

TRABAJO PRACTICO FINAL PROGRAMACION ORIENTADA A OBJETOS

INTRODUCCION

Bases de datos:

Precios claros es un conjunto de datasets realizado a partir de los precios provistos por el gobierno de los productos en las distintas cadenas de supermercados. Este parte de la iniciativa de la Dirección Nacional de Defensa del Consumidor de que los consumidores cuenten con mayor información a la hora de tomar decisiones con respecto a los productos que adquieren. Particularmente estos datos pertenecen al año 2020, lo cual es importante tener en consideración a la hora de mirar los precios dada su elevada diferencia con la actualidad.

La base de datos seleccionada cuenta con 6 datasets: uno de productos, uno con las sucursales y 3 con los precios de abril en distintos periodos, los cuales fueron relevados de la página del gobierno mediante técnicas de scraping por Open Data Córdoba, y por otro lado uno adicional que contiene las provincias y códigos de las mismas.

Link: <https://www.kaggle.com/datasets/linmqn/precios-claros-precios-de-argentina?select=sucursales.csv>

Variables:

A continuación se da una breve definición de cada variable en el dataset precios claros, el cual se utilizara a lo largo de toda la investigación.

Producto

- id** : productos_id : numero de identificación del producto.
- categoría** : categoría del producto (nuestros productos son los que forman parte de la canasta basica: Agua, Leche, Arroz, Harinas y Pastas).
- márca** : marca del producto.
- presentación** : presentación del producto en cantidad (numeros).
- unidad_presentacion** : unidad de medida de la presentación (cc, lt, kg o gr).
- precio** : precio del producto en el 2020.

Sucursal/Comercio

- sucursal_id** : numero de identificación de la sucursal del producto.
- sucursal_tipo** , **tipo_comercio** : tipo de comercio de la sucursal (Supermercado, Hipermercado o Autoservicio).
- cadena** : es la cadena de supermercado en la que se vende el producto. Este es el target de nuestro modelo, por ello seleccionamos 3 cadenas con distintas categorías: Disco dirigido a personas de alto poder adquisitivo, que es de alta categoría, Hipermercado Carrefour a personas poder adquisitivo medio y Supermercados DIA de menor poder adquisitivo).

Variable dependiente -> es el target que queremos, el cual depende del resto de variables: **cadena**. **Variables independientes** -> son las variables que nos van a permitir predecir el target, es decir, todo el resto de variables mencionadas: **categoría**, **márca**, **presentación**, **unidad_presentacion**, **precio** y **tipo_comercio**.

Objetivo:

El objetivo de este trabajo es poder analizar los precios de los diferentes tipos de productos de la canasta basica, teniendo en cuenta la categoría, marca, tipo de cadena, presentación y cadena en el que se realiza la venta del producto, para poder evaluar en que comercios y que marcas tienen los productos mas baratos.

Ademas, la idea es crear un modelo de regresion logistica que nos permita predecir en base a las variables a cual cadena pertenece el producto, ya sea Disco, Hipermercado Carrefour o Supermercados DIA. Dado que pertenecen a diferentes categorías de supermercados, sus consumidores objetivo son distintos, por lo que sus precios tambien varian.

Es por eso que el objetivo del modelo es predecir a cual de ellos pertenece el producto.

Hipotesis:

Respecto a las variables seleccionadas se espera poder ver que en el comercio de Disco, dado que es una cadena dirigido a personas de alto poder adquisitivo, sus productos tambien tendran mayores precios. Lo mismo sucedera son Supermercados DIA que al ser una cadena dirigida a consumidores de menor poder adquisitivo que los anteriores, se espera encontrar que sus precios tambien son menores, mientras que en Carrefour se espera encontrar un precio intermedio que en los dos anteriores.

Tranformaciones a realizar:

• Load & Merge

En primera instancia se van a cargar todos los datasets, para luego poder unirlos y filtrarlos segun los datos necesarios para el analisis, en este caso, las categorías de productos de canasta basica: Agua, Leche, Arroz, Harinas y Pastas, y la cadena de comercio: Disco, Hipermercado Carrefour y Supermercados DIA. Ademas, como para nuestro modelo necesitamos la columna de presentación separado de la unidad en la que esta medida, separaremos la columna en esta instancia. El dataframe resultante es el que sera utilizado para el trabajo, el cual se llamara **precios_claros**.

• Preprocesamiento

Una vez listo el dataset es necesario realizar un preprocesamiento de los datos. En esta instancia se evaluaran los posibles nulos en nuestras variables realizar el tratamiento adecuado de los mismos (drop, imputer, etc), y se evaluaran tambien los posibles duplicados y outliers.

• Analisis Inicial

Se muestra un analisis de nuestros datos, evaluando tendencias, frecuencias, etc, con el objetivo de tener un mejor conocimiento de los datos a tratar.

• Escalado

Se escalan las variables numericas **precio** y **presentación** para que sean comparables entre si, a pesar de que su rango de datos en primera instancia sea distinto.

• Encoding

Se tratan las variables categoriales, creando distintas columnas binarias con la presencia o no de la variable en el producto. Esto es fundamental para el modelo bucado. Ademas, aplicamos el encoding a la variable target.

• Modelo

Realizamos el primer modelo de regresion logistica y obtenemos un score mediante un cross validation, el cual se irá mejorando mediante diferentes métodos para lograr un accuracy mas elevado. Haremos una matriz de confusion y veremos accuracy, precision y recall. Ademas haremos la predicción de una instancia que no existe para probarlo.

```
In [1]: !load_ext autoreload
!autoreload 2

Importamos las librerias que vamos a necesitar.
```

```
In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

1 - LOAD & MERGE

```
In [3]: import os
files = [f for f in os.listdir("data") if f.endswith(".csv")]
print(files)

['productos_categoria.csv', 'clean_precios_claros.csv', 'precios_20200426_20200426.csv', 'precios_20200412_20200413.csv', 'sucursales', 'precios_20200419_20200419.csv']
```

```
In [4]: from load_csv import load_all_data
data = load_all_data("data/")
data.keys()
```

```
Out[4]: dict_keys(['productos_categoria', 'clean_precios_claros', 'precios_20200426_20200426', 'precios_20200412_20200413', 'sucursales', 'precios_20200419_20200419'])

Asigno los datasets en variables haciendole una copia para no sobrescribir el original.
```

Sucursales

```
In [5]: sucursales = data['sucursales'].copy()
sucursales.head()
```

id	comercio	bandera	banderaDescripcion	comercioRazonSocial	provincia	localidad	direccion	lat	lng
0	1-7	1	1	Super MAMI	Dinosaurio S.A.	AR-X	SALSIPIUEDES	E53 1011 None	-31.126667 -64.295260
1	10-1	10	1	Hipermercado Carrefour	INC S.A.	AR-B	San Isidro	Bernardo De Irigoyen 2647	-34.491345 -58.589025
2	10-1	10	1	Hipermercado Carrefour	INC S.A.	AR-B	Hurlingham	Av. Vergara 1910	-34.620610 -58.633769
3	10-11	10	1	Hipermercado Carrefour	INC S.A.	AR-B	Malvinas Argentinas	Av. Arturo Illia 3770	-34.528883 -58.701631
4	10-112	10	1	Hipermercado Carrefour	INC S.A.	AR-A	Salta	20 De Febrero 247	-24.789072 -65.413699

Productos

```
In [6]: productos = data['productos_categoria'].copy()
productos.head()
```

id	_id/soid	nombre	categoría	marca	presentacion	id
0	5cbcd9be7af152186cd0c837	Acetile de Girasol Cañuelas 1.5 Lt	Acetile	CAÑUELAS	1.5 lt	7792180001665
1	5cbcd9be7af152186cd0c837	Acetile de Girasol Cañuelas 900 Cc	Acetile	CAÑUELAS	900.0 cc	7792180001641
2	5cbcd9be7af152186cd0c837	Acetile de Girasol Cocinero 1.5 Lt	Acetile	COCINERO	1.5 lt	7790060023684
3	5cbcd9be7af152186cd0c837	Acetile de Girasol Cocinero 900 Ml	Acetile	COCINERO	900.0 ml	7790060021050
4	5cbcd9be7af152186cd0c837	Acetile de Girasol Natural 1.5 Lt	Acetile	NATURA	1.5 lt	7790272001029

Precios

Cargo los datasets de los precios de abril. Como se contaba con varios datasets del mes, se opto por concatenarlos uno debajo del otro para poder tratar los datos y realizar el analisis.

```
In [7]: precios = pd.concat([data['precios_20200412_20200413'].copy()], data['precios_20200419_20200419'].copy()], data['precios'].head())

Out[7]: precio producto_id sucursal_id
0 29.90 00000000001663 2-1-014
1 29.90 00000000002288 2-1-032
2 39.90 00000000002288 2-1-096
3 499.99 00000000205870 9-1-686
4 519.99 000000023870 9-2-248
```

- Merge

Unimos los distintos dataset para poder realizar el analisis.

Para poder realizarlo en primer lugar voy a establecer como indice de los dataframes los id, con el objetivo de que al unirlo las columnas no se reiptan. Y pasamos la columna id de productos a tipo objeto para poder unirlos con precios.

```
In [8]: precios = precios.set_index("producto_id")
sucursales = sucursales.set_index("id")

Filtrado de productos
```

Ademas, para no unir datos que no seran utilizados, en primera instancia filtraremos el dataset de productos para quedarnos unicamente con los lacteos.

```
In [9]: list[productos.categoria.unique()]

Out[9]: list[Array(['Acetile', 'Acetile de oliva', 'Aderezos', 'Agua', 'Arroz', 'Bebidas y niños', 'Bebidas con alcohol', 'Bebidas sin alcohol', 'Bebidas deportivas', 'Jugos en polvo', 'Cacao y café', 'Carnes congeladas', 'Otros congelados', 'Conservas', 'Endulzantes', 'Especias', 'Golosinas', 'Grasas', 'Harinas', 'Infusiones', 'Lacteos', 'Leche', 'Leche en polvo', 'Limpieza', 'Manteca', 'Mascotas', 'Mermeladas y dulces', 'Cereales azucarados', 'Panificados', 'Pan', 'Pastas', 'Perfumeria', 'Postres y reposteria', 'Premezclas', 'Quesos', 'Salsas', 'Sopas', 'Verduras congeladas', 'Yogur'], dtype=object)]
```

Vemos las categorías que tiene nuestro dataset, ya que para el analisis buscamos analizar los productos de la canasta basica. Al visualizar las diferentes categorías se observa que una de ellas pertenece a NANS. Antes de filtrar el dataset analizaremos estos NANS para ver si son relevantes para el trabajo.

```
In [10]: productos[productos.categoria.isna()]

Out[10]: _id/soid nombre categoría marca presentacion id
704 5cbcd9be7af152186cd0c837 Mix de Frutos Secos y Semillas Nutritivo Natur... NaN ARCOR 27.0 gr 7790580110413
998 5cbcd9be7af152186cd0c837 Crema de Avelana con Cacao Nutella 140 Gr NaN NUNO 140.0 gr 789024395232
```

Como este producto no son parte de la canasta basica, no es relevante que falte su categoría, ya que no vamos a utilizarlo para el analisis. Podemos ignorarlo o eliminarlo. Optaremos por eliminarlo.

```
In [11]: productos = productos.dropna(subset = ["categoria"])
productos.isna().sum()
```

```
Out[11]: id/soid 0
nombre 0
id 0
categoría 0
marca 0
presentacion 0
id 0
dtype: int64
```

Una vez eliminados, filtraremos por las categorías pertenecientes a la canasta basica.

```
In [12]: canasta_basica = ["Agua", "Leche", "Arroz", "Harinas", "Pastas"]
productos = productos[productos.categoria.isin(canasta_basica)]
productos.head()
```

id	_id/soid	nombre	categoría	marca	presentacion	id
32	5cbcd9be7af152186cd0c837	Agua Mineral con Gas Villavicencio 1.5 Lt	Agua	VILLAVICENCIO	1.5 lt	7798062547788
33	5cbcd9be7af152186cd0c837	Agua Mineral sin Gas Bajo Sodio Ser 1.5 Lt	Agua	SER	1.5 lt	779031508706
34	5cbcd9be7af152186cd0c837	Agua Mineral sin Gas Bonagua 1.5 Lt	Agua	BONAGUA	1.5 lt	7790895003875
35	5cbcd9be7af152186cd0c837	Agua Mineral sin Gas Bonagua 500 Cc	Agua	BONAGUA	500.0 cc	7790895003868
36	5cbcd9be7af152186cd0c837	Agua Mineral sin Gas Eco de Los Andes 500 Cc	Agua	ECO DE LOS ANDES	500.0 cc	7792799000011

Verificamos que las categorías en el nuevo df sean las seleccionadas

```
In [13]: list[productos.categoria.unique()]

Out[13]: list[Array(['Agua', 'Arroz', 'Harinas', 'Leche', 'Pastas'], dtype=object)]

Ahora si uno los dataframes,
```

```
In [14]: precios_claros = pd.merge(productos, precios, how = "left", left_on = 'id', right_on = 'producto_id')\
.merge(sucursales, how = "left", left_on = 'sucursal_id', right_on = 'id')
precios_claros
```

```
Out[14]: _id/soid nombre categoría marca presentacion id precio sucursal_id comercio
0 5cbcd9be7af152186cd0c83c Agua Mineral con Gas Villavicencio 1.5 Lt Agua VILLAVICENCIO 1.5 lt 7798062547788 70.00 10-1-112 1
1 5cbcd9be7af152186cd0c83c Agua Mineral con Gas Villavicencio 1.5 Lt Agua VILLAVICENCIO 1.5 lt 7798062547788 85.00 10-1-142 1
2 5cbcd9be7af152186cd0c83c Agua Mineral con Gas Villavicencio 1.5 Lt Agua VILLAVICENCIO 1.5 lt 7798062547788 69.00 10-1-178 1
3 5cbcd9be7af152186cd0c83c Agua Mineral con Gas Villavicencio 1.5 Lt Agua VILLAVICENCIO 1.5 lt 7798062547788 75.50 10-1-233 1
4 5cbcd9be7af152186cd0c83c Agua Mineral con Gas Villavicencio 1.5 Lt Agua VILLAVICENCIO 1.5 lt 7798062547788 73.00 10-1-29 1
```

```
... ..
14242 5cbcd9be7af152186cd0c83c Rissotto 4 Quesos Galle 200 Gr Arroz GALLO 200.0 gr 7790070411754 139.00 49-1-2 4
14243 5cbcd9be7af152186cd0c83c Rissotto 4 Quesos Galle 200 Gr Arroz GALLO 200.0 gr 7790070411754 110.90 5-1-3
14244 5cbcd9be7af152186cd0c83c Rissotto 4 Quesos Galle 200 Gr Arroz GALLO 200.0 gr 7790070411754 139.00 50-1-2 5
14245 5cbcd9be7af152186cd0c83c Rissotto 4 Quesos Galle 200 Gr Arroz GALLO 200.0 gr 7790070411754 123.55 6-1-26
14246 5cbcd9be7af152186cd0c83c Rissotto 4 Quesos Galle 200 Gr Arroz GALLO 200.0 gr 7790070411754 123.55 6-2-1
```

14247 rows x 19 columns

Antes de emépezar el analisis, miramos que columnas tiene para ver con cuales quedarnos para el analisis.

```
In [15]: precios_claros.columns

Out[15]: Index(['_id/soid', 'nombre', 'categoria', 'marca', 'presentacion', 'id', 'precio', 'sucursal_id', 'comercio', 'bandera', 'banderaDescripcion', 'comercioRazonSocial', 'provincia', 'localidad', 'direccion', 'lat', 'lng', 'sucursalNombre', 'sucursalTipo'], dtype='object')
```

Seleccionamos las columnas con las que vamos a trabajar.

```
In [16]: precios_claros = precios_claros[['id', 'categoria', 'marca', 'presentacion', 'precio', 'sucursal_id', 'sucursal_nombre', 'sucursal_tipo']]
precios_claros
```

```
Out[16]: id categoría marca presentacion precio sucursal_id sucursalTipo banderaDescripcion
0 7798062547788 Agua VILLAVICENCIO 1.5 lt 70.00 10-1-112 Hipermercado Hipermercado Carrefour
1 7798062547788 Agua VILLAVICENCIO 1.5 lt 85.00 10-1-142 Supermercado Hipermercado Carrefour
2 7798062547788 Agua VILLAVICENCIO 1.5 lt 69.00 10-1-178 Supermercado Hipermercado Carrefour
3 7798062547788 Agua VILLAVICENCIO 1.5 lt 75.50 10-1-233 Hipermercado Hipermercado Carrefour
4 7798062547788 Agua VILLAVICENCIO 1.5 lt 73.00 10-1-29 Hipermercado Hipermercado Carrefour
```

```
... ..
14242 7790070411754 Arroz GALLO 200.0 gr 139.00 49-1-2 Supermercado Super Tuti1 y Super Tuti 3
14243 7790070411754 Arroz GALLO 200.0 gr 110.90 5-1-3 Hipermercado California Supermarkets
14244 7790070411754 Arroz GALLO 200.0 gr 139.00 50-1-2 Supermercado Super Tuti2 y Super Tuto
14245 7790070411754 Arroz GALLO 200.0 gr 123.55 6-1-26 Supermercado Supermercados Comodin
14246 7790070411754 Arroz GALLO 200.0 gr 123.55 6-2-1 Hipermercado Maxi Comodin
```

14247 rows x 8 columns

Como el tipo del trabajo es poder predecir a que tipo de cadena pertenece el producto y en nuestro dataset actual hay demasiadas cadenas, filtraremos el dataset para algunas cadenas de comercios, las que se consideran mas relevantes e influyentes en el mercado. !!!Como considero cada cadena

```
In [17]: cadenas = ["Disco", "Hipermercado Carrefour", "Supermercados DIA"]
precios_claros = precios_claros[precios_claros.banderaDescripcion.isin(cadenas)]
precios_claros
```

```
Out[17]: id categoría marca presentacion precio sucursal_id sucursalTipo banderaDescripcion
0 7798062547788 Agua VILLAVICENCIO 1.5 lt 70.00 10-1-112 Hipermercado Hipermercado Carrefour
1 7798062547788 Agua VILLAVICENCIO 1.5 lt 85.00 10-1-142 Supermercado Hipermercado Carrefour
2 7798062547788 Agua VILLAVICENCIO 1.5 lt 69.00 10-1-178 Supermercado Hipermercado Carrefour
3 7798062547788 Agua VILLAVICENCIO 1.5 lt 75.50 10-1-233 Hipermercado Hipermercado Carrefour
4 7798062547788 Agua VILLAVICENCIO 1.5 lt 73.00 10-1-29 Hipermercado Hipermercado Carrefour
```

```
... ..
14172 7790236001263 Pastas LA SALTENA 400.0 gr 79.99 15-1-1059 Supermercados DIA
14173 7790236001263 Pastas LA SALTENA 400.0 gr 75.00 15-1-800 Autoservicio Supermercados DIA
14199 7790236001263 Pastas LA SALTENA 400.0 gr 77.00 9-2-4 Supermercado Disco
14200 7790236001263 Pastas LA SALTENA 400.0 gr 75.00 9-2-583 Supermercado Disco
14203 7790070411754 Arroz GALLO 200.0 gr 136.00 10-1-49 Hipermercado Hipermercado Carrefour
```

2486 rows x 8 columns

Establecemos mejores nombres a nuestras columnas, en base a los cuales se realizara el listado de variables del trabajo.

```
In [18]: precios_claros = precios_claros.rename(columns = { 'id': 'producto_id', 'sucursal_id': 'tipo_cadena', 'banderaDescripcion': 'cadena', 'comercioRazonSocial': 'comercio' })
precios_claros
```

Como para el analisis necesitamos la presentación de la unidad separada de la medida, primero separaremos la columna para realizar un ultimo filtrado previo al preprocesamiento de los datos.

```
In [19]: precios_claros[['presentacion', 'unidad_presentacion']] = precios_claros['presentacion'].str.split(' ', expand=True)
precios_claros
```

Cambiamos la variable a tipo numerico y verificamos que se haya cambiado efectivamente

```
In [20]: precios_claros[['presentacion']] = pd.to_numeric(precios_claros.presentacion)
precios_claros

Out[20]: producto_id object
categoría object
marca object
presentacion float64
precio float64
sucursal_id object
tipo_cadena object
cadena object
unidad_presentacion object
dtype: object
```

2 - PREPROCESAMIENTO

Veo las columnas que hay

```
In [21]: precios_claros.columns

Out[21]: Index(['producto_id', 'categoria', 'marca', 'presentacion', 'precio', 'sucursal_id', 'tipo_cadena', 'cadena', 'unidad_presentacion'], dtype='object')
```

Las dimensiones:

```
In [22]: dimensiones = precios_claros.shape
print(f"El dataset cuenta con {dimensiones[0]} registros y {dimensiones[1]} columnas")

El dataset cuenta con 2486 registros y 9 columnas

Miramos la tipologia de los datos de las variables. De igual manera seran tratados mas adelante.
```

```
In [23]: precios_claros.dtypes

Out[23]: producto_id object
categoría object
marca object
presentacion float64
precio float64
sucursal_id object
tipo_cadena object
cadena object
unidad_presentacion object
dtype: object
```

- Outliers

Miramos cuantos duplicados tiene el dataset y los miramos para ver si hay que eliminarlos para evitar el overfitting.

```
In [24]: print(f"El dataset parece tener {len(precios_claros)} = {len(precios_claros[precios_claros.duplicated()])} filas")

El dataset parece tener 1731 filas duplicadas

Cuando miramos esos datos vemos que en realidad no son duplicados, es decir, las primeras columnas contienen el mismo dato porque hacen referencia al mismo producto, pero en realidad hace referencia a ese producto en distintas sucursales con distintos precios. Dado esto, no vamos a eliminarlos pero es precisamente ese el objetivo de analisis en este trabajo, mirar los productos y sus precios comparandolos en distintas cadenas y tipos de cadena.
```

```
In [25]: precios_claros[precios_claros.duplicated()]

Out[25]: producto_id categoría marca presentacion precio sucursal_id tipo_cadena cadena unidad_presentacion
68 7798062547788 Agua VILLAVICENCIO 1.5 70.0 10-1-112 Hipermercado Hipermercado Carrefour
70 7798062547788 Agua VILLAVICENCIO 1.5 85.0 10-1-142 Supermercado Hipermercado Carrefour
73 7798062547788 Agua VILLAVICENCIO 1.5 73.0 10-1-29 Hipermercado Hipermercado Carrefour
75 7798062547788 Agua VILLAVICENCIO 1.5 88.0 10-1-48 Hipermercado Hipermercado Carrefour
76 7798062547788 Agua VILLAVICENCIO 1.5 83.0 10-1-53 Hipermercado Hipermercado Carrefour
```

```
... ..
14076 7790236001263 Pastas LA SALTENA 400.0 72.0 10-1-48 Hipermercado Hipermercado Carrefour
14107 7790236001263 Pastas LA SALTENA 400.0 75.0 15-1-1803 Supermercado Supermercados DIA
14136 7790236001263 Pastas LA SALTENA 400.0 69.0 10-1-147 Supermercado Hipermercado Carrefour
14137 7790236001263 Pastas LA SALTENA 400.0 85.0 10-1-142 Hipermercado Hipermercado Carrefour
14143 7790236001263 Pastas LA SALTENA 400.0 72.0 10-1-48 Hipermercado Hipermercado Carrefour
```

755 rows x 9 columns

Igualmente podemos analizar que productos con su respectiva sucursal tienen repetidos los precios. Esto se da debido a que se tomaron los precios 3 veces en el mes. Como sacar el promedio podria ser alejado a la realidad por la existencia de outliers, nos quedaremos con el primer precio tomado de cada producto en el mes.

```
In [26]: precios_claros.drop_duplicates(subset = ['producto_id', 'sucursal_id'], keep = "first", inplace = True)

In [27]: print(f"La cantidad de registros sin duplicados es de {len(precios_claros)}")

La cantidad de registros sin duplicados es de 1729
```

```
In [28]: precios_claros.drop(columns = ['producto_id', 'sucursal_id'], inplace = True)

Verificamos su eliminacion
```

```
In [29]: precios_claros.columns

Out[29]: Index(['categoria', 'marca', 'presentacion', 'precio', 'tipo_cadena', 'cadena', 'unidad_presentacion'], dtype='object')
```

- Faltantes

Porcentaje de faltantes de variables por columna. Nuestro dataset no contiene faltantes, por lo que no hay que hacer ningun tipo de tratamiento especial.

```
In [30]: round(precios_claros.groupby(["unidad_presentacion", "categoria"]).presentacion.mean()).round(2)

Out[30]: categoría 0.0
marca 0.0
presentacion 0.0
precio 0.0
tipo_cadena 0.0
cadena 0.0
unidad_presentacion 0.0
dtype: float64
```

- Outliers

Antes de poder empezar es fundamental analizar la posible existencia de outliers en nuestros datos para poder tratarlos de manera adecuada.

```
In [31]: precios_claros.precio.describe().round(2)

Out[31]: count 1729.00
mean 84.51
std 45.82
min 9.59
25% 54.75
50% 71.00
75% 92.00
max 224.00
Name: precio, dtype: float64
```

Mediante el boxplot podemos observar que el cuadrado azul es donde hay mayor concentracion y es el rango intercuantil. Lo que esta fuera de los bigotes son los que parecen ser outliers o podrian ser malas imputaciones.

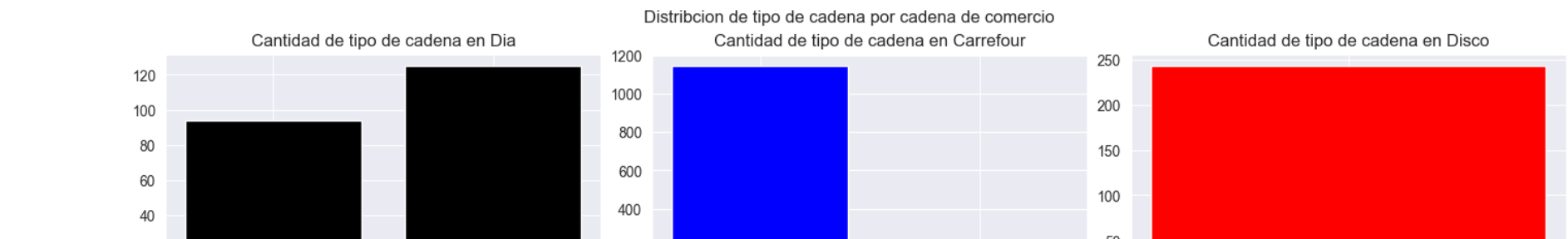
```
In [32]: fig, ax = plt.subplots(1,2, figsize=(14,5))

ax[0].set_title("Boxplot de precio")
sns.boxplot(data = precios_claros, x = "precio", ax = ax[0])

ax[1].set_title("Boxplot de presentacion")
sns.boxplot(data = precios_claros, x = "presentacion", ax = ax[1])
```

```
Out[32]: <AxesSubplot: title='center': 'Boxplot de presentacion', xlabel='presentacion'>

Boxplot de precio
```

Se puede ver que Supermercados dia tiene autoservicios y en su mayoría supermercados, Carrefour en su mayoría Hi permercados pero tambien cuenta con supermercados, mientras que en Disco son todos supermercados

4 - ESCALADO

Aplicamos el escalado a las variables con el objetivo de estandarizarlas, es decir, de ponerlas en el mismo rango de valores para que sean comparables entre si

In [47]: `precios_claros.dtypes`

Out[47]: `categoria object
marca object
presentacion float64
precio float64
tipo_cadena object
cadena object
unidad_presentacion object
dtype: object`

In [48]: `numericas = precios_claros.select_dtypes(exclude=['object']).columns`

`print("Hacer un escalado a la variable:")
for variable in numericas:
 print(f"t-{variable}")`

Hacer un escalado a las variables:
- precio

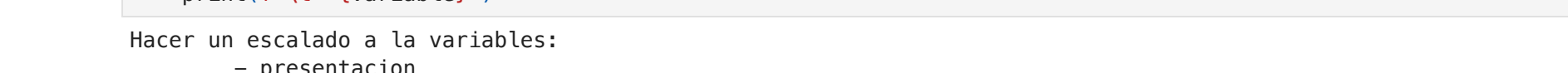
Variable presentacion

Como se ve en los graficos, esta variable no tiene una distribucion normal, esto se da porque los productos estan en diferentes unidades de medida. Por ello, previo a escalar la variable, podemos pasar los cc a lt y los gr a kg.

In [49]: `variable = 'presentacion'`

`fig, ax = plt.subplots(1,2,figsize=(14,5))
ax[0].set_title("Distribucion de (variable)")
sns.histplot(data = precios_claros, x = f"(variable)", kde=True, ax = ax[0])
ax[1].set_title("Boxplot de(variable)")
sns.boxplot(data = precios_claros, x = f"(variable)", ax=ax[1])`

Out[49]: `<AxesSubplot: title='(center): 'Boxplot de(presentacion)', xlabel='presentacion'>`



In [50]: `precios_claros.presentation.describe().round(2)`

Out[50]:

count	1729.00
mean	319.85
std	238.63
min	1.00
25%	1.50
50%	500.00
75%	500.00
max	750.00

Name: presentacion, dtype: float64

Como 1 litro equivale a 1000 centilitros cúbicos y 1 kilo gramo equivale a 1000 gramos, dividimos la presentacion de aquellos productos que estan medidos en cc y gr por 1000 y cambiamos su unidad de medida.

In [51]: `precios_claros.loc[precios_claros['unidad_presentacion'] == "cc", 'presentacion'] = precios_claros.presentation.precios_claros['unidad_presentacion'] = precios_claros['unidad_presentacion'].str.replace("cc", "lt", regex = F`

In [52]: `precios_claros.loc[precios_claros['unidad_presentacion'] == "gr", 'presentacion'] = precios_claros.presentation.precios_claros['unidad_presentacion'] = precios_claros['unidad_presentacion'].str.replace("gr", "kg", regex = F`

Vemos que ahora las unidades de medida son litros y kg.

In [53]: `precios_claros.unidad_presentacion.value_counts(dropna = False)`

Out[53]:

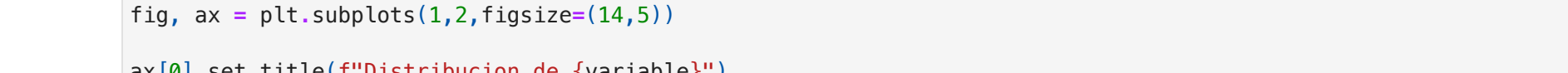
kg	1203
lt	526

Name: unidad_presentacion, dtype: int64

In [54]: `variable = 'presentacion'`

`fig, ax = plt.subplots(1,2,figsize=(14,5))
ax[0].set_title("Distribucion de (variable)")
sns.histplot(data = precios_claros, x = f"(variable)", kde=True, ax = ax[0])
ax[1].set_title("Boxplot de(variable)")
sns.boxplot(data = precios_claros, x = f"(variable)", ax=ax[1])`

Out[54]: `<AxesSubplot: title='(center): 'Boxplot de(presentacion)', xlabel='presentacion'>`



En esta instancia se puede ver nuevamente que los que parece ser malas imputaciones es porque el producto es de mayor presentacion.

In [55]: `precios_claros[precios_claros.presentation > 3]`

Out[55]:

categoria	marca	presentacion	precio	tipo_cadena	cadena	unidad_presentacion	
886	Agua	KIN	6.00	125.00	Hipermarcado	Hipermarcado Carrefour	lt
886	Agua	KIN	6.00	121.49	Autoservicio	Supermercados DIA	lt
887	Agua	KIN	6.00	123.39	Autoservicio	Supermercados DIA	lt
888	Agua	KIN	6.00	93.00	Supermercado	Supermercados DIA	lt
891	Agua	KIN	6.00	125.00	Supermercado	Disco	lt
...
1345	Agua	VILLA DEL SUR	6.25	175.00	Hipermarcado	Supermercados Carrefour	lt
1384	Agua	VILLA DEL SUR	6.25	197.99	Supermercado	Supermercados DIA	lt
1385	Agua	VILLA DEL SUR	6.25	198.59	Autoservicio	Supermercados DIA	lt
1386	Agua	VILLA DEL SUR	6.25	175.00	Autoservicio	Supermercados DIA	lt
1423	Agua	VILLA DEL SUR	6.25	187.00	Supermercado	Disco	lt

84 rows x 7 columns

Ahora si le podemos aplicar el escalador.

In [56]: `from sklearn.preprocessing import RobustScaler
escalador = RobustScaler()
variable = 'presentacion'
precios_claros[variable] = escalador.fit_transform(precios_claros[variable])
print("Como los datos en (variable) tienen outliers, imputamos el RobustScaler, que escala los valores restando la mediana y dividiendolo por el rango intercuartil")`

Como los datos en presentacion tienen outliers, imputamos el RobustScaler, que escala los valores restando la mediana y dividiendolo por el rango intercuartil

Out[56]:

categoria	marca	presentacion	precio	tipo_cadena	cadena	unidad_presentacion	
0	Agua	VILLAVICENCIO	2.0	70.0	Hipermarcado	Hipermarcado Carrefour	lt
1	Agua	VILLAVICENCIO	2.0	85.0	Supermercado	Hipermarcado Carrefour	lt
2	Agua	VILLAVICENCIO	2.0	69.0	Supermercado	Hipermarcado Carrefour	lt
3	Agua	VILLAVICENCIO	2.0	75.5	Hipermarcado	Hipermarcado Carrefour	lt
4	Agua	VILLAVICENCIO	2.0	73.0	Hipermarcado	Hipermarcado Carrefour	lt

Variable precio

A igual que en el caso anterior, aplicaremos el robust scaler.

In [57]: `variable = 'precio'`

`fig, ax = plt.subplots(1,2,figsize=(14,5))
ax[0].set_title("Distribucion de (variable)")
sns.histplot(data = precios_claros, x = f"(variable)", kde=True, ax = ax[0])
ax[1].set_title("Boxplot de(variable)")
sns.boxplot(data = precios_claros, x = f"(variable)", ax=ax[1])`

Out[57]: `<AxesSubplot: title='(center): 'Boxplot de(precio)', xlabel='precio'>`



In [58]: `from sklearn.preprocessing import RobustScaler
escalador = RobustScaler()
variable = 'precio'
precios_claros[variable] = escalador.fit_transform(precios_claros[variable])
print("Como los datos en (variable) tienen outliers, imputamos el RobustScaler, que escala los valores restando la mediana y dividiendolo por el rango intercuartil")`

Como los datos en precio tienen outliers, imputamos el RobustScaler, que escala los valores restando la mediana y dividiendolo por el rango intercuartil

Out[58]:

categoria	marca	presentacion	precio	tipo_cadena	cadena	unidad_presentacion	
0	Agua	VILLAVICENCIO	2.0	-0.026846	Hipermarcado	Hipermarcado Carrefour	lt
1	Agua	VILLAVICENCIO	2.0	0.375839	Supermercado	Hipermarcado Carrefour	lt
2	Agua	VILLAVICENCIO	2.0	-0.053691	Supermercado	Hipermarcado Carrefour	lt
3	Agua	VILLAVICENCIO	2.0	0.120805	Hipermarcado	Hipermarcado Carrefour	lt
4	Agua	VILLAVICENCIO	2.0	0.053691	Hipermarcado	Hipermarcado Carrefour	lt

5 - ENCODING

In [59]: `categoricas = precios_claros.select_dtypes(exclude=['float', 'integer']).columns
print("Hacer un encoding a las variables:")
for variable in categoricas:
 print(f"t-{variable}")`

Hacer un encoding a las variables:
- categoria
- marca
- tipo_cadena
- cadena
- unidad_presentacion

Variable categoria

Lo que hace este bloque de codigo expandir la columna cuenta a columnas binarias para cada categoria dentro de la variable, asignandolas 1 y 0 dependiendo de si el producto coincide con la categoria o no. Mediante este codigo pudimos automatizar la busqueda de las diferentes categorias y la asignacion del nombre de la columna. Se repite el codigo para todas las variables categoricas menos al target.

In [60]: `categoricas = precios_claros.categoria.unique()
print("Categorias unicas: ")
for categoria in categoricas:
 print(f"t-{categoria}")`

Categorias unicas:
- Agua
- Arroz
- Harinas
- Leche
- Pastas

In [61]: `from sklearn.preprocessing import OneHotEncoder
ohe = OneHotEncoder(sparse=False)
df = precios_claros
variable = 'categoria'
ohe.fit(df[[variable]])
display(ohe.categories_)
display(ohe.get_feature_names_out())
variable_encoded = ohe.transform(df[[variable]])
df[ohe.get_feature_names_out()] = variable_encoded
df.drop(columns = variable, inplace = True)
print("Como (variable) es una variable multicategorica, creo una columna binaria para cada categoria")
df.head()`

[array(['Agua', 'Arroz', 'Harinas', 'Leche', 'Pastas'], dtype=object)]
array(['categoria_Agua', 'categoria_Arroz', 'categoria_Harinas', 'categoria_Leche', 'categoria_Pastas'], dtype=object)
Como categoria es una variable multicategorica, creo una columna binaria para cada categoria

Out[61]:

	marca	presentacion	precio	tipo_cadena	cadena	unidad_presentacion	categoria_Agua	categoria_Arroz	categoria_Harinas	categoria_Leche	categoria_Pastas
0	VILLAVICENCIO	2.0	-0.026846	Hipermarcado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0
1	VILLAVICENCIO	2.0	0.375839	Supermercado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0
2	VILLAVICENCIO	2.0	-0.053691	Supermercado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0
3	VILLAVICENCIO	2.0	0.120805	Hipermarcado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0
4	VILLAVICENCIO	2.0	0.053691	Hipermarcado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0

Variable marca

In [62]: `marcas = precios_claros.marca.unique()
print("Marcas unicas: ")
for marca in marcas:
 print(f"t-{marca}")`

Marcas unicas:
- VILLA DEL SUR
- SER
- BONAQUA
- ECO DE LOS ANDES
- KIN
- NESTLE
- VILLA DEL SUR
- GLACIAR
- GALLO
- MOLINOS ALA
- LUCCHETTI
- MAIZENA
- PUREZA
- CAÑUELAS
- FAVORITA
- QUAKER
- ARCOR
- LA SERENISIMA
- ARMONIA
- GIACOMO
- DON VICENTE
- MATARAZZO
- KNORR
- LA SALTEÑA

In [63]: `from sklearn.preprocessing import OneHotEncoder
ohe = OneHotEncoder(sparse=False)
df = precios_claros
variable = 'marca'
ohe.fit(df[[variable]])
display(ohe.categories_)
display(ohe.get_feature_names_out())
variable_encoded = ohe.transform(df[[variable]])
df[ohe.get_feature_names_out()] = variable_encoded
df.drop(columns = variable, inplace = True)
print("Como (variable) es una variable multicategorica, creo una columna binaria para cada categoria")
df.head()`

[array(['ARCOR', 'ARMONIA', 'BONAQUA', 'CAÑUELAS', 'DON VICENTE', 'ECO DE LOS ANDES', 'FAVORITA', 'GALLO', 'GIACOMO', 'GLACIAR', 'KIN', 'KNORR', 'LA SALTEÑA', 'LA SERENISIMA', 'LUCCHETTI', 'MAIZENA', 'MATARAZZO', 'MOLINOS ALA', 'NESTLE', 'PUREZA', 'QUAKER', 'SER', 'VILLA DEL SUR', 'VILLAVICENCIO'], dtype=object)]
array(['marca_ARCOR', 'marca_ARMONIA', 'marca_BONAQUA', 'marca_CAÑUELAS', 'marca_DON VICENTE', 'marca_ECO DE LOS ANDES', 'marca_FAVORITA', 'marca_GALLO', 'marca_GIACOMO', 'marca_GLACIAR', 'marca_KIN', 'marca_KNORR', 'marca_LA SALTEÑA', 'marca_LA SERENISIMA', 'marca_LUCCHETTI', 'marca_MAIZENA', 'marca_MATARAZZO', 'marca_MOLINOS ALA', 'marca_NESTLE', 'marca_PUREZA', 'marca_QUAKER', 'marca_SER', 'marca_VILLA DEL SUR', 'marca_VILLAVICENCIO'], dtype=object)
Como marca es una variable multicategorica, creo una columna binaria para cada categoria

Out[63]:

	presentacion	precio	tipo_cadena	cadena	unidad_presentacion	categoria_Agua	categoria_Arroz	categoria_Harinas	categoria_Leche	categoria_Pastas
0	2.0	-0.026846	Hipermarcado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0
1	2.0	0.375839	Supermercado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0
2	2.0	-0.053691	Supermercado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0
3	2.0	0.120805	Hipermarcado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0
4	2.0	0.053691	Hipermarcado	Hipermarcado Carrefour	lt	1.0	0.0	0.0	0.0	0.0

5 rows x 34 columns

Variable tipo_cadena

In [64]: `tipos_cadenas = precios_claros.tipo_cadena.unique()
print("Tipos de cadena unicas: ")
for tipo_cadena in tipos_cadenas:
 print(f"t-{tipo_cadena}")`

Tipos de cadena unicas:
- Hipermarcado
- Supermercado
- Autoservicio

In [65]: `from sklearn.preprocessing import OneHotEncoder
ohe = OneHotEncoder(sparse=False)
df = precios_claros
variable = 'tipo_cadena'
ohe.fit(df[[variable]])
display(ohe.categories_)
display(ohe.get_feature_names_out())
variable_encoded = ohe.transform(df[[variable]])
df[ohe.get_feature_names_out()] = variable_encoded
df.drop(columns = variable, inplace = True)
print("Como (variable) es una variable multicategorica, creo una columna binaria para cada categoria")
df.head()`

[array(['Autoservicio', 'Hipermarcado', 'Supermercado'], dtype=object)]
array(['tipo_cadena_Autoservicio', 'tipo_cadena_Hipermarcado', 'tipo_cadena_Supermercado'], dtype=object)
Como tipo_cadena es una variable multicategorica, creo una columna binaria para cada categoria

Out[65]:

	presentacion	precio	cadena	unidad_presentacion	categoria_Agua	categoria_Arroz	categoria_Harinas	categoria_Leche	categoria_Pastas
0	2.0	-0.026846	Hipermarcado	Carrefour	lt	1.0	0.0	0.0	0.0
1	2.0	0.375839	Supermercado	Carrefour	lt	1.0	0.0	0.0	0.0
2	2.0	-0.053691	Supermercado	Carrefour	lt	1.0	0.0	0.0	0.0
3	2.0	0.120805	Hipermarcado	Carrefour	lt	1.0	0.0	0.0	0.0
4	2.0	0.053691	Hipermarcado	Carrefour	lt	1.0	0.0	0.0	0.0

5 rows x 36 columns

Variable unidad_presentacion

En este caso, como la variable ya es binaria, no es necesario expandirla, podemos hacer el encoding en la misma columna, asignandole 0 a los lt y 1 a los kg.

In [66]: `unidades_presentacion = precios_claros.unidad_presentacion.unique()
print("Unidades de presentacion unicas: ")
for unidad_presentacion in unidades_presentacion:
 print(f"t-{unidad_presentacion}")`

Unidades de presentacion unicas:
- lt
- kg

In [67]: `hot_encoder = OneHotEncoder(sparse = False, drop = 'if_binary', categories=[['lt', 'kg']]) #lt va a ser cero y kg va a ser uno
df = precios_claros
variable = 'unidad_presentacion'
df[variable] = hot_encoder.fit_transform(df[[variable]])
print("Como tengo una columna con dos categorias, hago el one-hot pero transformando la columna en binaria")
df.head()`

Out[67]:

	presentacion	precio	cadena	unidad_presentacion	categoria_Agua	categoria_Arroz	categoria_Harinas	categoria_Leche	categoria_Pastas
0	2.0	-0.026846	Hipermarcado	Carrefour	0.0	1.0	0.0	0.0	0.0
1	2.0	0.375839	Hipermarcado	Carrefour	0.0	1.0	0.0	0.0	0.0
2	2.0	-0.053691	Hipermarcado	Carrefour	0.0	1.0	0.0	0.0	0.0
3	2.0	0.120805	Hipermarcado	Carrefour	0.0	1.0	0.0	0.0	0.0
4	2.0	0.053691	Hipermarcado	Carrefour	0.0	1.0	0.0	0.0	0.0

5 rows x 36 columns

Target cadena

Ahora aplicamos el Label Encoder a tipo_cadena, que es nuestro target. Siendo 0 Disco, 1 Supermercado Carrefour y 2 Supermercados Dia.

In [68]: `cadenas = precios_claros.cadena.unique()
print("Unidades de presentacion unicas: ")
for cadena in cadenas:
 print(f"t-{cadena}")`

Unidades de presentacion unicas:
- Hipermarcado Carrefour
- Supermercados DIA
- Disco

In [69]: `from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df = precios_claros
variable = 'cadena'
df[variable] = le.fit_transform(df[[variable]])
df.head()`

Out[69]:

	presentacion	precio	cadena	unidad_presentacion	categoria_Agua	categoria_Arroz	categoria_Harinas	categoria_Leche	categoria_Pastas
0	2.0	-0.026846	1	0.0	1.0	0.0	0.0	0.0	0.0
1	2.0	0.375839	1	0.0	1.0	0.0	0.0	0.0	0.0
2	2.0	-0.053691	1	0.0	1.0	0.0	0.0	0.0	0.0
3	2.0	0.120805	1	0.0	1.0	0.0	0.0	0.0	0.0
4	2.0	0.053691	1	0.0	1.0	0.0	0.0	0.0	0.0

5 rows x 36 columns

6 - MODELO

Hacemos un primer modelo de regresion logistica mediante cross validation, ajustando el parametro para obtener un mejor accuracy. Lo que hace es:

- divide el conjunto de datos en un conjunto de entrenamiento y un conjunto de prueba.
- ajusta el modelo, hace predicciones y calcula una puntuación
- repite el proceso 10 veces en total
- genera el promedio de las 10 puntuaciones

In [70]: `from sklearn.model_selection import cross_validation
from sklearn.linear_model import LogisticRegression
X = precios_claros.drop(columns = 'cadena')
y = precios_claros['cadena']
model = LogisticRegression(max_iter=1000)
cv_results = cross_validate(model, X, y, cv=10)
base_model_score = cv_results['test_score'].mean()
print(f"En este primer modelo se obtuvo un accuracy de {round(base_model_score, 4)}")
En este primer modelo se obtuvo un accuracy de 0.8005`

Ahora aplicamos el train_test_split y le indicamos el tamaño de muestra que queremos para que se entrene el modelo y se prube, en este caso el tamaño de prueba es 0.3, por lo tanto el de entrenamiento es de 0.7

In [71]: `from sklearn.model_selection import LogisticRegression
from sklearn.linear_model import LogisticRegression
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 10)
model = LogisticRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
train_test_score = model.score(X_test, y_test)
print(f"Realizamos un train y test para mejorar nuestro modelo y obtuvimos un score de: {round(train_test_score, 4)}")
Realizamos un train y test para mejorar nuestro modelo y obtuvimos un score de: 0.8343, pudiendo ver que mejoró nuestro modelo en 0.0338`

/Users/usuario/PycharmProjects/pythonProjects/P00/lib/python3.9/site-packages/sklearn/linear_model/_logistic.py:444: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_ = 1
n_iter_ = 1

Repetimos el train_test_split anterior, pero en esta instancia estratificamos nuestro target, es decir, buscamos que cada conjunto contenga aproximadamente el mismo porcentaje de muestras de cada clase objetivo que el conjunto completo.

In [72]: `X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, stratify = y, random_state = 10)
model = LogisticRegression()
model.fit(X_train, y_train)
y_pred = model.predict(X_test)
train_test_score = model.score(X_test, y_test)
print(f"Realizamos nuestra variable target obteniendo un accuracy {round(stratify_score, 4)}, mejorado un 0.0338
Realizamos nuestra variable target obteniendo un accuracy 0.8536, mejorado un 0.0193 con respecto al anterior en 0.0198`

/Users/usuario/PycharmProjects/pythonProjects/P00/lib/python3.9/site-packages/sklearn/linear_model/_logistic.py:444: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.
Increase the number of iterations (max_iter) or scale the data as shown in:
https://scikit-learn.org/stable/modules/preprocessing.html
Please also refer to the documentation for alternative solver options:
https://scikit-learn.org/stable/modules/linear_model.html#logistic-regression
n_iter_ = 1
n_iter_ = 1

Utilizamos la funcion para establecer cuales son las variables con mayor influencia en nuestro modelo. A estas variables les llamamos variables "fuertes" y mediante un umbral seleccionamos aquellas que nos quedaremos en nuestro modelo.

In [73]: `from sklearn.model_selection import cross_val_score
from sklearn.inspection import permutation_importance
print("Antes de la permutacion, el score del modelo es: {round(stratify_score, 4)}")
model = LogisticRegression().fit(X, y)
permutation_score = permutation_importance(model, X, y, n_repeats=100)
importance_df = pd.DataFrame(np.vstack((X.columns, permutation_score.importances.mean())).T, columns = ['feature', '`

