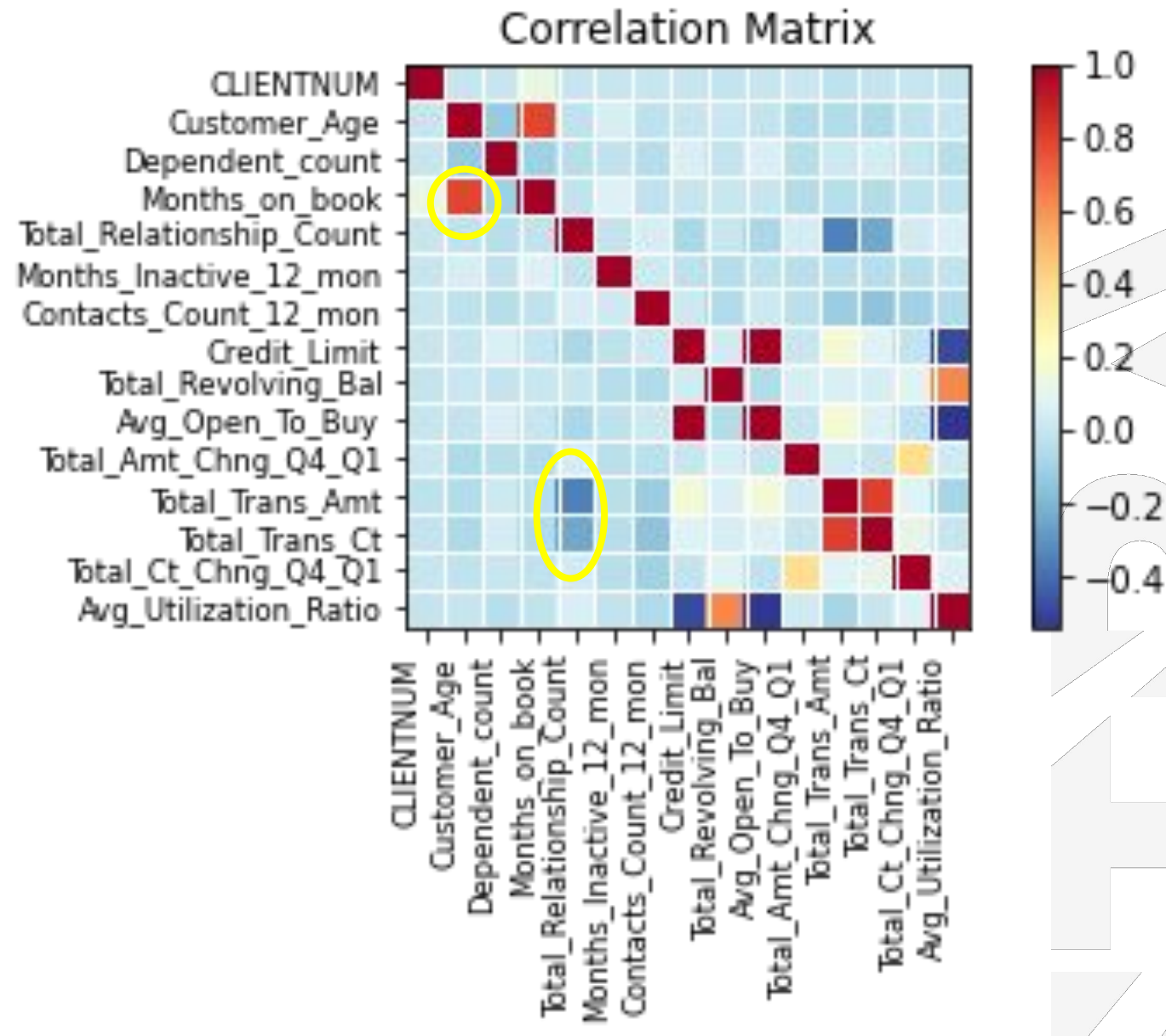
	CLIENTNUM	Customer_Age	Dependent_count	Months_on_book	Total_Relationship_Count	Months_Inactive_12_mon	Contacts_Count_12_mon
count	1.012700e+04	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000
mean	7.391776e+08	46.325960	2.346203	35.928409	3.812580	2.341167	2.455317
std	3.690378e+07	8.016814	1.298908	7.986416	1.554408	1.010622	1.106225
min	7.080821e+08	26.000000	0.000000	13.000000	1.000000	0.000000	0.000000
25%	7.130368e+08	41.000000	1.000000	31.000000	3.000000	2.000000	2.000000
50%	7.179264e+08	46.000000	2.000000	36.000000	4.000000	2.000000	2.000000
75%	7.731435e+08	52.000000	3.000000	40.000000	5.000000	3.000000	3.000000
max	8.283431e+08	73.000000	5.000000	56.000000	6.000000	6.000000	6.000000

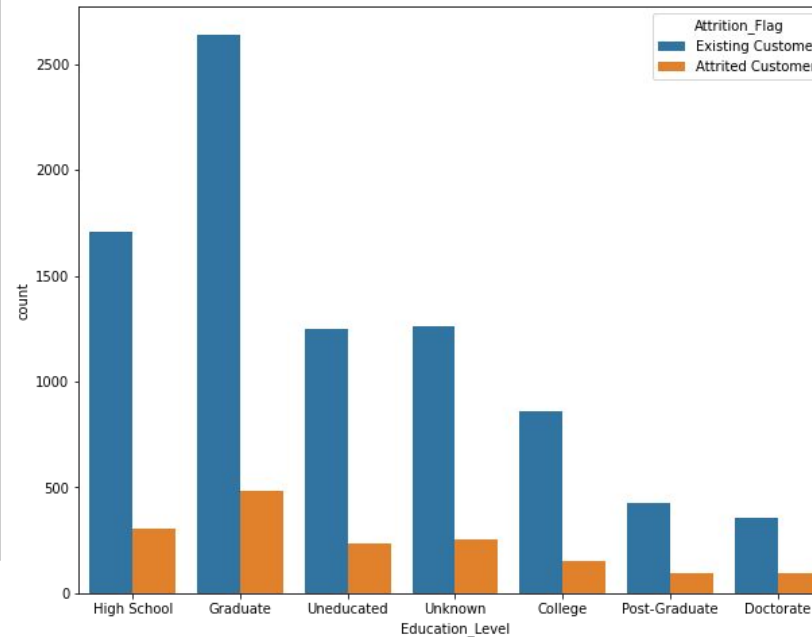
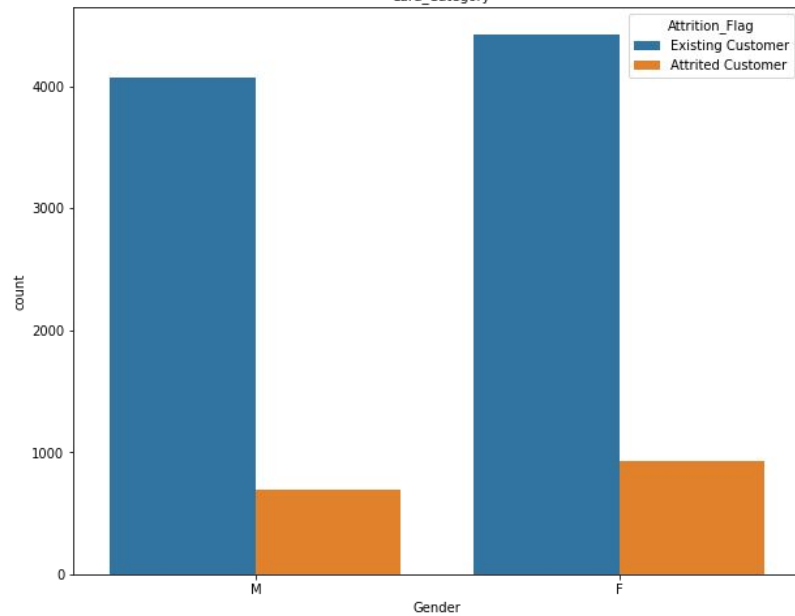
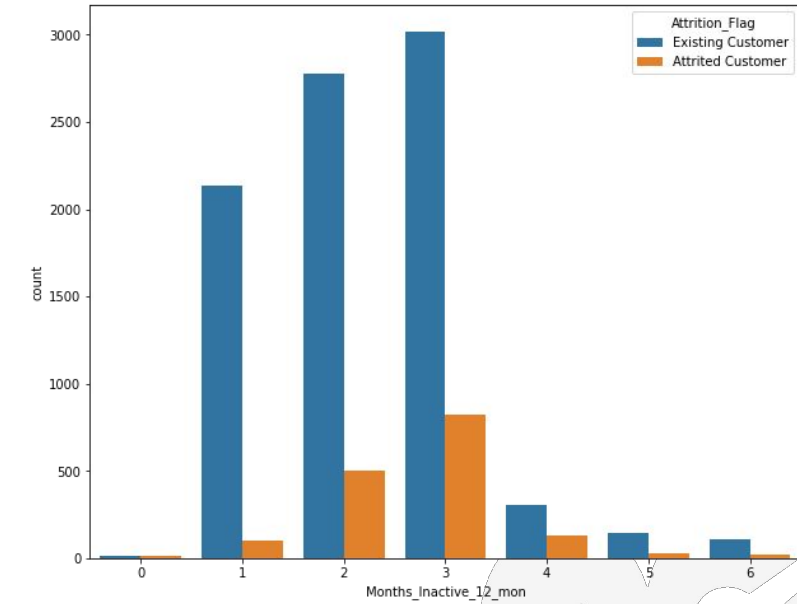
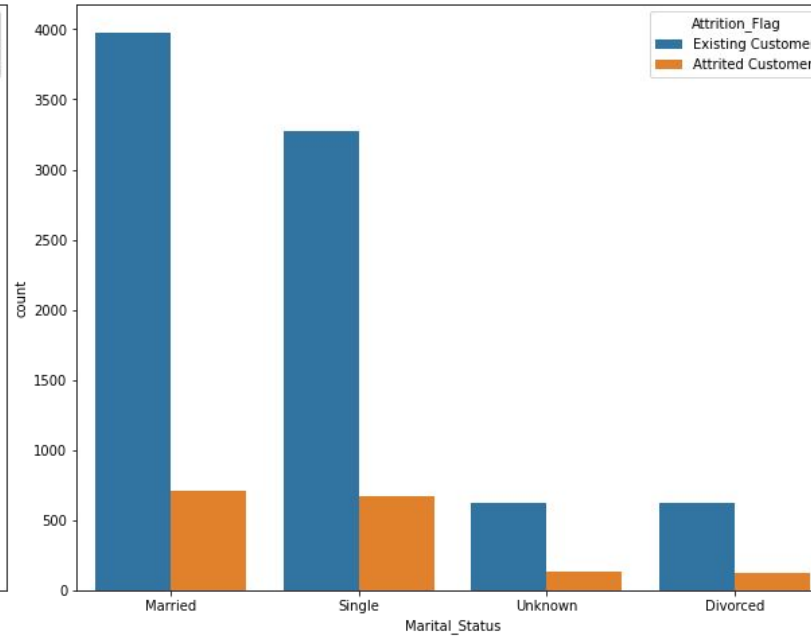
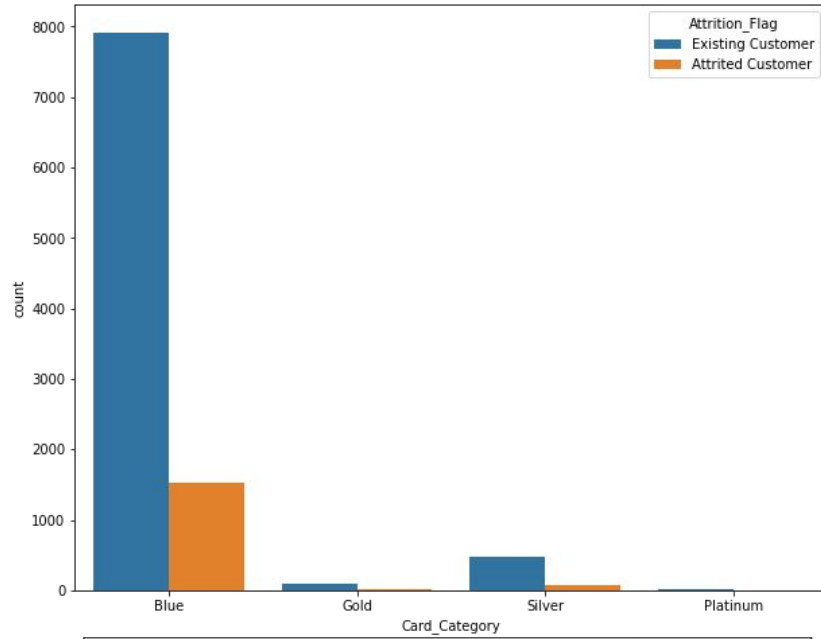
Credit_Limit	Total_Revolving_Bal	Avg_Open_To_Buy	Total_Amt_Chng_Q4_Q1	Total_Trans_Amt	Total_Trans_Ct	Total_Ct_Chng_Q4_Q1	Avg_Utilization_Ratio
10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000	10127.000000
8631.953698	1162.814061	7469.139637	0.759941	4404.086304	64.858695	0.712222	0.274894
9088.776650	814.987335	9090.685324	0.219207	3397.129254	23.472570	0.238086	0.275691
1438.300000	0.000000	3.000000	0.000000	510.000000	10.000000	0.000000	0.000000
2555.000000	359.000000	1324.500000	0.631000	2155.500000	45.000000	0.582000	0.023000
4549.000000	1276.000000	3474.000000	0.736000	3899.000000	67.000000	0.702000	0.176000
11067.500000	1784.000000	9859.000000	0.859000	4741.000000	81.000000	0.818000	0.503000
34516.000000	2517.000000	34516.000000	3.397000	18484.000000	139.000000	3.714000	0.999000

# Data discovery

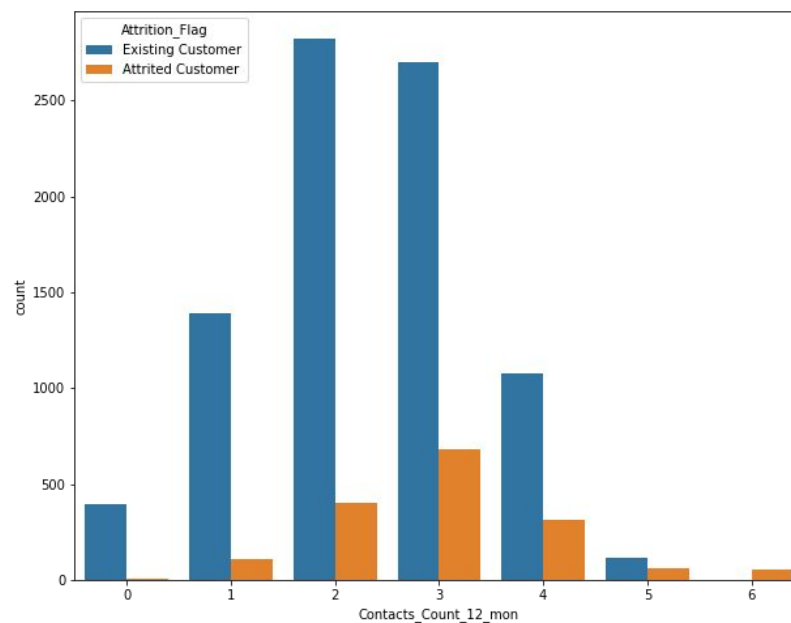
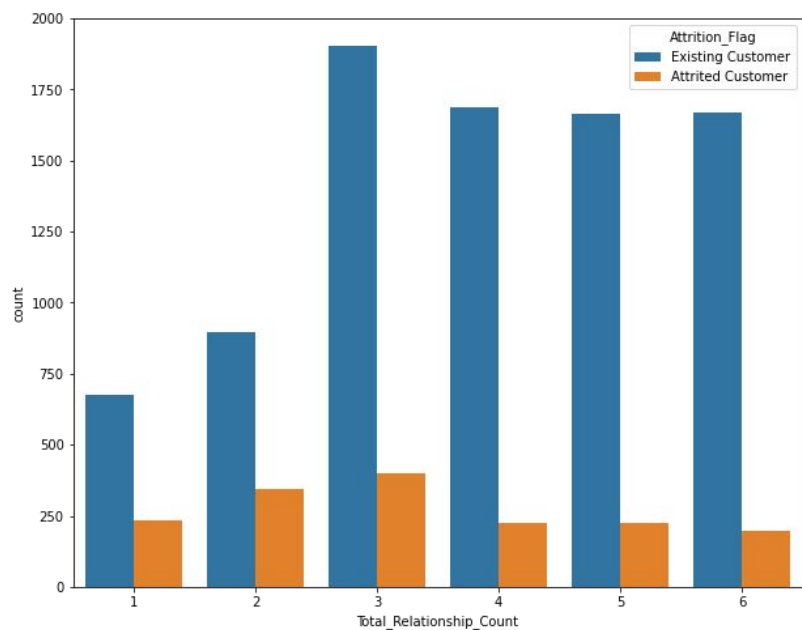
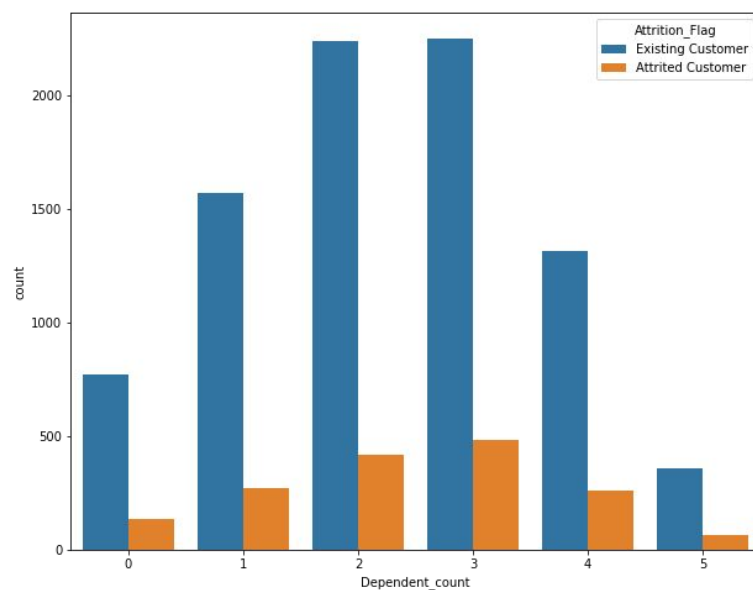
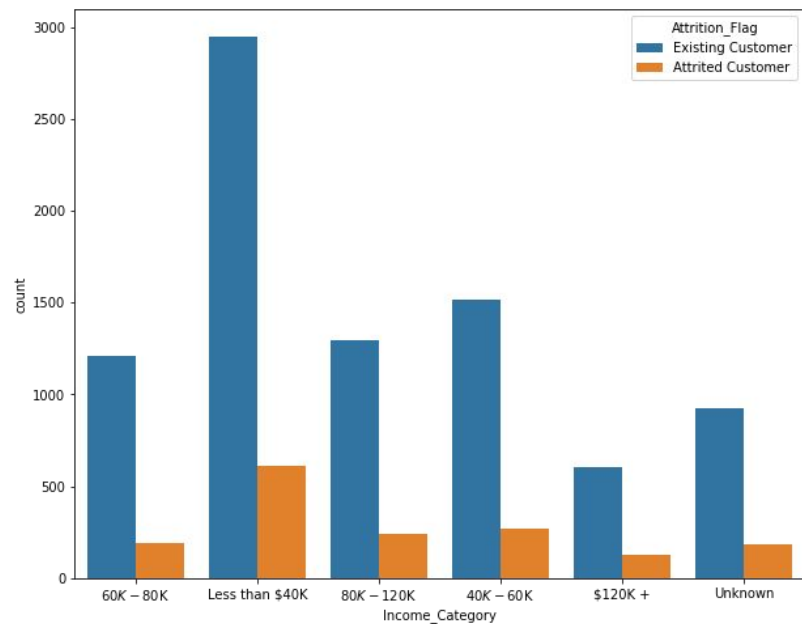


## Matriz de correlación de Pearson



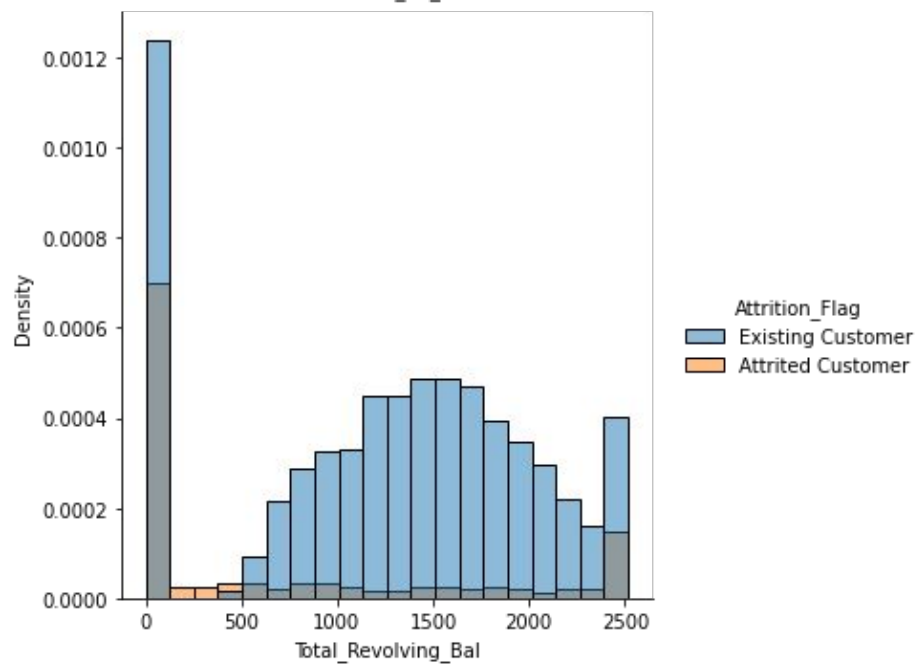
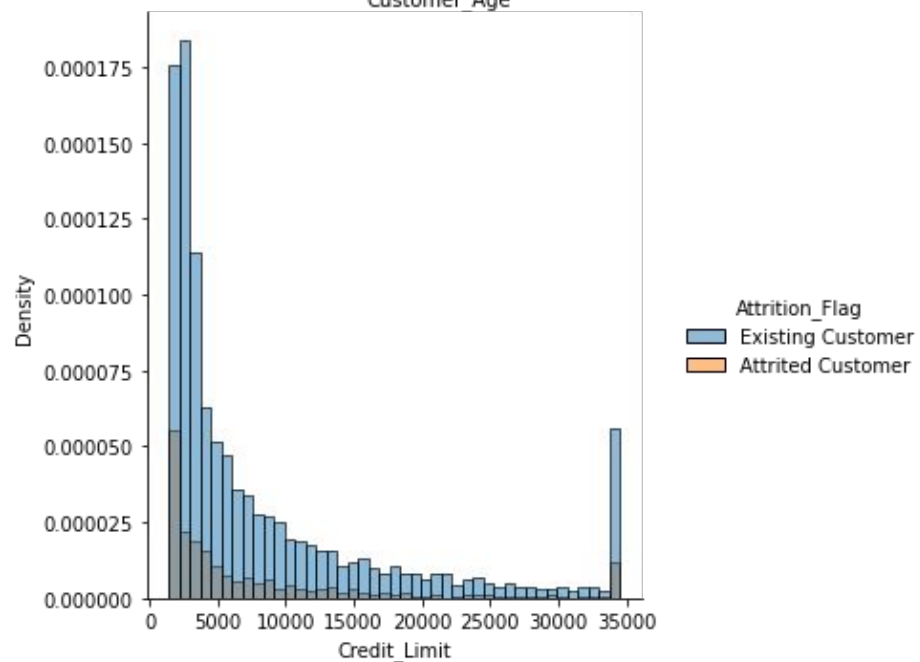
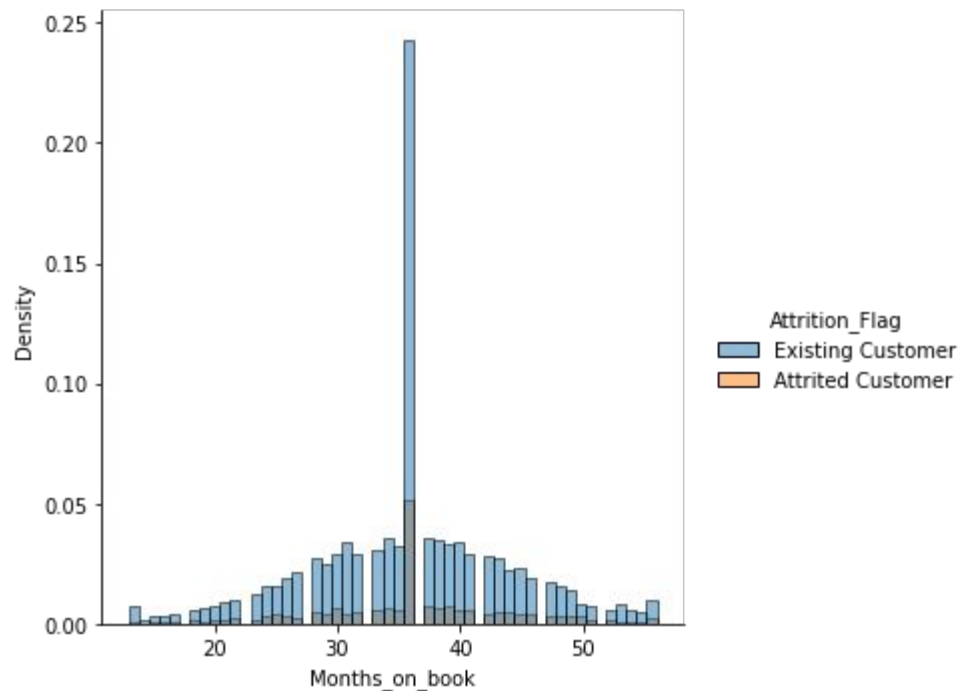
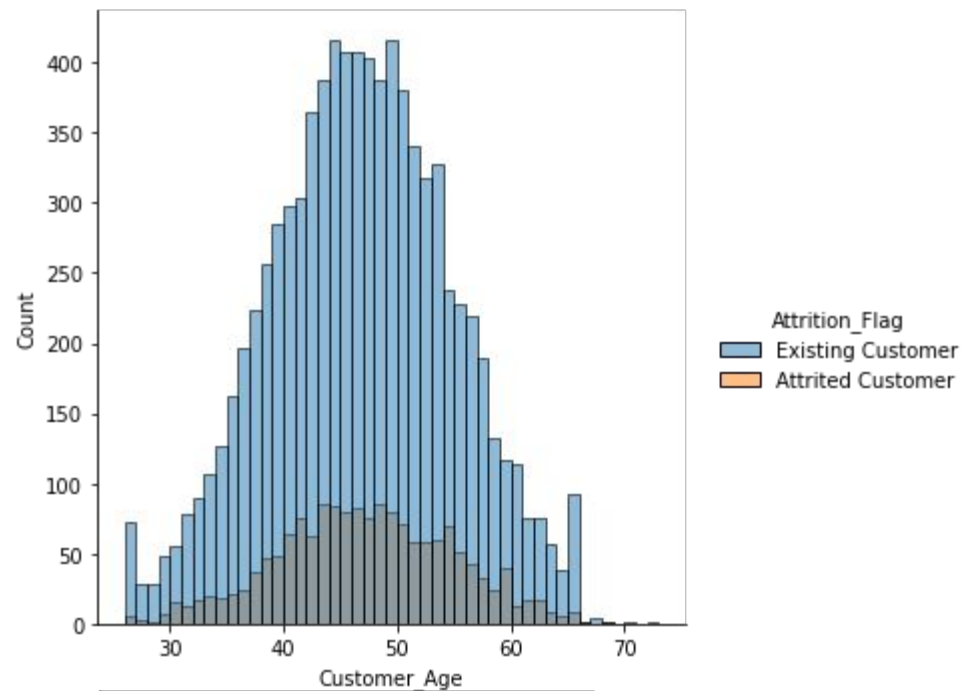


Card category: blue  
Marital status: married / single  
Months inactive 12 mon: 2 / 3  
Gender: M / F  
Education level: High school / graduate

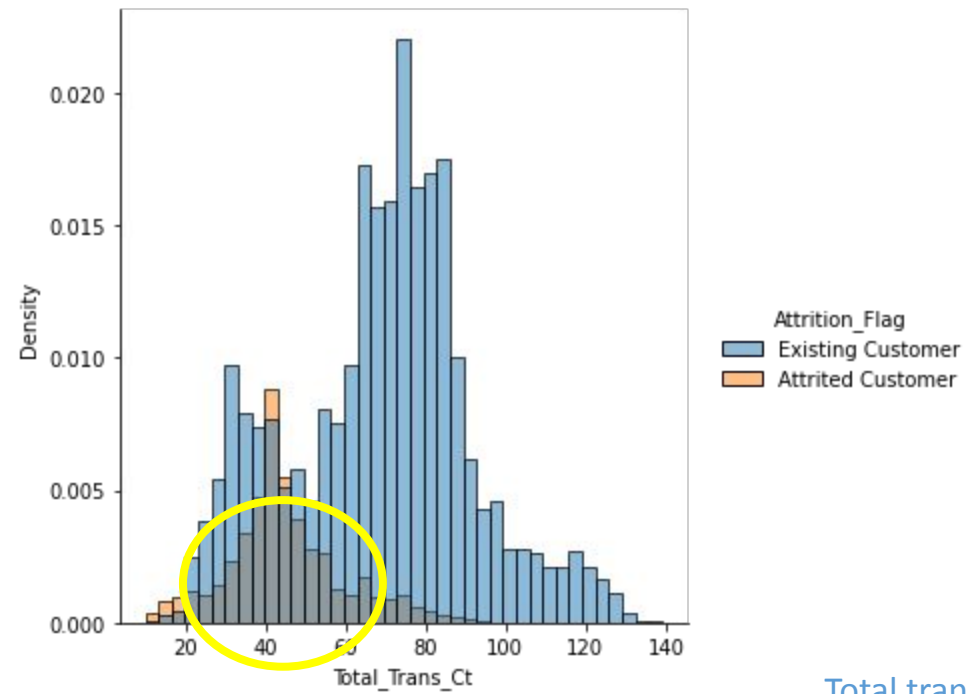
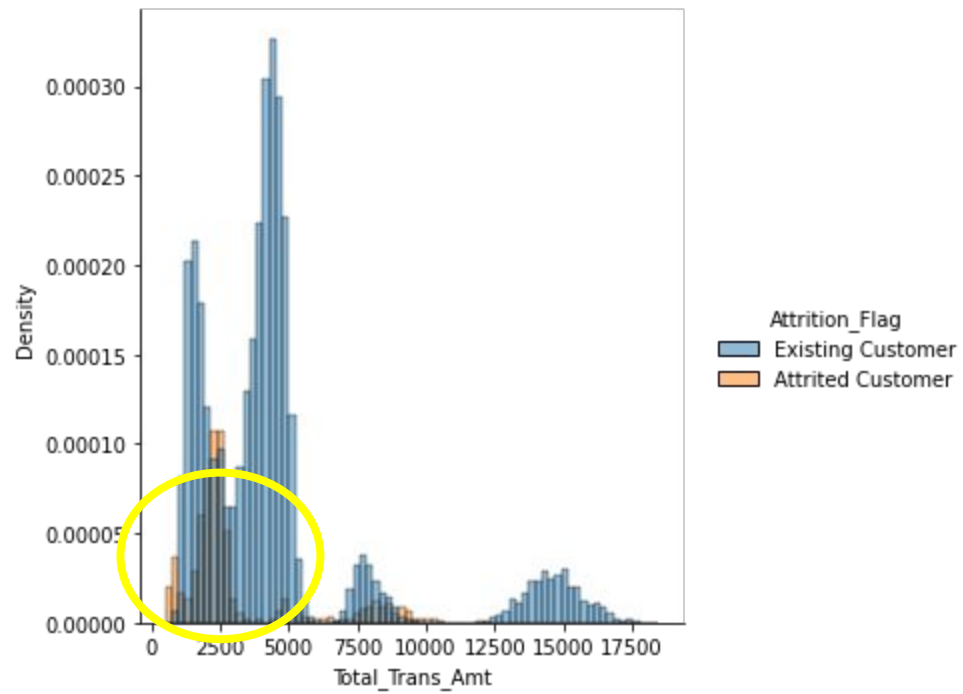


Less than \$40K  
 Dependent count: 2 / 3  
 Total relationship count: 2 / 3  
 Contact count: 2 / 3 / 4

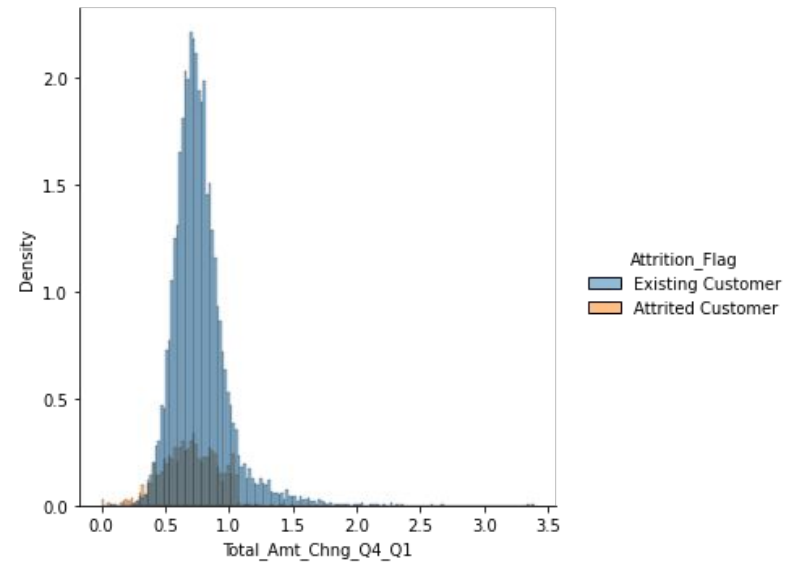
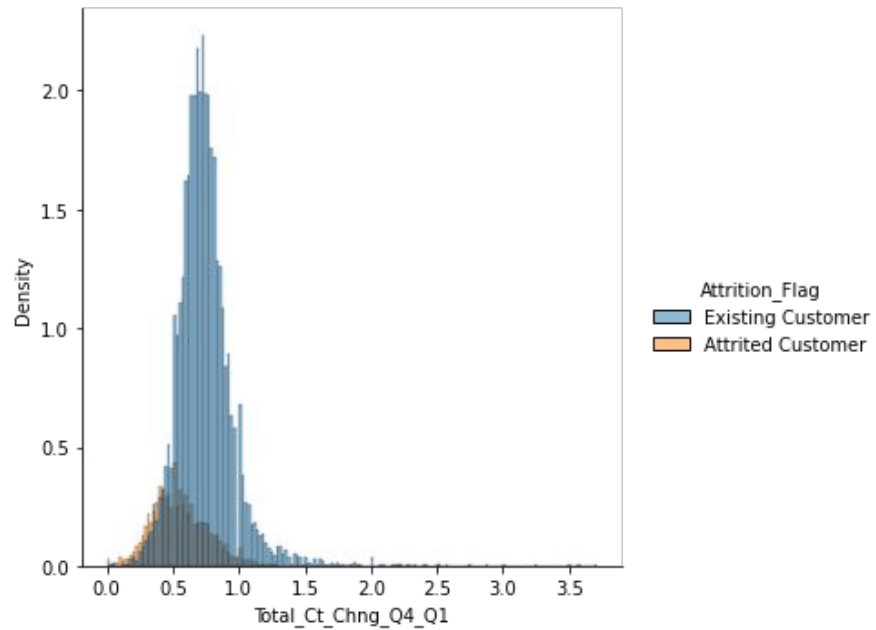






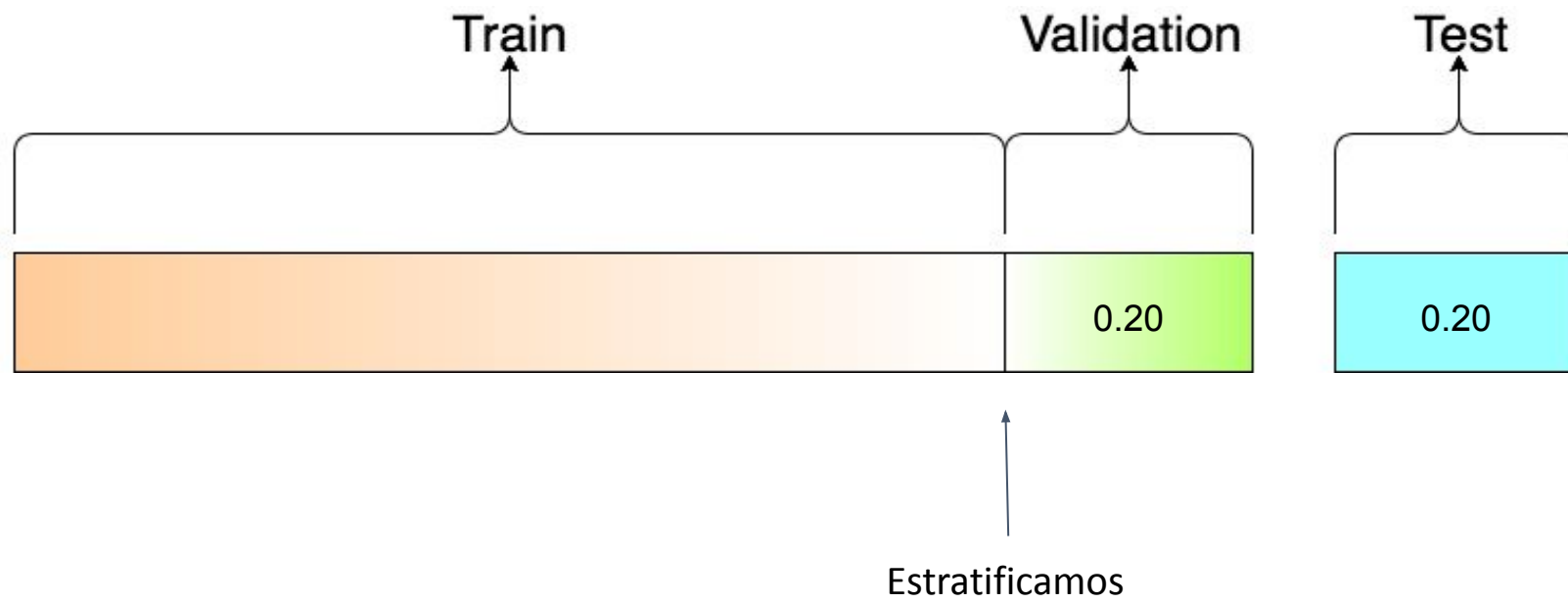


Total trans amt: 0-3000  
Total trans ct: 0-1000



# Partición



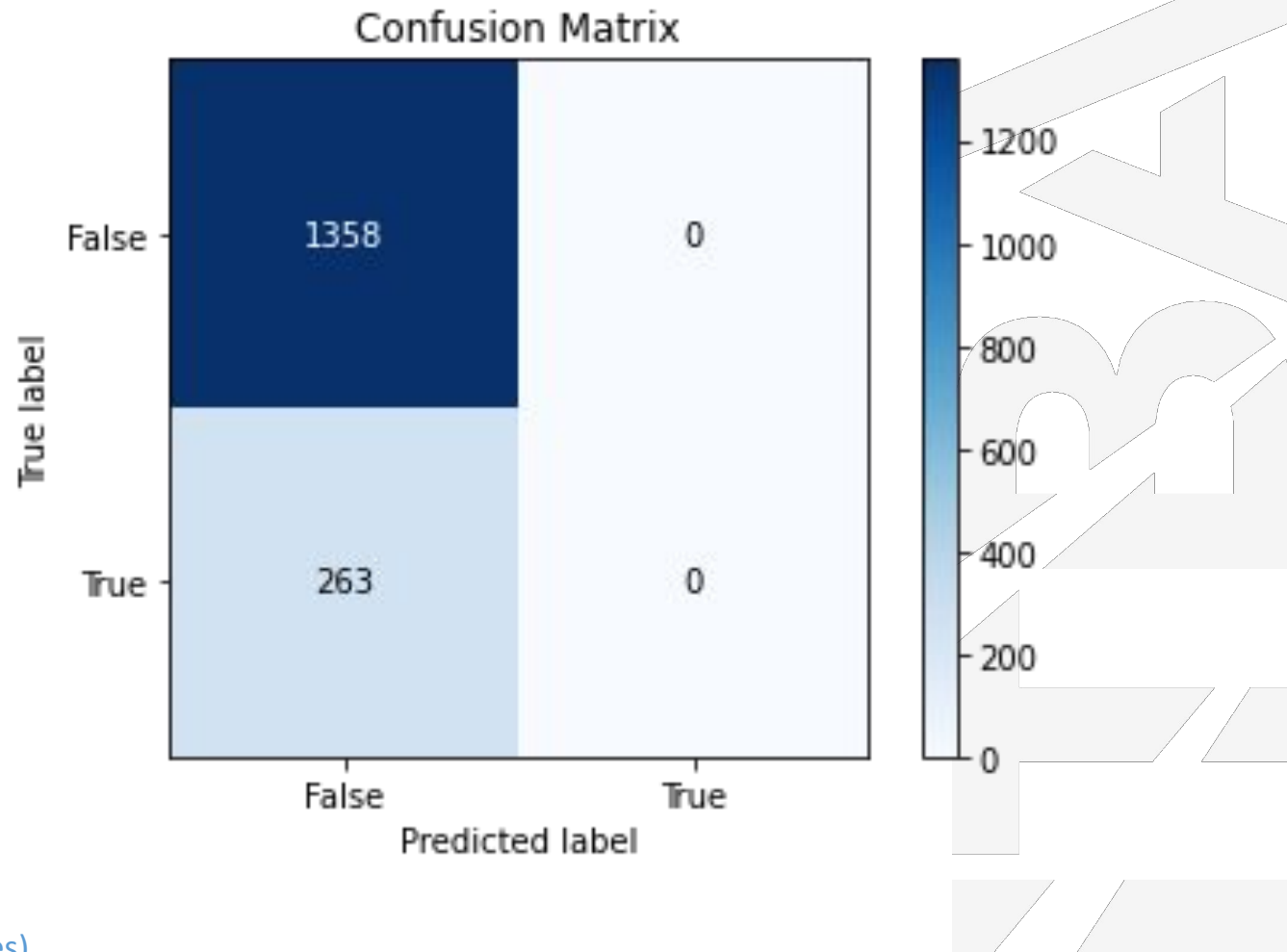


Transformamos a **False (0)** la clase mayoritaria (Existing Customer) y a **True (1)** la minoritaria.

y	Existing Customer	0.837921
	Attrited Customer	0.162079
y_train	Existing Customer	0.837963
	Attrited Customer	0.162037
y_val	Existing Customer	0.837754
	Attrited Customer	0.162246

Accuracy de predecir la clase mayoritaria:  
0.83775447254781

raw frequencies:  
[5430 1050]  
class weights:  
[0.59668508 3.08571429]  
[3240. 3240.]



$W_i = \text{cantidad de observaciones} / (\text{cantidad de clases} \times \text{frecuencia de clases})$

# Pipeline





#### Winsorizer

- Cola derecha: mean + 3\* std
- Cola izquierda: mean - 3\* std

	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book	Total_Relationship_Count	Months_Inactive_12_mon	Contacts_Count_12_mon	Credit_Limit	Total_Revolving_Bal	Avg_Open_To_Buy
8951	52.0000	0	2	2	0	2	0	37	1	3.0000	1.0000	34516.0000	1369	33147.0000
7232	49.0000	0	2	2	0	2	0	38	3	4.0000	3.0000	34516.0000	0	34516.0000
8861	55.0000	0	3	3	2	4	0	44	2	2.0000	3.0000	6455.0000	1837	4618.0000
2112	53.0000	0	4	2	0	2	0	36	3	2.0000	2.0000	3924.0000	2517	1407.0000
4361	46.0000	1	3	1	0	3	0	38	4	2.0000	2.0000	1781.0000	1315	466.0000

Total_Amt_Chng_Q4_Q1	Total_Trans_Amt	Total_Trans_Ct	Total_Ct_Chng_Q4_Q1	Avg_Utilization_Ratio	Customer_Age_right	Dependent_count_right	Months_on_book_right	Total_Relationship_Count_right	Months_Inactive_12_mon_right
0.7640	7745.0000	82	0.6400	0.0400	0.0000	0.0000	0.0000	0.0000	0.0000
0.4820	1592.0000	35	0.3460	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0.7170	8001.0000	99	0.7070	0.2850	0.0000	0.0000	0.0000	0.0000	0.0000
0.5110	1417.0000	30	0.6670	0.6410	0.0000	0.0000	0.0000	0.0000	0.0000
0.7040	4775.0000	82	0.7290	0.7380	0.0000	0.0000	0.0000	0.0000	0.0000

Contacts_Count_12_mon_right	Credit_Limit_right	Total_Revolving_Bal_right	Avg_Open_To_Buy_right	Total_Amt_Chng_Q4_Q1_right	Total_Trans_Amt_right	Total_Trans_Ct_right	Total_Ct_Chng_Q4_Q1_right	Avg_Utilization_Ratio_right
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000

# Modelos





# Métricas de clasificación

Al ser datos asimétricos(nos interesa más la clase positiva), utilizamos el F score para medir la capacidad de clasificación del modelo.

La precisión es cuántos de los que yo identifique como positivos eran verdaderamente positivos

El recall es cuántos de los positivos yo capturo, cuantos los identifico.

El F-score es la media geométrica de precisión y recordación. Se encuentra entre la precisión y la recuperación.

Valor beta default 1 dándole mismo peso a recall como a la precisión

$$\text{precision} = \frac{tp}{tp + fp}$$

$$\text{recall} = \frac{tp}{tp + fn}$$

$$F_1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

# Métricas de performance

AUPROC: para medir la bondad global del modelo, para medir qué tan bueno es el score utilizamos la curva de AUPROC.

AUROC evalúa FPR (bondad para captar negativos) y TNR (bondad para captar positivos) con la misma ponderación. La curva AUPROC a diferencia de la AUROC es sensible al desbalance y predice con precisión la clase minoritaria. Se utiliza para datos asimétricos y desbalanceados como lo es en nuestro caso.

Es más importante medir la clase minoritaria (1) ya que buscamos encontrar las razones de por que se pierden los clientes. Entonces utilizamos Area Under the Precision-Recall Curve.

AUROC utiliza el  $FPR = FP / (FP + TN)$ . En estos casos TN es grande entonces baja su valor y lo vuelve poco informativo.

AUPROC trabaja con el TPR y la precisión relacionada con los datos positivos (minorías ambos) cambian sus ejes.

# LGBMClassifier

`num_leaves`: número máximo de hojas

`min_child_samples`: número mínimo de data necesaria en una hoja

`max_depth`: número máximo de profundidad del árbol

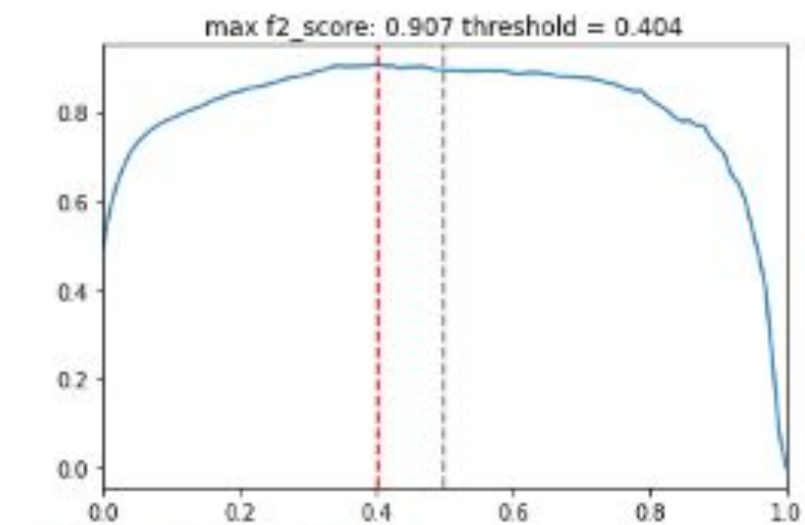
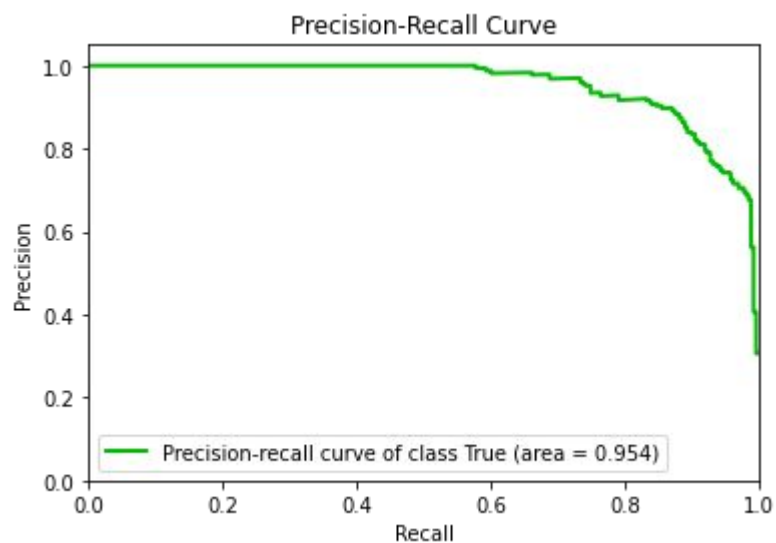
`objective`: "binary": clasificación binaria

`class_weight`: balanced

```
param_gridLGBM = {'LGBM__num_leaves': [5,10,20],  
                  'LGBM__min_child_samples': [5,10,15],  
                  'LGBM__max_depth': [5,10,20],  
                  }
```

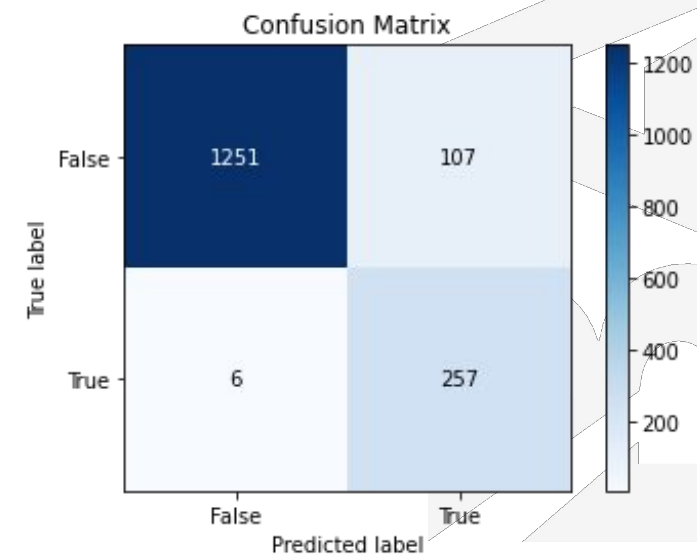
```
gridsearch1.best_params_
```

```
{'LGBM__max_depth': 5, 'LGBM__min_child_samples': 5, 'LGBM__num_leaves': 5}
```



NNN threshold = 0.404  
~~accuracy: 0.938~~  
 f2: 0.907  
 AUPRC: 0.954  
 NNN threshold = 0.500  
 accuracy: 0.941  
 f2: 0.894  
 AUPRC: 0.954

precision = 0.7060  
 recall = 0.9772  
 f1 = 0.8198



# Decision Tree Classifier

**min\_samples\_leaf**: número mínimo de registros requeridos para que se cree una hoja

**max\_depth**: número máximo de profundidad del árbol

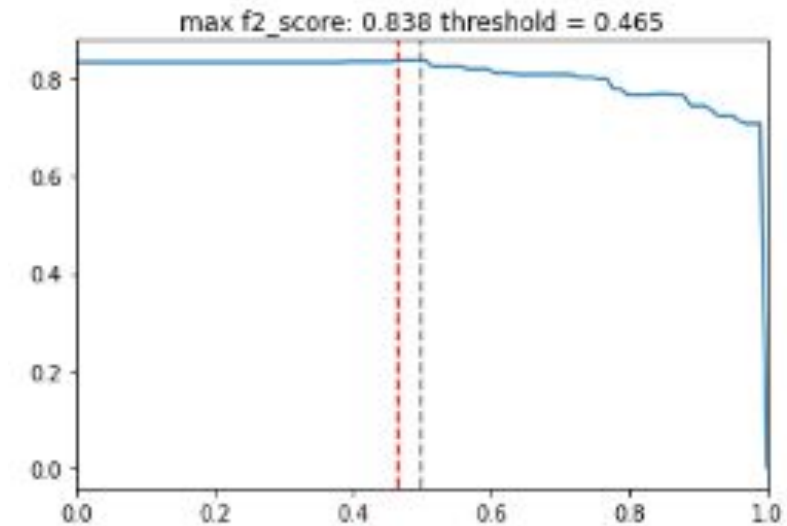
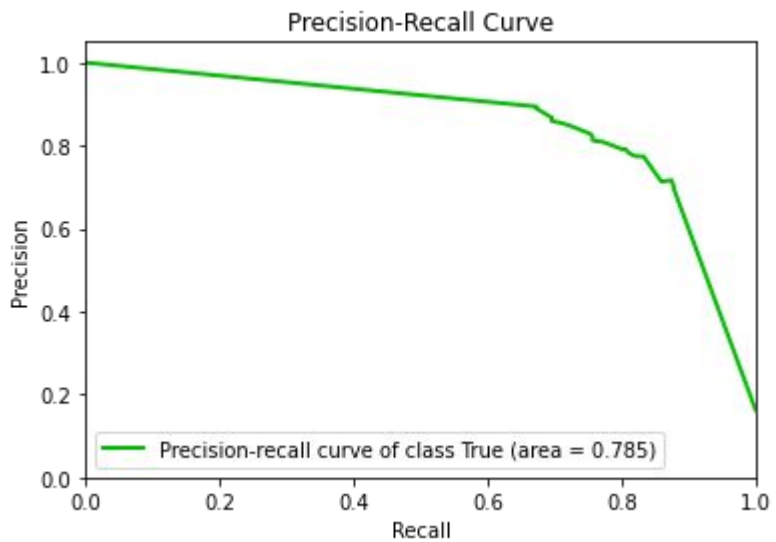
**criterion**: función de medida de calidad del corte

**class\_weight**: balanced

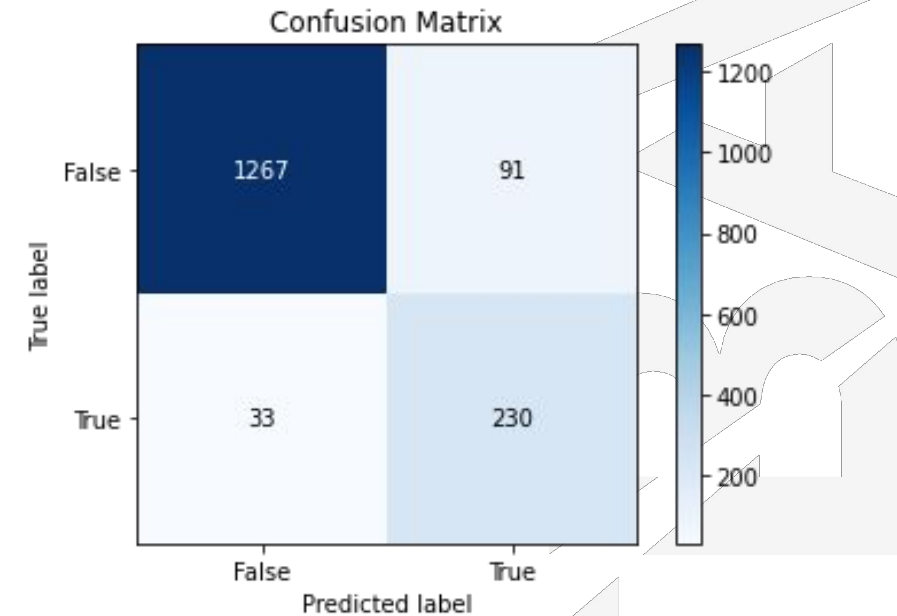
```
parameters={  
    'max_depth': [2, 5, 20],  
    'min_samples_leaf': [5, 10, 20, 100],  
    'criterion': ["gini", "entropy"],  
    "class_weight":["balanced"]  
}
```

```
search3.best_params_  
{  
    'class_weight': 'balanced',  
    'criterion': 'gini',  
    'max_depth': 20,  
    'min_samples_leaf': 5  
}
```

precision = 0.7165  
recall = 0.8745  
f1 = 0.7877



```
### threshold = 0.465  
accuracy: 0.924  
f2: 0.838  
AUPRC: 0.785  
### threshold = 0.500  
accuracy: 0.924  
f2: 0.838  
AUPRC: 0.785
```



# Extra Trees Classifier

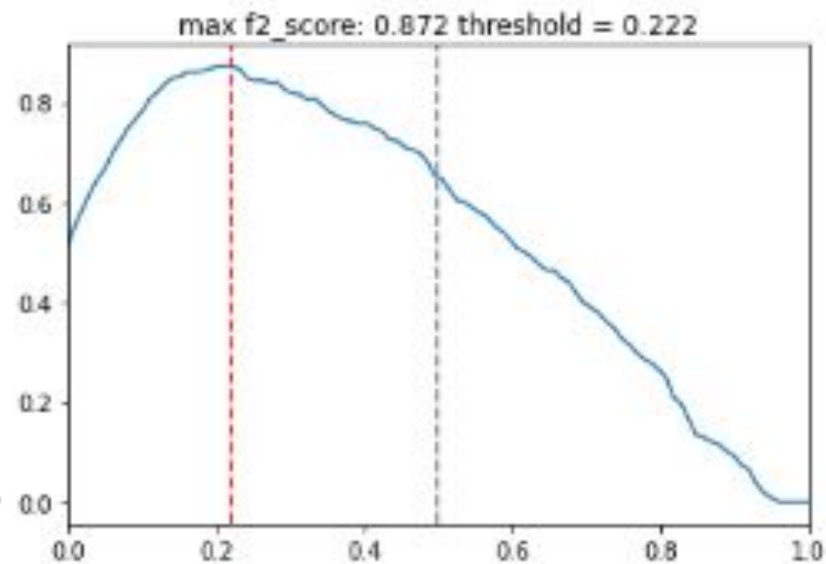
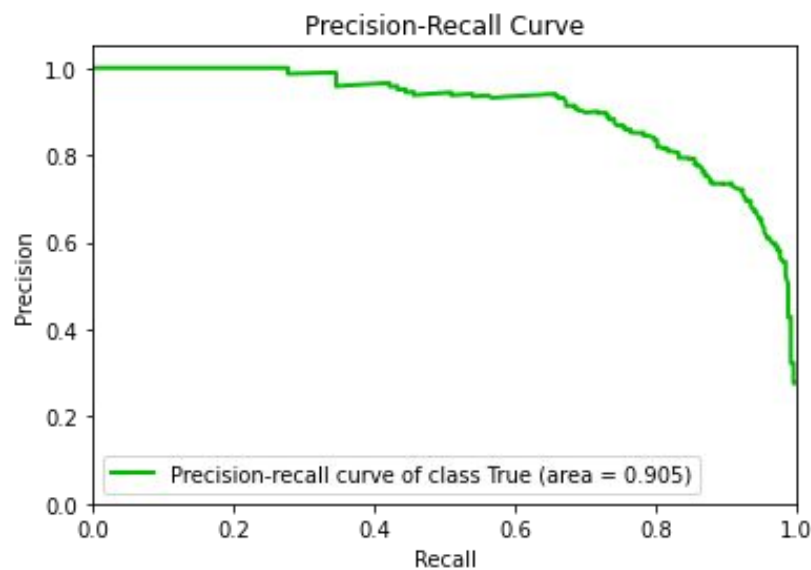
**max\_depth**: número máximo de profundidad del árbol

**criterion**: función de medida de calidad del corte

**class\_weight**: balanced

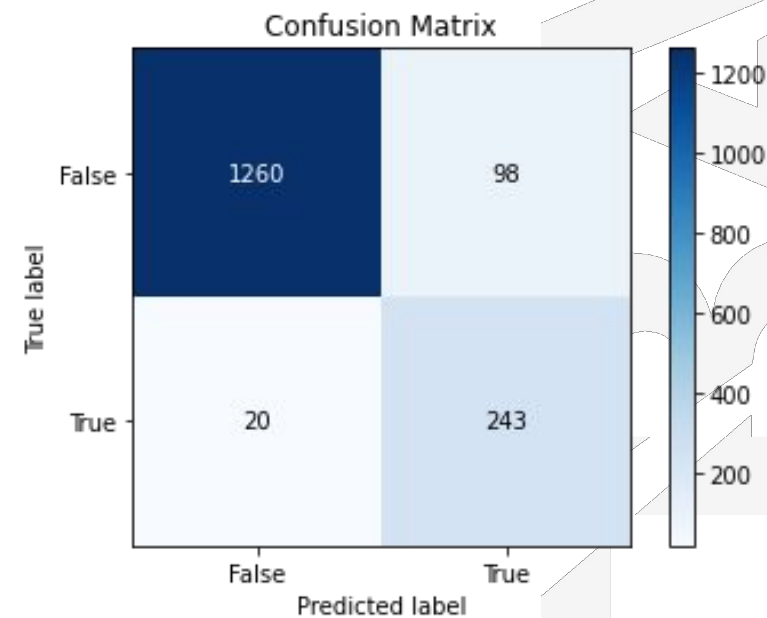
```
parameters = {  
    "class_weight":["balanced"],  
    'max_depth':[25, 40, 2],  
    'criterion' : ["gini", "entropy"]  
}
```

```
gridsearch4.best_params_  
  
{'class_weight': 'balanced', 'criterion': 'entropy', 'max_depth': 25}
```



```
### threshold = 0.222
accuracy: 0.927
f2: 0.872
AUPRC: 0.905
### threshold = 0.500
accuracy: 0.930
f2: 0.654
AUPRC: 0.905
```

```
precision = 0.7126
recall = 0.9240
f1 = 0.8046
```

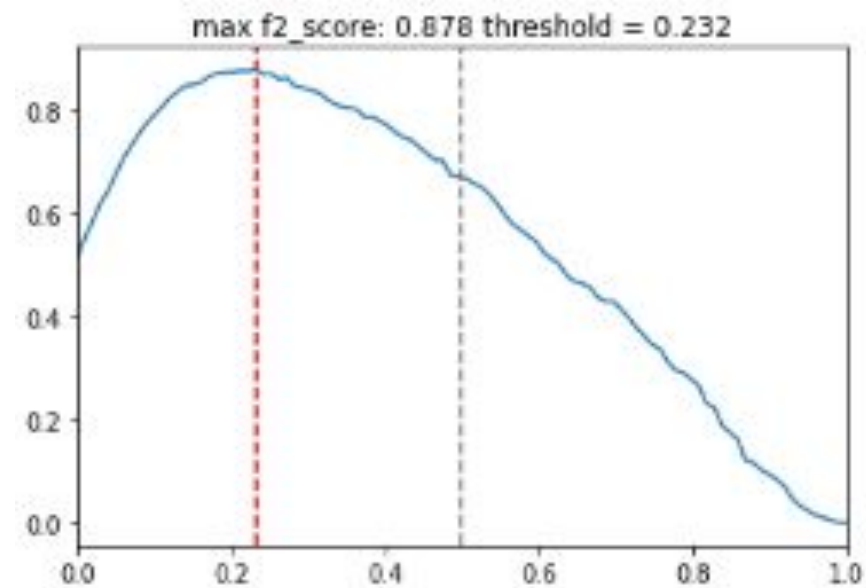
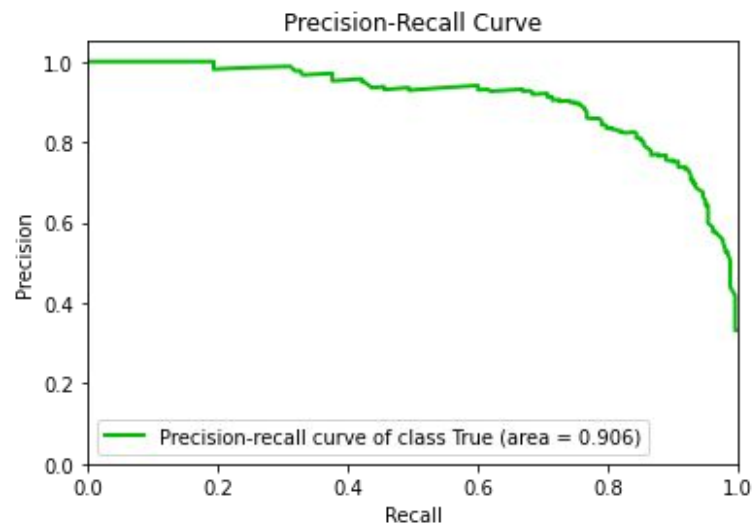




	params	AUPRC	EXT__class_weight	EXT__criterion	EXT__max_depth
4	{'EXT__class_weight': 'balanced', 'EXT__criter...	0.907509	balanced	entropy	40
0	{'EXT__class_weight': 'balanced', 'EXT__criter...	0.901658	balanced	gini	25
1	{'EXT__class_weight': 'balanced', 'EXT__criter...	0.900429	balanced	gini	40
3	{'EXT__class_weight': 'balanced', 'EXT__criter...	0.898103	balanced	entropy	25
2	{'EXT__class_weight': 'balanced', 'EXT__criter...	0.711208	balanced	gini	2
5	{'EXT__class_weight': 'balanced', 'EXT__criter...	0.684194	balanced	entropy	2

	feature	importance			
16	Total_Trans_Ct	0.1713			
15	Total_Trans_Amt	0.1259			
12	Total_Revolving_Bal	0.1048			
17	Total_Ct_Chng_Q4_Q1	0.0735			
8	Total_Relationship_Count	0.0598			
14	Total_Amt_Chng_Q4_Q1	0.0489	31	Total_Ct_Chng_Q4_Q1_right	0.0021
18	Avg_Utilization_Ratio	0.0482	23	Months_Inactive_12_mon_right	0.0021
9	Months_Inactive_12_mon	0.0419	30	Total_Trans_Ct_right	0.0000
10	Contacts_Count_12_mon	0.0400	19	Customer_Age_right	0.0000
11	Credit_Limit	0.0341	27	Avg_Open_To_Buy_right	0.0000
13	Avg_Open_To_Buy	0.0312	22	Total_Relationship_Count_right	0.0000
0	Customer_Age	0.0311	26	Total_Revolving_Bal_right	0.0000
7	Months_on_book	0.0303	25	Credit_Limit_right	0.0000
2	Dependent_count	0.0278	21	Months_on_book_right	0.0000
3	Education_Level	0.0265	20	Dependent_count_right	0.0000
5	Income_Category	0.0257	32	Avg_Utilization_Ratio_right	0.0000
4	Marital_Status	0.0235			
1	Gender	0.0192			
29	Total_Trans_Amt_right	0.0113			
24	Contacts_Count_12_mon_right	0.0086			
6	Card_Category	0.0082			
28	Total_Amt_Chng_Q4_Q1_right	0.0040			

1. Valores bajos
2. Si tiene sentido para el problema
3. No cambia el resultado para peor



```

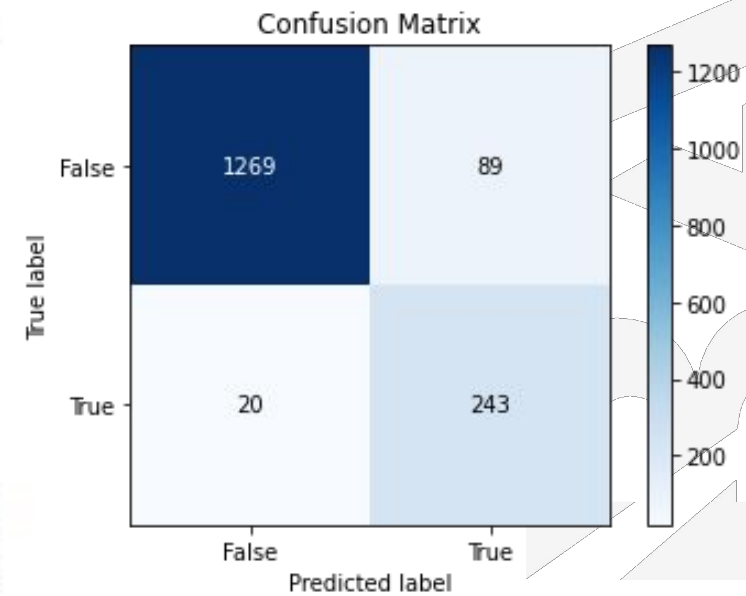
VW threshold = 0.232
accuracy: 0.933
F2: 0.878
AUPRC: 0.906
VW threshold = 0.500
accuracy: 0.932
F2: 0.671
AUPRC: 0.906

```

```

precision = 0.7319
recall = 0.9248
f1 = 0.8168

```



# Random Forest Classifier

`n_estimators`: número de árboles

`RFC__max_features`: número de variables a considerar cuando se hace el mejor corte

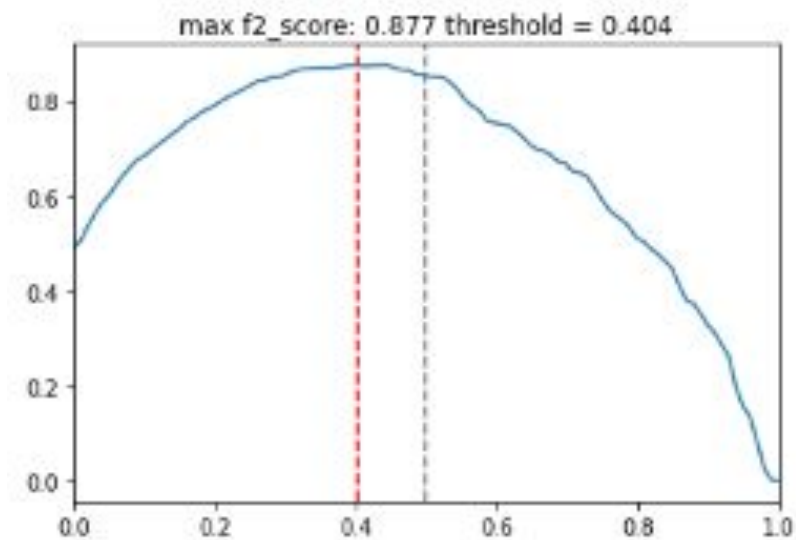
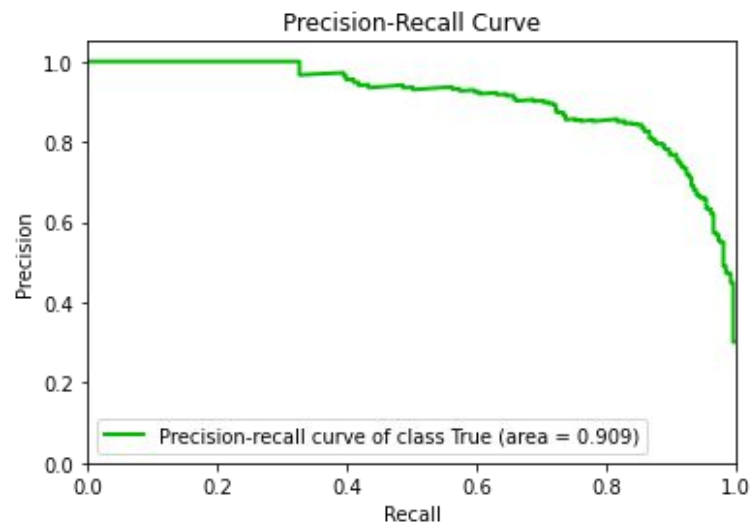
`RFC__max_depth`: número máximo de profundidad del árbol

`RFC__criterion`: función de medida de calidad del corte

`class_weight`: "balanced"

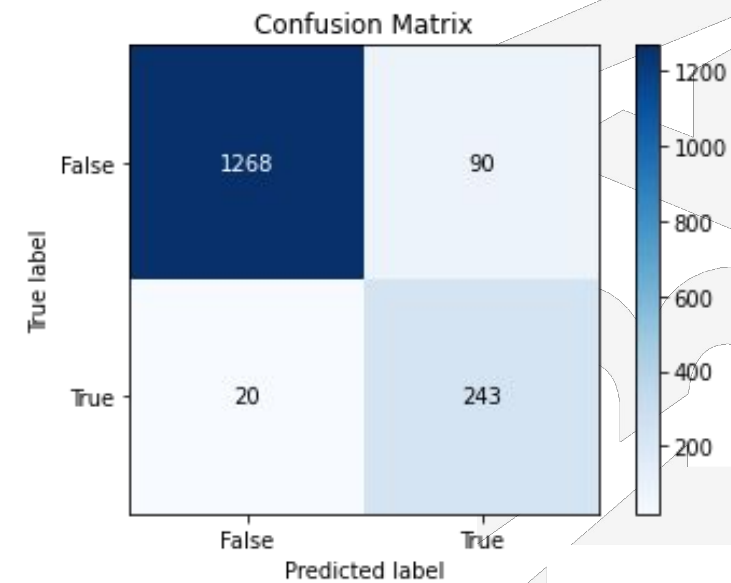
```
param_gridRFC={  
    'RFC__n_estimators': [200, 500],  
    'RFC__max_features': ['auto', 'sqrt'],  
    'RFC__max_depth' : [2, 5, 9],  
    'RFC__criterion' :['gini', 'entropy']  
}
```

```
gridsearch2.best_params_  
  
{'criterion': 'entropy',  
  'max_depth': 9,  
  'max_features': 'log2',  
  'n_estimators': 500}
```



```
### threshold = 0.404
accuracy: 0.932
f2: 0.877
AUPRC: 0.909
### threshold = 0.500
accuracy: 0.948
f2: 0.855
AUPRC: 0.909
```

```
precision = 0.7297
recall = 0.9240
f1 = 0.8154
```



	params	AUPRC	RFC__criterion	RFC__max_depth	RFC__max_features	RFC__n_estimators
22	{'RFC__criterion': 'entropy', 'RFC__max_depth'...	0.9204	entropy	9	sqrt	200
21	{'RFC__criterion': 'entropy', 'RFC__max_depth'...	0.9189	entropy	9	log2	500
23	{'RFC__criterion': 'entropy', 'RFC__max_depth'...	0.9172	entropy	9	sqrt	500
11	{'RFC__criterion': 'gini', 'RFC__max_depth': 9...	0.9168	gini	9	sqrt	500
9	{'RFC__criterion': 'gini', 'RFC__max_depth': 9...	0.9164	gini	9	log2	500
20	{'RFC__criterion': 'entropy', 'RFC__max_depth'...	0.9162	entropy	9	log2	200

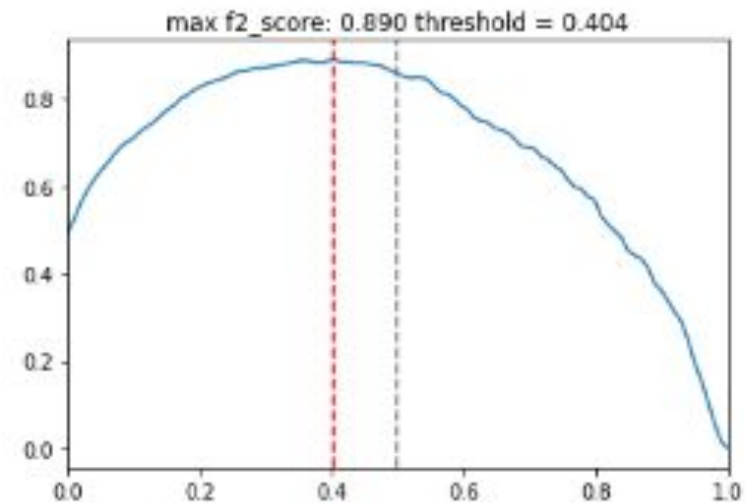
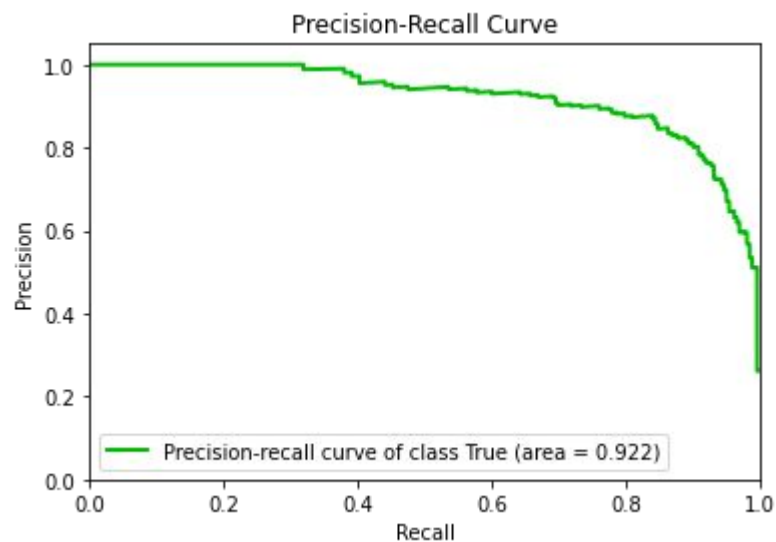


# Feature importance

1. Valores bajos
2. Si tiene sentido para el problema
3. No cambia el resultado para peor

	feature	importance
15	Total_Trans_Amt	0.2150
16	Total_Trans_Ct	0.2088
12	Total_Revolving_Bal	0.1027
17	Total_Ct_Chng_Q4_Q1	0.1010
18	Avg_Utilization_Ratio	0.0676
14	Total_Amt_Chng_Q4_Q1	0.0603
8	Total_Relationship_Count	0.0457
13	Avg_Open_To_Buy	0.0302
11	Credit_Limit	0.0295
9	Months_Inactive_12_mon	0.0294
10	Contacts_Count_12_mon	0.0230
0	Customer_Age	0.0214
7	Months_on_book	0.0168
1	Gender	0.0073
29	Total_Trans_Amt_right	0.0068
2	Dependent_count	0.0066
5	Income_Category	0.0065
3	Education_Level	0.0056
4	Marital_Status	0.0052
24	Contacts_Count_12_mon_right	0.0047

28	Total_Amt_Chng_Q4_Q1_right	0.0028
6	Card_Category	0.0021
31	Total_Ct_Chng_Q4_Q1_right	0.0009
23	Months_Inactive_12_mon_right	0.0003
19	Customer_Age_right	0.0000
30	Total_Trans_Ct_right	0.0000
27	Avg_Open_To_Buy_right	0.0000
22	Total_Relationship_Count_right	0.0000
26	Total_Revolving_Bal_right	0.0000
25	Credit_Limit_right	0.0000
21	Months_on_book_right	0.0000
20	Dependent_count_right	0.0000
32	Avg_Utilization_Ratio_right	0.0000



```

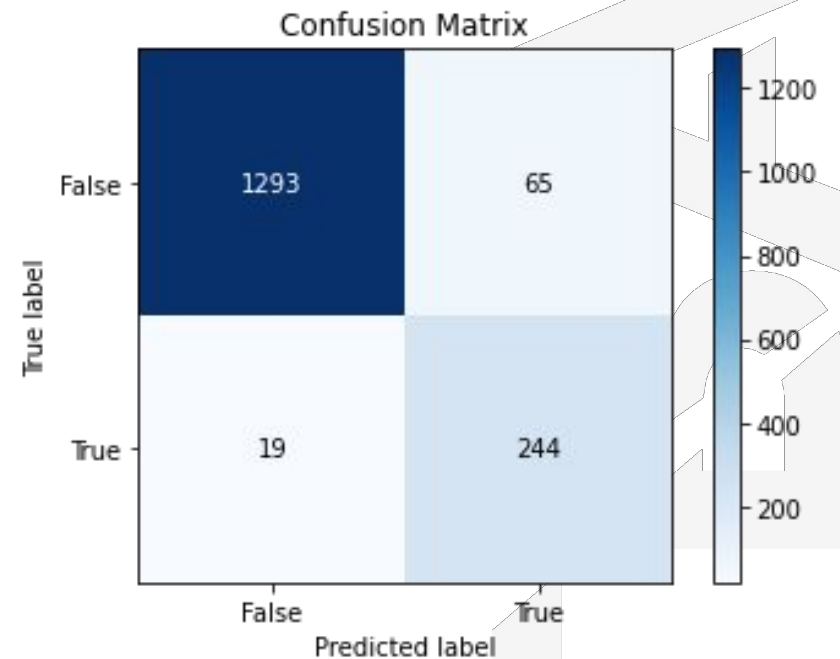
### threshold = 0.404
accuracy: 0.940
f2: 0.890
AUPRC: 0.922
### threshold = 0.500
accuracy: 0.952
f2: 0.859
AUPRC: 0.922

```

```

precision = 0.7562
recall = 0.9316
f1 = 0.8348

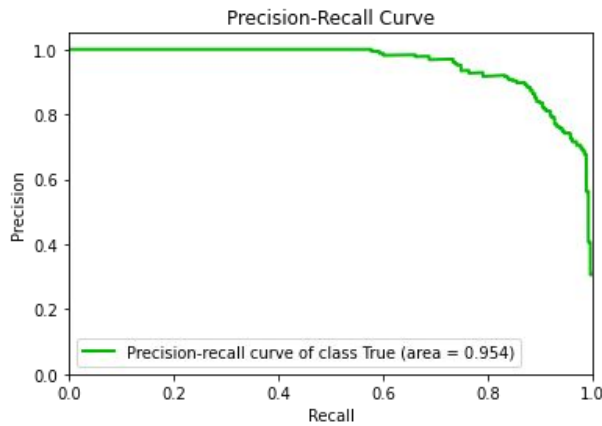
```





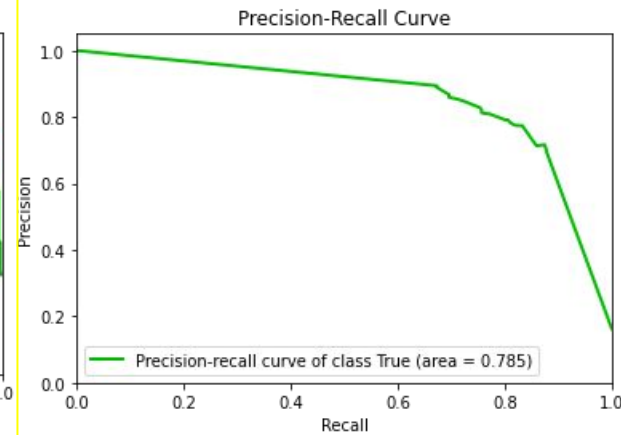
# Comparamos modelos

## LGBMClassifier



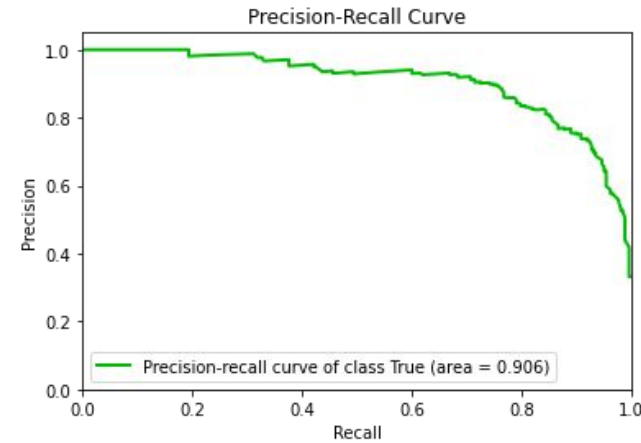
```
precision = 0.7060  
recall = 0.9772  
f1 = 0.8198
```

## Decision Tree Classifier



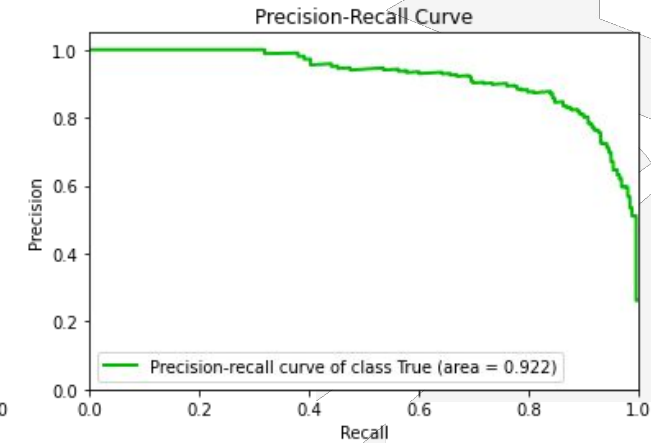
```
precision = 0.7165  
recall = 0.8745  
f1 = 0.7877
```

## Extra Trees Classifier



```
precision = 0.7319  
recall = 0.9240  
f1 = 0.8168
```

## Random Forest Classifier



```
precision = 0.7562  
recall = 0.9316  
f1 = 0.8348
```

# Conclusión





Card category: blue  
Marital status: married / single  
Months inactive 12 mon: 2 / 3  
Gender: M / F  
Education level: High school / graduate

Less than \$40K  
Dependent count: 2 / 3  
Total relationship count: 2 / 3  
Contact count: 2 / 3 / 4

Total trans amt: 0-3000  
Total trans ct: 0-1000

CLIENTNUM	Attrition_Flag	Customer_Age	Gender	Dependent_count	Education_Level	Marital_Status	Income_Category	Card_Category	Months_on_book	Total_Relationship_Count	Months_Inactive_12_mon	Contacts_Count_12_mon	Credit_Limit
9616 751155783	Existing Customer	47	F	2	High School	Married	\$40K - \$60K	Blue	40	2	1	3	4074.0000
7865 714998508	Existing Customer	53	M	3	Uneducated	Married	\$120K +	Blue	36	2	2	0	3742.0000
73 820582308	Existing Customer	42	M	5	Uneducated	Married	\$80K - \$120K	Blue	37	6	2	2	22913.0000
596 720370533	Attrited Customer	55	M	3	Uneducated	Married	\$60K - \$80K	Blue	44	3	2	2	2323.0000
114 711844758	Existing Customer	48	M	3	Graduate	Single	\$80K - \$120K	Blue	35	6	1	0	13551.0000

Total_Revolving_Bal	Avg_Open_To_Buy	Total_Amt_Chng_Q4_Q1	Total_Trans_Amt	Total_Trans_Ct	Total_Ct_Chng_Q4_Q1	Avg_Utilization_Ratio
1868	2206.0000	0.6880	15005	118	0.7610	0.4590
1454	2288.0000	0.8940	5326	62	0.6760	0.3890
1528	21385.0000	0.4140	1394	35	0.5220	0.0670
0	2323.0000	0.7370	804	15	0.5000	0.0000
1294	12257.0000	0.7910	1388	37	1.0560	0.0950

CLIENTNUM	Attrition_Flag
751155783	False
714998508	False
820582308	False
720370533	True
711844758	False

MUCHAS GRACIAS