



Universidad de Antioquia

Ingeniería de Sistemas

Introducción a la inteligencia artificial

PROYECTO ENTREGA 2

PRESENTA:

Jonatan Stiven Restrepo Lora

Tutor Principal:
Raul Ramos Pollan

Resumen

Estos registros se guardaron y se sometieron a una evaluación con el propósito de

Leemos los archivos, el primero contiene toda la información, el segundo los resultados y procedemos a mirar qué columnas tienen en común para hacer una unión por el campo en común.

['customer_ID',	0000009966d5597952...	0.0	2.0	1.0	0.0	4.0	0.0	1.0	CR	0	NULL	6.0
'S_2',	0000009966d5597952...	0.0	2.0	1.0	0.0	4.0	0.0	1.0	CR	0	NULL	6.0
'P_2',	0000009966d5597952...	0.0	2.0	1.0	0.0	4.0	0.0	1.0	CR	0	NULL	6.0
	0000009966d5597952...	0.0	2.0	1.0	0.0	4.0	0.0	1.0	CR	0	NULL	6.0

`[00000000000000789C...| 0.0| 7.0| 1.0| 0.0| 4.0| 0.0| 1.0| 0.0| 0|0|0|1| 6.0]`

```
data.select('customer_ID', 'B_30', 'B_38', 'D_114', 'D_116', 'D_117', 'D_120', 'D_126', 'D_63', 'D_64', 'D_66', 'D_68').show()
```

```
customer TD R 30 R 38 D 114 D 116 D 117 D 120 D 126 D 63 D 64 D 66 D 68
```

Guardamos el conjunto de datos en un excel para poder trabajar con el más f3cil y poder hacer la uni3n con los labels

Merge Datasets

Una vez los datos divididos, procedemos a hacer un merge de la muestra extraída y **train_labels** para así poder agregar el objetivo correspondiente al **customer_ID** de nuestro conjunto de datos para proceder a guardarlo dado que este sera el dataset con el que estaremos trabajando.

```
print(f'cantidad de registros: {len(pd_label)}')
```

cantidad de registros: 458913

```
data_full = pd_muestra.merge(pd_label, on='customer_ID')
data_full.head()
```

		customer_ID	S_2	P_2	D_39	B_1	B_2	R_1	S_3	D_41	B_3	...	D_137	D_138
0	000adf2938f771f75a581b65107024eddeae70684778c0...	2017-04-25	0.885999	0.009377	0.008894	1.006219	0.008716	0.065198	0.001650	0.009946	...		NaN	NaN
1	001a152e1893ab8372e7c9627c9de2e024399f2660d5d8...	2017-11-10	0.750890	0.037378	0.045959	1.001278	0.000312	0.087636	0.004212	0.001730	...		NaN	NaN
2	001e2ceaf1421f1477c0de9ba1c9357b9d278f7b670ab7...	2018-03-26	0.607227	0.002678	0.001271	0.812964	0.005415	NaN	0.005049	0.006548	...		NaN	NaN
3	0034f7e366a41d2500643c7dd0faa6302ce944743ccdf5...	2018-01-07	0.775546	0.003486	0.234352	0.040793	0.256024	0.172966	0.000081	0.225639	...		NaN	NaN
4	00394e07aa3f71174f8bedfd16d64f194c80ad9445e17f...	2017-04-27	0.358659	0.001814	0.028545	1.004951	0.009079	0.765318	0.003546	0.009843	...		NaN	NaN

5 rows × 191 columns

```
print(f'cantidad de registros: {len(data_full)}')
```

cantidad de registros: 6041

Una vez comprobado que el merge lo altero el DF y se procede a guardarlos en un csv y txt

```
data_full.to_excel('data_full.xlsx', index=False)
```

```
data_full.to_csv('data_full.txt', sep='\t', index=False)
```

Figura 2: Muestra completa

Limpeza de Datos

En el proceso de limpieza de datos, comenzamos por identificar las variables categóricas, a pesar de que ya se habían especificado en el conjunto de datos original como 'B_30', 'B_38', 'D_114', 'D_116', 'D_117', 'D_120', 'D_126', 'D_63', 'D_64', 'D_66', 'D_68'. Estas variables se dividieron en dos categorías: nominales y ordinales. En primer lugar, se llevó a cabo la limpieza y transformación de las variables nominales.

Dentro de las variables nominales, se identificaron 'S_2', 'D_63', 'D_64'. Se observó que 'S_2' representaba fechas, por lo que se procedió a ajustar el tipo de dato utilizando la biblioteca pandas. Para las demás variables nominales, se completaron los valores faltantes utilizando la moda de los datos.

Posteriormente, se aplicó el método de codificación One-Hot Encoding para organizar estas variables categóricas.

```
df = pd.get_dummies(df, columns=['D_63', 'D_64'])
```

```
df.rename(columns={'D_64_1': 'D_64_I'}, inplace=True)
```

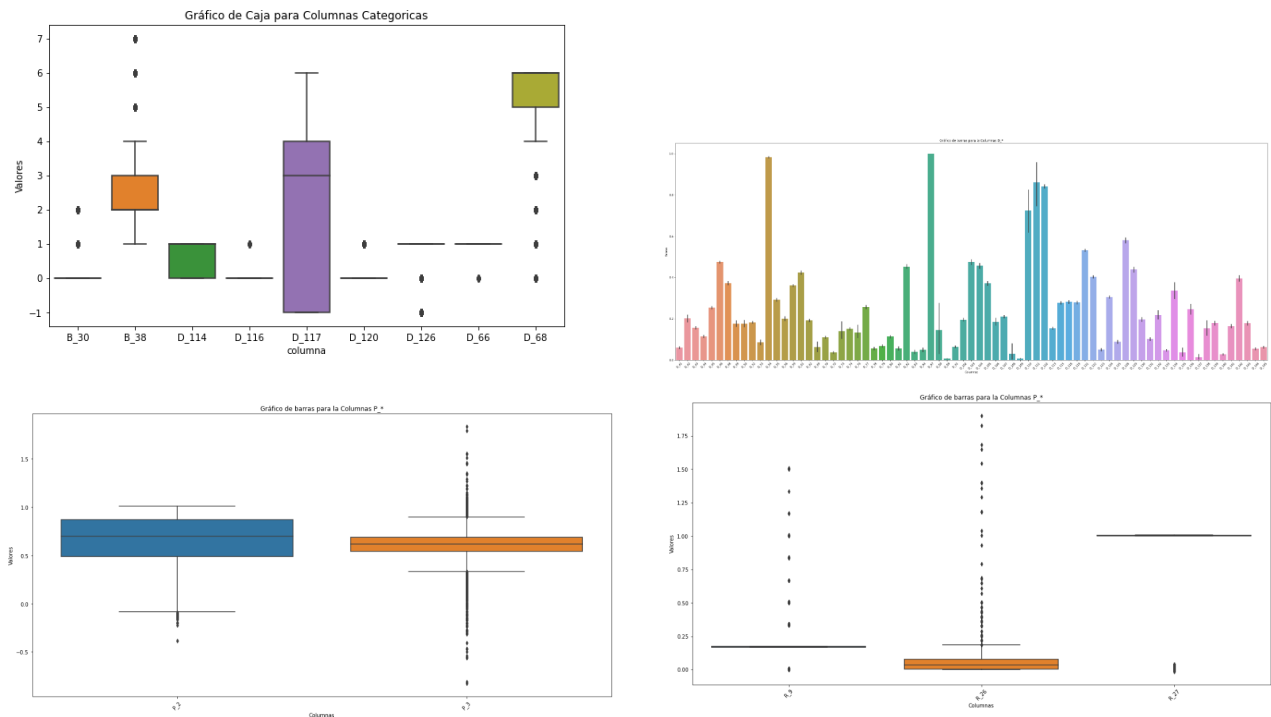
```
df[['D_63_CO', 'D_63_CR', 'D_63_CL', 'D_63_XZ', 'D_63_XZ', 'D_63_XM', 'D_63_XL', 'D_64_O', 'D_64_U', 'D_64_R',  
    'D_64_I', 'target']]
```

	D_63_CO	D_63_CR	D_63_CL	D_63_XZ	D_63_XZ	D_63_XM	D_63_XL	D_64_O	D_64_U	D_64_R	D_64_I	target
0	1	0	0	0	0	0	0	1	0	0	0	0
1	1	0	0	0	0	0	0	0	1	0	0	0
2	1	0	0	0	0	0	0	1	0	0	0	1
3	1	0	0	0	0	0	0	1	0	0	0	0
4	0	0	1	0	0	0	0	1	0	0	0	1
...
6036	1	0	0	0	0	0	0	0	1	0	0	0
6037	1	0	0	0	0	0	0	0	1	0	0	1
6038	1	0	0	0	0	0	0	0	0	1	0	1
6039	1	0	0	0	0	0	0	0	1	0	0	0
6040	0	1	0	0	0	0	0	1	0	0	0	0

6041 rows x 12 columns

Figura 3: Manejo de variables categoricas

Luego, procedimos a analizar el resto de las variables para verificar si contenían valores atípicos, con el fin de determinar el método más apropiado para completar los valores faltantes. Si se identificaban datos atípicos y, después de examinar la distribución de los datos, se determinaba que llenar los valores faltantes con la moda era la mejor opción para evitar una desviación significativa.



Cuadro 2: Analisis de los grupos

Conclusiones

Se ha completado la separación y limpieza de los datos de trabajo. Los siguientes pasos a seguir serán la evaluación de las relaciones entre las variables y el objetivo con el fin de determinar qué algoritmo de clasificación puede ser más adecuado para abordar la pregunta en cuestión.

- [Folder de trabajo de los notebooks](#)

Nota: los notebooks fueron trabajados en local, si los va a ejecutar en colab asegurece que la ruta de los archivos sea la correcta.

- [video explicatorio](#)