

UNIVERSIDAD COMPLUTENSE DE MADRID
FACULTAD DE ESTUDIOS ESTADÍSTICOS



TAREA EVALUABLE

MÓDULO: MACHINE LEARNING - CLASIFICACIÓN BINARIA

Jorge González Perea 51553561G

Máster en Big Data, Data Science & Inteligencia Artificial

Curso académico 2024-2025

Índice

1. Introducción.	2
2. Sección 2 - Análisis y depuración de los datos.	3
2.1. Variables numéricas.	3
2.2. Variables categóricas.	5
3. Modelo de árbol de decisiones.	7
3.1. Resultados del modelo de árbol.	10
4. Modelo de <i>Random Forest</i>.	13
5. Modelo de <i>XGBoost</i>.	17
6. Comparación de modelos y conclusiones.	20
6.1. Comparación.	20
6.2. Conclusiones.	21
7. Anexo.	22
7.1. Librerías.	22
7.2. Funciones.	22
7.3. Imágenes.	23
7.4. Códigos para generar imágenes.	23

1. Introducción.

En esta práctica se va a tratar una base de datos con 4340 registros de información perteneciente a distintos vehículos de segunda mano. Las variables disponibles son las siguientes:

- **Price**: variable numérica. Precio del vehículo.
- **Levy**: variable numérica. Cantidad de impuestos a pagar.
- **Manufacturer**: variable categórica. Empresa fabricante del vehículo.
- **Prod. year**: variable numérica. Año de fabricación.
- **Category**: variable categórica. Tipo de vehículo.
- **Leather interior**: Variable categórica. Indica si el interior es de cuero o no.
- **Fuel type**: variable categórica. Indica el tipo de combustible del vehículo.
- **Engine volume**: variable categórica. Volumen de desplazamiento de los cilindros del motor. También indica si el motor tiene acoplado un turbo.
- **Mileage**: variable categórica. Kilometraje del vehículo.
- **Cylinders**: variable numérica. Número de cilindros del bloque.
- **Gear box Type**: variable categórica. Tipo de transmisión del vehículo.
- **Drive wheels**: variable categórica. Indica las ruedas que traccionan.
- **Wheel**: variable categórica. Indica el lado en el que el volante está situado.
- **Color**: variable categórica. Color del vehículo.
- **Airbags**: variable numérica. Número de airbags en el vehículo.

El objetivo es emplear las 14 variables restantes para poder predecir la variable objetivo **Color**, de forma que se pueda saber si conviene pintar el coche de color blanco o no. Para ello se va a llevar a cabo un análisis y una limpieza de la base de datos. Posteriormente se van a desarrollar distintos modelos de árboles de regresión que ayude con la toma de decisiones para la variable objetivo: un árbol simple, un *Random Forest* y *XGBoost*. Por último, estos modelos serán comparados con el fin de justificar cuál de ellos es mejor para esta BBDD. Todo ello será llevado a cabo con *Python* y con las librerías de funciones que se pueden encontrar en la sección 7 (anexo).

La metodología para el desarrollo de los modelos es la siguiente en los tres casos:

- Primero se crea un clasificador base con parámetros obtenidos de forma arbitraria y que minimicen el sobreajuste sin pasarse. Este estimador ya se puede usar como modelo.
- El segundo paso consiste en hacer una búsqueda de los mejores parámetros mediante la función *GridSearchCV* sobre el estimador base del paso anterior. El mejor modelo se obtiene con el método *grid_search.best_estimator*.
- Otro modelo prometedor se puede encontrar haciendo una búsqueda manual mediante *boxplots* en los resultados de *GridSearchCV*.
- La cuarta y última manera empleada para buscar parámetros óptimos es visualizado la dependencia de la precisión de los modelos en función de estos. Para ello se crean gráficas manualmente y se eligen los mejores parámetros.

2. Sección 2 - Análisis y depuración de los datos.

En primer lugar, para poder trabajar con la base de datos de forma más eficiente es imprescindible familiarizarse con ella y llevar a cabo un análisis descriptivo de las variables. Tras este primer análisis se procede con una limpieza de los datos, ya que puede haber registros *missing* o valores anómalos. Por último se lleva a cabo la transformación de variables (*feature engineering*) de forma que la tabla de datos se ha simplificado lo máximo posible mediante codificaciones y eliminación de información redundante.

Con el siguiente código se lee el archivo de datos y se muestran los formatos de las distintas variables:

```
1 datos = pd.read_excel("datos_tarea25.xlsx")
2 print(datos.dtypes)
```

Listing 1: Lectura de datos y análisis inicial.

El método `datos.dtypes` proporciona los formatos de las variables del *DataFrame*. En este caso se observan algunas incongruencias:

Price	int64
Levy	object
Manufacturer	object
Prod. year	int64
Category	object
Leather interior	object
Fuel type	object
Engine volume	object
Mileage	object
Cylinders	int64
Gear box type	object
Drive wheels	object
Wheel	object
Color	object
Airbags	int64

Figura 1: Formato de las variables

Tal y como se puede ver, hay variables categóricas como **Levy**, **Engine volume** o **Mileage** que deberían ser numéricas. Por otro lado, otras son numéricas a pesar de poder ser codificadas tal y como se verá más adelante.

Para trabajar con esta base de datos se van a analizar todas las variables por separado, de forma que se puedan controlar con seguridad los formatos y los valores perdidos o anómalos de cada una.

2.1. Variables numéricas.

En esta sección se van a analizar las variables numéricas o aquellas que deberían ser numéricas. En este caso son **Price**, **Levy**, **Prod. year**, **Engine volume**, **Mileage**, **Cylinders**, **Airbags**. Algunas están en formato numérico y otras en formato de texto tal y como se indica en la imagen 1, por lo que será conveniente convertir todas a formato numérico antes de comenzar con el verdadero análisis.

En el caso de la variable **Price** no se observan datos faltantes o perdidos, pero sí un registro con un valor anormalmente alto (se puede observar en la imagen 24 del anexo)

```
1 var = 'Price'
2 print(len(datos[var].unique()))
3 print(datos[var].isna().sum())
4 datos[var].hist(bins = 30, color = 'mediumorchid')
5 plt.ylim(top = 20)
6 plt.show()
```

Listing 2: Tratamiento de la variable **Price**.

Por el contrario, la variable **Levy** contiene valores faltantes identificados con el carácter ". Por ello, se cambia el formato de la columna a número y se sustituyen los registros faltantes por *NaN*:

```

1  var = 'Levy'
2
3  print(len(datos[var].unique()))
4  print(datos[var].isna().sum())
5  print(datos[var].unique())
6
7  datos[var] = datos[var].replace("-", np.nan)
8  num_missing = datos[var].isna().sum()
9  levy = [float(dato) for dato in datos[var].tolist()]
10 datos[var] = levy
11 porcentaje = round(100*num_missing/len(datos[var].tolist()), 2)
12 print(f"\nLa variable Levy tiene {num_missing} valores perdidos (un {porcentaje}%)")
13
14 datos[var].hist(bins = 30, color = 'mediumorchid')
15 plt.ylim(top = 50)
16 plt.show()

```

Listing 3: Tratamiento de la variable **Levy**.

La siguiente variable a analizar es **Prod. year**. Esta columna se encuentra en el formato correcto y sólo contiene un *outlier* que no supone un valor perdido/anómalo.

```

1  var = 'Prod. year'
2
3  print(len(datos[var].unique()))
4  print(datos[var].isna().sum())
5  datos[var].hist(bins = 30, color = 'mediumorchid')
6  plt.ylim(top = 20)
7  plt.show()

```

Listing 4: Tratamiento de la variable **Prod. year**.

En la columna **Engine volume** se tienen 58 valores únicos que se pueden pasar a formato numérico. Para ello es necesario separar las cadenas de aquellos vehículos cuyo motor contiene un turbocompresor. Esta información se recoge en una nueva variable categórica y binaria llamada **Turbo**, que contiene un 1 o un 0 en el caso de que el vehículo registrado contenga turbo o no respectivamente.

```

1  var = 'Engine volume'
2  print(len(datos[var].unique()))
3  print(datos[var].unique())
4  engine_volume = []
5  turbo = []
6  for dato in datos[var]:
7      cadenas = dato.split()
8      engine_volume.append(float(cadenas[0]))
9      if len(cadenas) > 1:
10         turbo.append(str(1))
11     else:
12         turbo.append(str(0))
13 datos[var] = engine_volume
14 datos["Turbo"] = turbo
15
16 datos[var].hist(color = 'mediumorchid')

```

Listing 5: Tratamiento de la variable **Engine volume**.

La variable **Mileage** contiene la unidad de medida en cada registro (la cadena "*km*"), por lo que se elimina la unidad y se transforman los datos a formato numérico. Al generar el histograma de los datos se observa un valor anómalo extremadamente alto (con un valor de más de 10^9 *km*, siendo la vida media de un vehículo unos 3×10^5 *km*). Por tanto, es necesario sustituir este valor por *NaN* para su posterior imputación aleatoria.

Por último, las variables **Cylinders** y **Airbags** no presentan valores anómalos o perdidos y están en el formato correcto, por lo que no es necesario modificar nada.

El siguiente y último paso para el tratamiento de las variables numéricas consiste en imputar los datos faltantes. Para ello se hace uso de la función *ImputacionCuant* de la librería *FuncionesMineria.py* (el código

de la función se puede encontrar en el anexo). El método de imputación consiste en la selección de valores aleatorios en función de la distribución de la variable. Las únicas columnas que contienen valores nulos son **Levy** (644) y **Mileage** (1):

```

1  for col in list(datos.columns):
2      print(f"{col}: {datos[col].isna().sum()}")
3  print()
4
5  columnas = ["Levy", "Mileage"]
6  for var in columnas:
7      datos[var] = ImputacionCuant(datos[var], 'aleatorio')
8
9  for col in columnas:
10     print(f"{col}: {datos[col].isna().sum()}")

```

Listing 6: Imputación de datos faltantes.

2.2. Variables categóricas.

En esta sección se van a analizar las columnas no numéricas de la base de datos para comprobar si contienen valores no válidos o están en formatos no correspondientes. Con el siguiente código se pueden comprobar los valores únicos de cada una:

```

1  numericas = datos.select_dtypes(include = ['int', 'int32', 'int64', 'float', 'float32',
2      'float64']).columns.tolist()
3  categoricas = [col for col in datos.columns.tolist() if col not in numericas]
4
5  for var in categoricas:
6      uv = len(datos[var].unique())
7      print(f"La variable {var} tiene {uv} valores únicos:")
8      print(datos[var].unique())
9      print()

```

Listing 7: Variables categóricas.

Estas variables contienen una baja cantidad de valores únicos, por lo que podrían ser codificadas. Sin embargo, el desarrollo de los modelos de árboles hace que no sea necesario ese proceso en el tratamiento de los datos, ya que a la hora de clasificar los datos mediante ramas, estas actúan como codificación. Sin embargo, las variables **Leather interior** y **Color** sí se han modificado para una mayor claridad en la clasificación de las mismas:

- La variable **Leather interior** se puede codificar para que los registros contengan 1 en lugar de *Yes* y 0 en lugar de *No*.
- La variable **Gear box type** se puede codificar para que los registros contengan 1 en lugar de *Automatic* y 0 en lugar de *Tiptronic*. También se ha cambiado el nombre de la columna a **Automatic gear box**.
- La variable **Wheel** se puede codificar para que los registros contengan 1 en lugar de *Leftwheel* y 0 en lugar de *Right – handdrive*. También se ha cambiado el nombre de la columna a **Left wheel**.
- La variable **Color** se puede codificar para que los registros contengan 1 en lugar de *White* y 0 en lugar de *Black*.

El código es el siguiente:

```

1  two_value_variables = ['Leather interior', 'Gear box type', 'Wheel', 'Color']
2  ones = ['Yes', 'Automatic', 'Left wheel', 'White']
3
4  for var in two_value_variables:
5      encoded_data = []
6      for dato in datos[var]:
7          if dato in ones:
8              encoded_data.append('1')
9          else:
10             encoded_data.append('0')
11     datos[var] = encoded_data

```

```

12  datos = datos.rename(columns={'Gear box type': 'Automatic gear box',
13                                'Wheel': 'Left wheel'})
14  print(datos.dtypes)
15  datos.to_excel("datos_tarea25_clean.xlsx")

```

Listing 8: Codificación de variables categóricas.

El resultado es un *DataFrame* con los datos limpios y los siguientes formatos:

Price	int64
Levy	float64
Manufacturer	object
Prod. year	int64
Category	object
Leather interior	object
Fuel type	object
Engine volume	float64
Mileage	float64
Cylinders	int64
Automatic gear box	object
Drive wheels	object
Left wheel	object
Color	object
Airbags	int64
Turbo	object

Figura 2: Formato de las variables tras la depuración.

3. Modelo de árbol de decisiones.

El primer modelo a crear es un árbol de decisiones que permita clasificar y predecir la variable objetivo **Color**. En el siguiente listado de código se lee el nuevo archivo de datos con la información depurada y se muestra el *DataFrame*. También se imprime por pantalla la frecuencia absoluta de cada valor de la variable objetivo.

```
1 df = pd.read_excel('datos_tarea25_clean.xlsx')
2 display(df.head())
3 print(f'\nLa frecuencia de cada clase es: \n{df.Color.value_counts()}')
```

Listing 9: Lectura de datos y análisis de la variable objetivo.

A continuación es necesario separar la variable objetivo de las predictoras. Para ello se separan los datos en una muestra de entrenamiento y una de prueba siguiendo el criterio 80 – 20 y se comprueba que ambas submuestras siguen una distribución parecida.

```
1 variables = df.columns.tolist()
2 objetivo = ['Color']
3 predictoras = [col for col in variables if col not in objetivo]
4 X = df[predictoras]
5 y = df[objetivo]
```

Listing 10: Separación de las variables y de la muestra.

El primer modelo de árbol simple es el siguiente. Para crearlo, se han escogido los valores de los parámetros de forma que se minimiza el sobreajuste (la diferencia entre las predicciones de *train* y *test* se hace mínima) sin que conlleve pérdidas en la eficacia de las predicciones.

```
1 arbol1 = DecisionTreeClassifier(min_samples_split = 200,
2 criterion = 'entropy', max_depth = 25)
3 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
4 random_state=42)
5 print(f'Frecuencia en train es: \n{y_train.value_counts(normalize=True)}')
6 print(f'\nFrecuencia en test es: \n{y_test.value_counts(normalize=True)}')
7 arbol1.fit(X_train, y_train)
8 y_train_pred = arbol1.predict(X_train)
9 y_test_pred = arbol1.predict(X_test)
10 print(f'Accuracy para train de: {accuracy_score(y_train,y_train_pred)}')
11 print(f'Accuracy para test de: {accuracy_score(y_test,y_test_pred)}')
12 print(f'Diferencia de precisión entre train y test:
13 {accuracy_score(y_train,y_train_pred)-accuracy_score(y_test,y_test_pred)}')
14 accuracy_arbol_base = accuracy_score(y_test, y_test_pred)
```

Listing 11: Primer modelo de árbol de regresión.

El siguiente paso consiste en buscar los parámetros y las métricas para la evaluación de la bondad del modelo. En este caso se han elegido distintas profundidades máximas entre 5 y 50 para compensar la pérdida de capacidad predictiva con el riesgo de sobreajuste, así como número mínimo de casos (*min_samples_split*) entre 10 y 300. Los criterios elegidos son el de Gini (desigualdad) y el de entropía (aleatoriedad). Las métricas objetivo son las incluidas en la variable *scoring_metrics*.

```
1 params = {
2     'max_depth': [5, 10, 15, 20, 25, 30, 40, 50],
3     'min_samples_split': [10, 15, 30, 50, 100, 200, 300],
4     'criterion': ["gini", "entropy"]}
5 scoring_metrics = ['accuracy', 'precision_macro', 'recall_macro', 'f1_macro']
6 grid_search = GridSearchCV(estimator=arbol1, param_grid=params, cv=4,
7 scoring = scoring_metrics, refit='accuracy')
8 grid_search.fit(X_train, y_train)
9 results = pd.DataFrame(grid_search.cv_results_)
10 best_model = grid_search.best_estimator_
11 print('Parámetros del mejor modelo de árbol por GridSearchCV:')
12 print(best_model)
13 y_train_pred = best_model.predict(X_train)
14 y_test_pred = best_model.predict(X_test)
15 accuracy_arbol_gscv = accuracy_score(y_test, y_test_pred)
16 print(f'Precisión del mejor modelo de árbol por GridSearchCV: {accuracy_arbol_gscv}')
```

Listing 12: Resultados del proceso *GridSearchCV*.

Se obtiene que el mejor modelo, según *GridSearchCV* es aquel con profundidad **max_depth** = 10, divisiones **min_samples_split** = 15 y **criterion** = *entropy*. Al repetir este proceso para cada una de las métricas se puede obtener una estimación del mejor modelo en cada caso, que es aquel que maximiza estas. Para ello se hace una nueva validación cruzada cambiando los parámetros en la variable *refit*:

```

1  # Accuracy:
2  grid_search = GridSearchCV(estimator=arbol1, param_grid=params, cv=4,
3                           scoring = scoring_metrics, refit='accuracy')
4  grid_search.fit(X_train, y_train)
5  best_model = grid_search.best_estimator_
6  print("El mejor modelo para 'Accuracy' es:")
7  print(grid_search.best_estimator_)
8  print()
9  # Precision:
10 grid_search = GridSearchCV(estimator=arbol1, param_grid=params, cv=4,
11                            scoring = scoring_metrics, refit='precision_macro')
12 grid_search.fit(X_train, y_train)
13 best_model = grid_search.best_estimator_
14 print("El mejor modelo para 'Precision' es:")
15 print(grid_search.best_estimator_)
16 print()
17 # Recall:
18 grid_search = GridSearchCV(estimator=arbol1, param_grid=params, cv=4,
19                            scoring = scoring_metrics, refit='recall_macro')
20 grid_search.fit(X_train, y_train)
21 best_model = grid_search.best_estimator_
22 print("El mejor modelo para 'Recall' es:")
23 print(grid_search.best_estimator_)
24 print()
25 # F1:
26 grid_search = GridSearchCV(estimator=arbol1, param_grid=params, cv=4,
27                            scoring = scoring_metrics, refit='f1_macro')
28 grid_search.fit(X_train, y_train)
29 best_model = grid_search.best_estimator_
30 print("El mejor modelo para 'F1' es:")
31 print(grid_search.best_estimator_)

```

Listing 13: Búsqueda del mejor modelo para cada métrica.

Los resultados de esta búsqueda se pueden encontrar en la siguiente tabla:

Métrica	Criterio	max_depth	min_samples_split
Accuracy	<i>entropy</i>	20	10
Precision	<i>entropy</i>	20	10
Recall	<i>entropy</i>	10	10
F1	<i>entropy</i>	10	10

Tabla 1: Mejores modelos para cada métrica según el método *grid_search.best_estimator_*

A partir de esta tabla se puede empezar a considerar un modelo que siga el criterio de entropía, con una profundidad máxima de entre 10 y 20 y un número mínimo de *splits* de 10. Sin embargo, conviene comprobar el comportamiento de las distintas métricas en función de los parámetros *max_depth*, *min_samples_split* con el fin de evitar el sobreajuste o la pérdida de capacidad predictiva ya mencionados.

En las gráficas de la figura 3 se puede observar que, para todas las métricas, se alcanza un máximo local en el promedio cuando la profundidad máxima del árbol es igual a 10. Además, en el caso de 10 *splits*, el máximo absoluto se alcanza para una profundidad de 20 nodos. A juzgar por las gráficas, los valores óptimos para maximizar las métricas sin perder información (esto es, eligiendo una profundidad o un número de *splits* muy bajos) sería **max_depth** = 10 y **min_samples_split** = 10. Esta elección de parámetros se asimila a la búsqueda *grid search* del mejor estimador para la métrica **accuracy** y no se aleja de los valores de las búsquedas restantes (ver tabla 1).

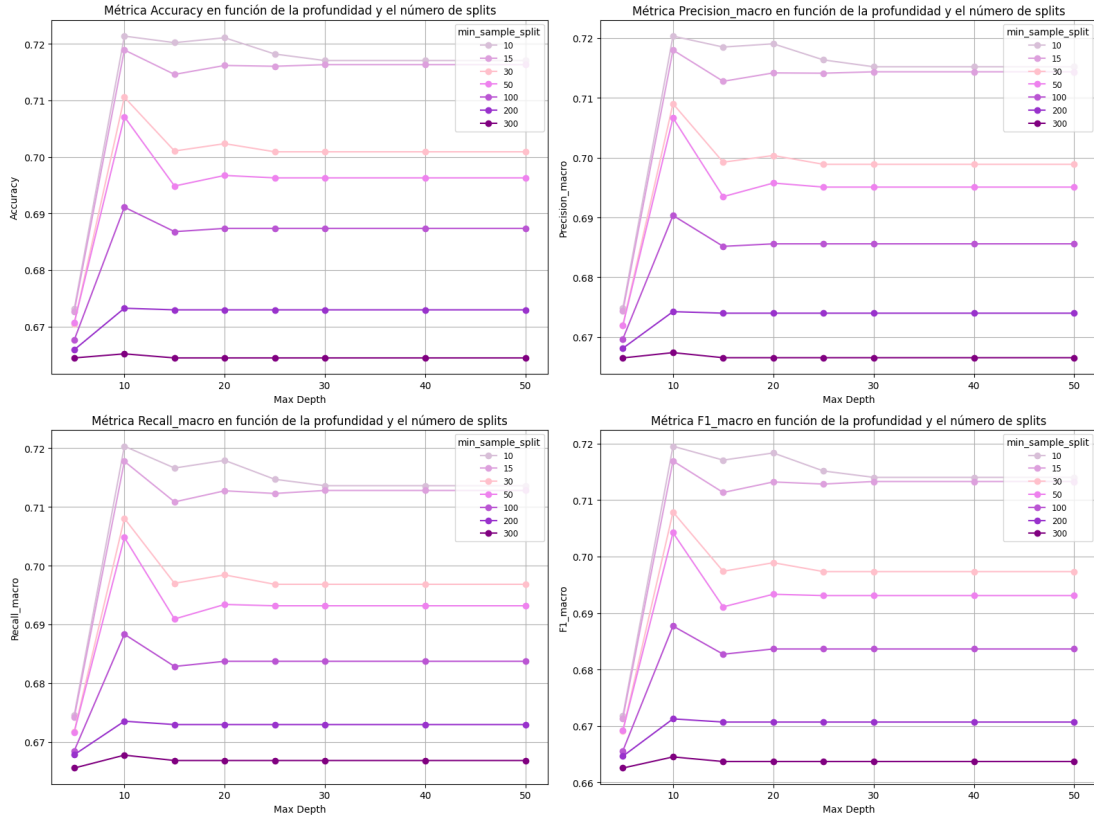


Figura 3: Evolución de las métricas del árbol en función de los parámetros.

Por lo tanto, se van a considerar todos los modelos con una profundidad máxima de entre 10 y 20 nodos (incluidos), así como un mínimo de divisiones de 10 y que sigan el criterio de aleatoriedad (entropía):

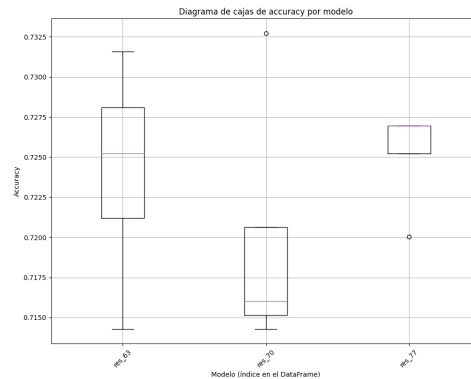


Figura 4: Diagrama de cajas para accuracy del árbol. Se muestran aquellos modelos con $10 < \text{max_depth} < 20$ y $\text{min_samples_split} = 10$.

De todos ellos, el mejor candidato es el que tiene índice 63 (*res_63*) ya que, aunque no tiene la mayor mediana, posee un rango intercuartílico no demasiado amplio y no tiene valores anómalos (no como otros candidatos parecidos, como el resultado *res_77*). El modelo 63 posee los parámetros **max_depth** = 10 y **min_samples_split** = 10.

```

1 print(results.iloc[63].params)
2 modelo_63 = grid_search.cv_results_['params'][63]
3 best_model = grid_search.estimator.set_params(**modelo_63)
4 y_test_pred = best_model.predict(X_test)
5 accuracy_arbol_bmanual = accuracy_score(y_test, y_test_pred)
6 print(f'Precisión del mejor modelo de árbol: {accuracy_arbol_bmanual}')

```

Listing 14: Elección del mejor modelo por búsqueda manual.

Por otro lado, si se ordenan los resultados y se genera otro diagrama de cajas, se puede obtener una nueva estimación de cuál es el mejor modelo según *GridSearchCV*. En este caso, el mejor modelo obtenido (*res_2*) coincide con el mejor modelo del *boxplot* de la figura 4:

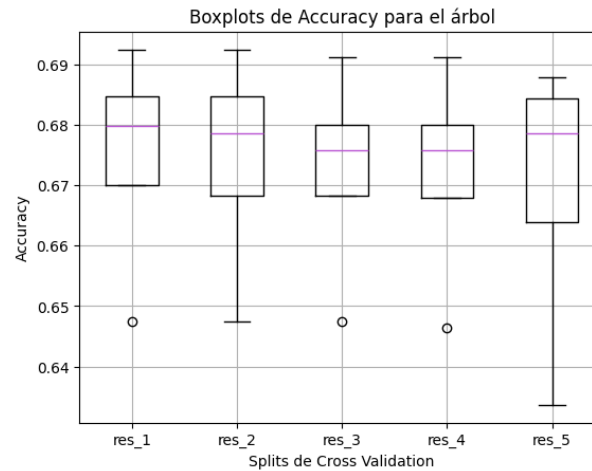


Figura 5: Diagrama de cajas para accuracy del árbol. Se muestran aquellos modelos con $10 < \text{max_depth} < 20$ y $\text{min_samples_split} = 10$.

3.1. Resultados del modelo de árbol.

La siguiente tabla consiste en una recopilación de los mejores modelos según los distintos tipos de búsqueda de parámetros realizadas (mediante *GridSearchCV*, manual o con gráficas, mediante *boxplot*, etc).

Modelo	criterion	max_depth	min_samples_split
Árbol base	entropy	25	200
<i>GridSearchCV</i>	entropy	20	10
Manual	entropy	10	10
<i>Boxplot</i>	entropy	10	10

Tabla 2: Parámetros de los mejores modelos de árbol encontrados.

Y los resultados para **accuracy** de estos modelos se pueden encontrar en la siguiente imagen:

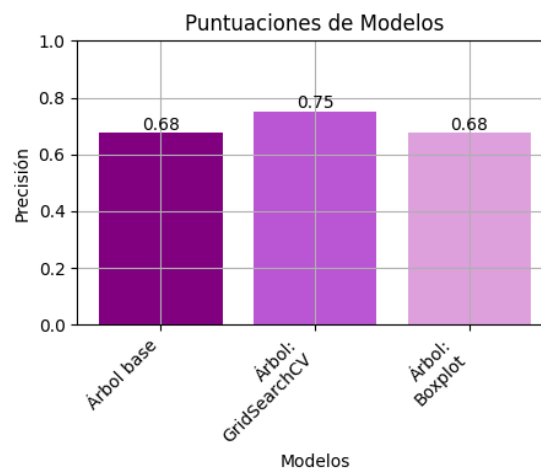


Figura 6: Resultados de los modelos de árbol.

Por lo tanto, el mejor modelo es aquel con $\text{max_depth} = 20$ y $\text{min_samples_split} = 10$. Los resultados para ambas submuestras según este modelo son los siguientes:

Matriz de Confusión:					
[[1735 150]					
[127 1460]]					
Medidas de Desempeño:					
	precision	recall	f1-score	support	
0	0.93	0.92	0.93	1885	
1	0.91	0.92	0.91	1587	
accuracy			0.92	3472	
macro avg	0.92	0.92	0.92	3472	
weighted avg	0.92	0.92	0.92	3472	

Matriz de Confusión:					
[[346 96]					
[119 307]]					
Medidas de Desempeño:					
	precision	recall	f1-score	support	
0	0.74	0.78	0.76	442	
1	0.76	0.72	0.74	426	
accuracy			0.75	868	
macro avg	0.75	0.75	0.75	868	
weighted avg	0.75	0.75	0.75	868	

Figura 7: Medidas de bondad del ajuste en la submuestra de entrenamiento (izqa) y de prueba (dcha).

Y las curvas ROC para ambas submuestras son las siguientes:

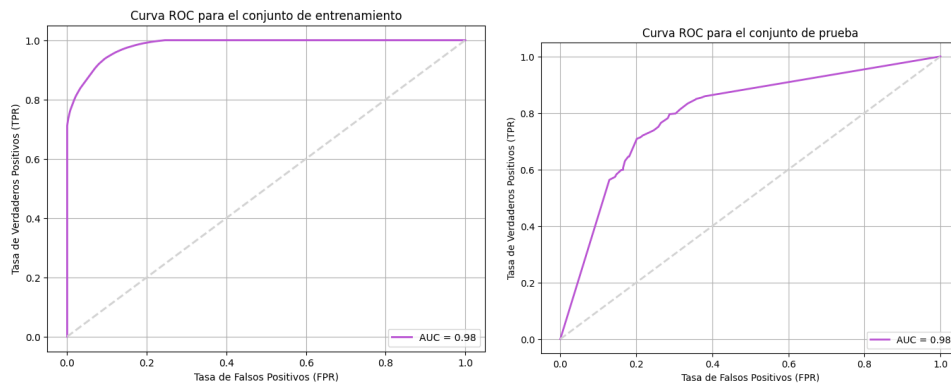


Figura 8: Curvas ROC para la submuestra de entrenamiento (izqa) y de prueba (dcha).

La importancia de las variables para el árbol ganador es la siguiente:



Figura 9: Importancia de las variables para la toma de decisiones en el árbol de la figura ??.

En esta gráfica se puede observar que las variables **Price**, **Mileage** y **Engine volume** son las que más contribuyen al desarrollo del árbol (con contribuciones menos importantes pero notables de **Levy**, **Prod. year**, etcétera).

Y por último, el árbol de clasificación asociado, con profundidad 20, un número mínimo de divisiones de 10 y según el criterio de entropía, se genera visualmente con el siguiente código. Es un árbol en el que el nodo padre se divide en numerosos nodos hijo, dando lugar a un esquema muy profundo y difícil de interpretar::

```

1 plt.figure(figsize=(20, 15))
2 plot_tree(best_model, feature_names=X.columns.tolist(), filled=True, proportion = True)
3 plt.show()

```

Listing 15: Código para generar el árbol.

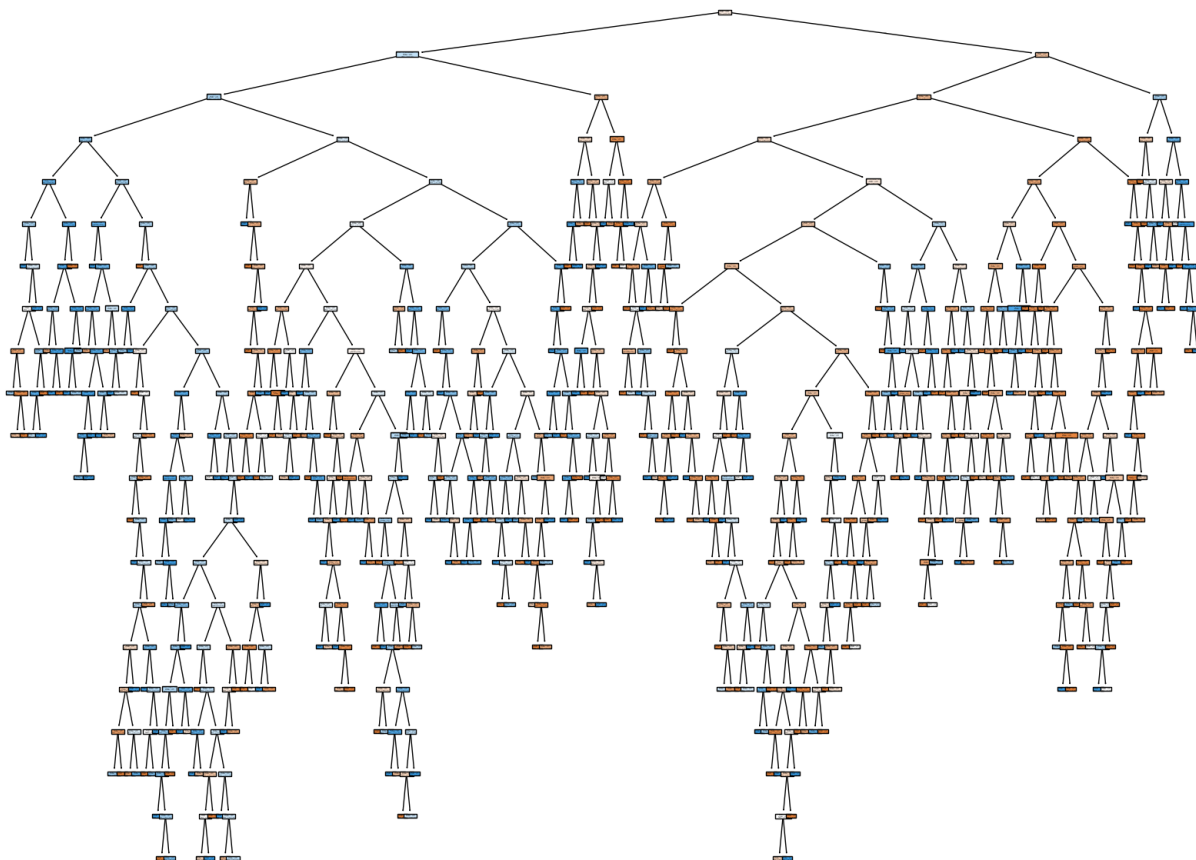


Figura 10: Árbol de decisiones asociado al mejor estimador elegido.

Nota: Las normas del árbol, que contienen toda la información sobre los nodos de este, se obtienen ejecutando el código del siguiente listado. Debido a la extensión del texto, estas no se incluyen en este documento, pero sí en un archivo .txt adjunto.

```

1 rules = export_text(best_model, feature_names = list(X.columns), show_weights = True)
2 with open('decision_tree_rules.txt', 'w', encoding='utf-8') as f:
3     f.write(rules)
4 print(rules)

```

Listing 16: Normas de clasificación del árbol de la figura 10.

4. Modelo de *Random Forest*.

En esta sección se va a desarrollar un modelo de *Random Forest*. Para ello es necesario llevar a cabo un proceso de ensamblado de cierto número de árboles conocido como *Bagging* (*Bootstrap Averaging*). En este proceso se utilizan distintas submuestras de los datos para el desarrollo de los árboles, alternando entre diferentes parámetros como el tamaño de dichas submuestras, el número de árboles desarrollados, si hay reemplazamiento de las submuestras o no, etc. Este modelo de múltiples árboles permite paliar las desventajas que los árboles simples conllevan, como la poca fiabilidad del ajuste o la eficacia de las predicciones.

Con un modelo básico de *Random Forest* con una elección arbitraria de parámetros que den lugar a un sobreajuste bajo se consigue un valor de **accuracy** parecido, tal y como se puede ver con las siguientes líneas. Estas incluyen también valores de precisión para las predicciones de ambos conjuntos (entrenamiento y prueba):

```
1 RF_model = RandomForestClassifier(n_estimators = 100, bootstrap = True, max_depth = 25,
2     min_samples_split = 10, criterion = 'entropy', min_samples_leaf = 10,
3     random_state=123)
4 RF_model.fit(X_train, y_train)
5 y_pred_rf = RF_model.predict(X_test)
6 accuracy_rf = accuracy_score(y_test, y_pred_rf)
7 print(f'Precisión del modelo con RF estándar: {accuracy_rf}')
8 y_train_pred = RF_model.predict(X_train)
9 y_test_pred = RF_model.predict(X_test)
10 print(f'Se tiene un accuracy para train de: {accuracy_score(y_train,y_train_pred)}')
11 print(f'Se tiene un accuracy para test de: {accuracy_score(y_test,y_test_pred)}')
```

Listing 17: Primer modelo de *Random Forest*.

Al igual que en la sección anterior, se puede hacer una búsqueda de los parámetros óptimos:

```
1 params = {
2     'n_estimators': [50, 100, 200],
3     'max_depth': [10, 15, 20, 25],
4     'bootstrap': [True, False],
5     'min_samples_leaf': [5, 10, 20, 30],
6     'min_samples_split': [2, 5, 10, 15, 20, 30, 50, 100],
7     'criterion': ["gini", "entropy"]}
8
9 scoring_metrics = ['accuracy', 'precision_macro', 'recall_macro', 'f1_macro']
10 grid_search_RF = GridSearchCV(estimator=RF_model, param_grid=params, cv=4,
11     scoring = scoring_metrics, refit='accuracy')
12 grid_search_RF.fit(X_train, y_train)
```

Listing 18: Búsqueda *grid search* para el modelo de *Random Forest*.

El mejor modelo dado por la validación cruzada es el siguiente, que además es fruto de la combinación de los parámetros de la tabla 3. El código también incluye los resultados de las predicciones de dicho modelo:

```
1 best_model_RF = grid_search_RF.best_estimator_
2 print(grid_search_RF.best_estimator_)
3 y_pred_rf = best_model_RF.predict(X_test)
4 accuracy_rf_gscv = accuracy_score(y_test, y_pred_rf)
5 print(f'Precisión del modelo: {accuracy_rf_gscv}')
6
7 y_train_pred = best_model_RF.predict(X_train)
8 y_test_pred = best_model_RF.predict(X_test)
9 print(f'Se tiene un accuracy para train de: {accuracy_score(y_train,y_train_pred)}')
10 print(f'Se tiene un accuracy para test de: {accuracy_score(y_test,y_test_pred)}')
```

Listing 19: Mejor modelo por validación cruzada de *GridSearchCV* para *Random Forest*.

Y si se ordenan de forma manual los resultados de la validación cruzada, se puede obtener el siguiente diagrama de cajas para la elección de modelos:

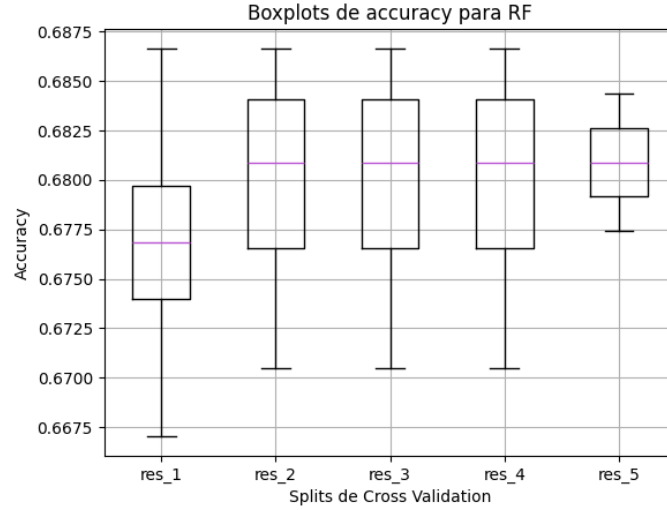


Figura 11: Diagrama de cajas para **accuracy** del modelo de *Random Forest*.

Nuevamente, es posible generar gráficas que ayuden a la elección de parámetros para maximizar **accuracy**. En este caso esas gráficas tienen la siguiente forma:

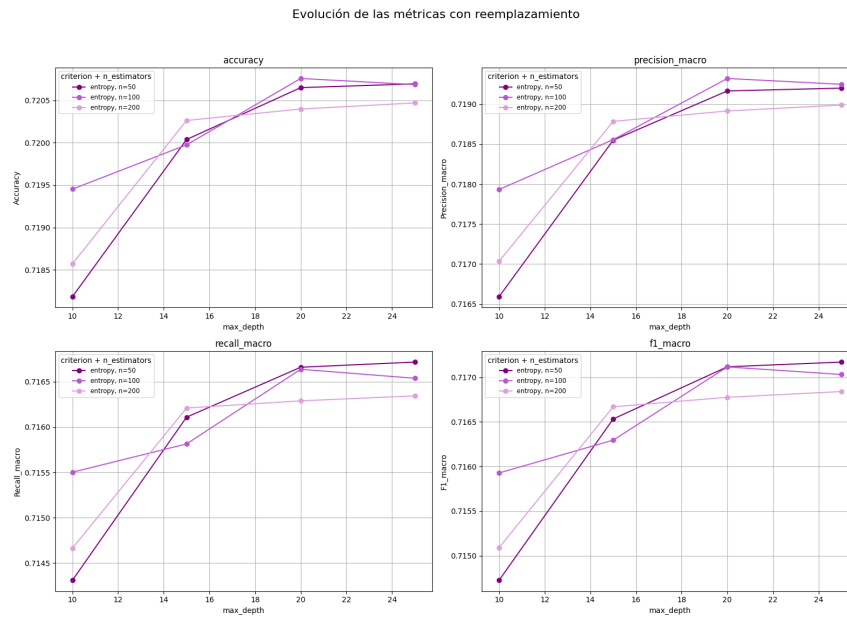


Figura 12: Evolución de la precisión del modelo de RF en función de los parámetros sin reemplazamiento.

Se puede observar que los resultados tienen una representación diferente y menos intuitiva que en los anteriores casos. Por otro lado, si se representan las mismas medidas para *bootstrap = False* se obtiene:

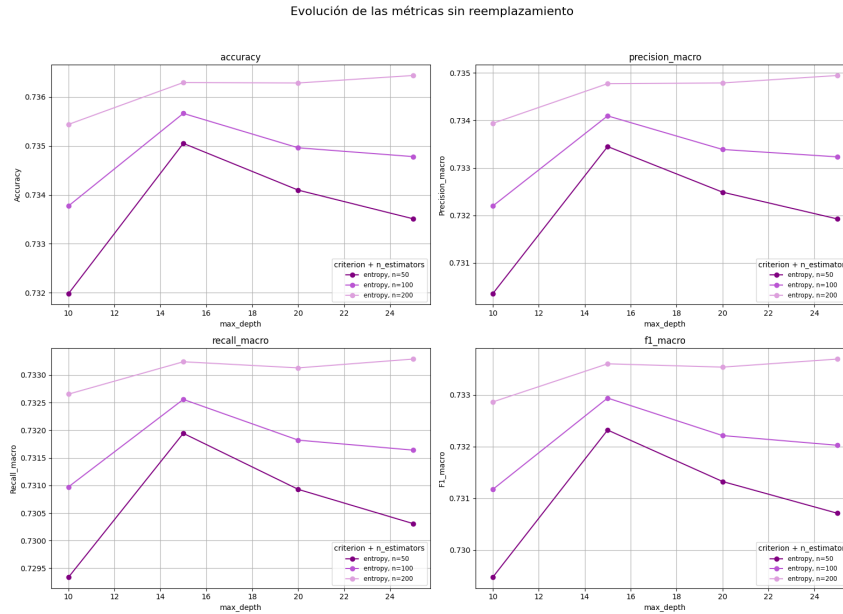


Figura 13: Evolución de la precisión del modelo de RF en función de los parámetros sin reemplazamiento.

Como en la imagen 13 hay mejores resultados, y estos no varían mucho entre 50 y 200 estimadores, se eligen los parámetros que maximicen la precisión sin elevar el tiempo de computación de forma excesiva.

De todos los modelos que cumplen con las condiciones de los máximos de las figuras 12 y 13, los mejores modelos son los siguientes:

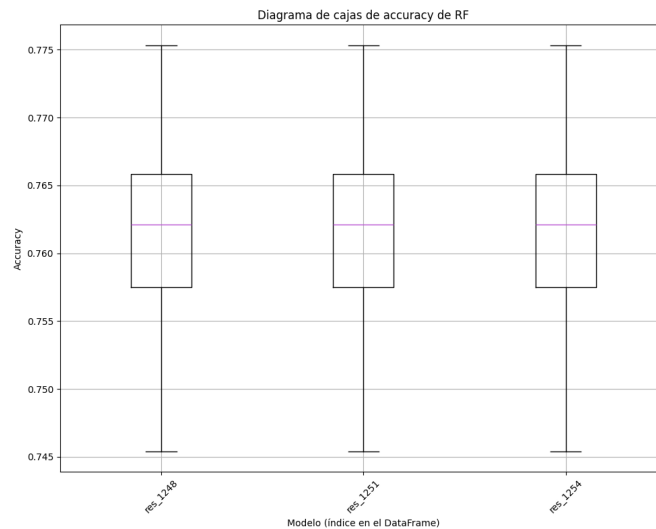


Figura 14: Diagrama de cajas de accuracy para RF. Se muestran modelos con **n_estimators** = 50, **max_depth** = 15 y media de **accuracy** alta.

Los resultados de estos modelos son idénticos y sus parámetros sólo difieren en **min_samples_split**. Se puede escoger cualquiera de los tres:

```
1 print(results.iloc[1254].params)
2 modelo_1254 = grid_search_RF.cv_results_['params'][1254]
3 best_model = grid_search_RF.estimator.set_params(**modelo_1254)
4 y_train_pred = best_model.predict(X_train)
5 y_test_pred = best_model.predict(X_test)
6 accuracy_modelo_1254 = accuracy_score(y_test, y_test_pred)
7 print(f'Precisión del modelo 1254 por búsqueda manual: {accuracy_modelo_1254}')
```

Listing 20: Mejor modelo por búsqueda manual de RF.

4.1. Resultados del modelo de *Random Forest*.

Los mejores modelos obtenidos son:

Modelo	n	max_depth	min_samples_leaf	min_samples_split	bootstrap
<i>Base</i>	100	25	10	10	True
<i>GridSearchCV</i>	50	15	5	2	False
<i>Boxplot</i>	50	10	30	100	True
Manual	50	15	5	10	False

Tabla 3: Búsqueda de parámetros para *Random Forest*. Todos con **criterion** = *entropy*.

La precisión de cada modelo se puede ver en la imagen 15. Más adelante se verá una comparación con la figura 6 para comprobar que la precisión de los modelos aumenta para el algoritmo de *Random Forest*.

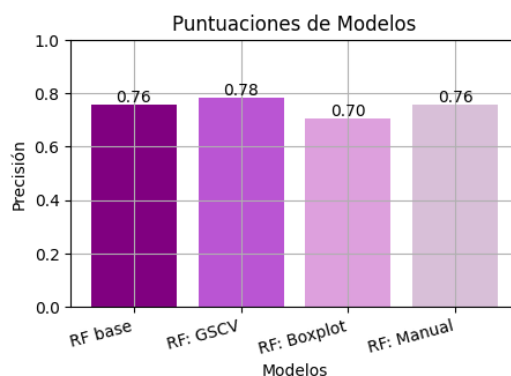


Figura 15: Resultados de los modelos de RF.

Tal y como se puede observar, el mejor modelo es aquel encontrado mediante *GridSearchCV* y el método *best_estimator*. La importancia de las variables es otra medida que conviene comprobar una vez escogido el modelo ganador. En este caso las variables más importantes son las mismas que en el caso del árbol simple, aunque en otro orden.

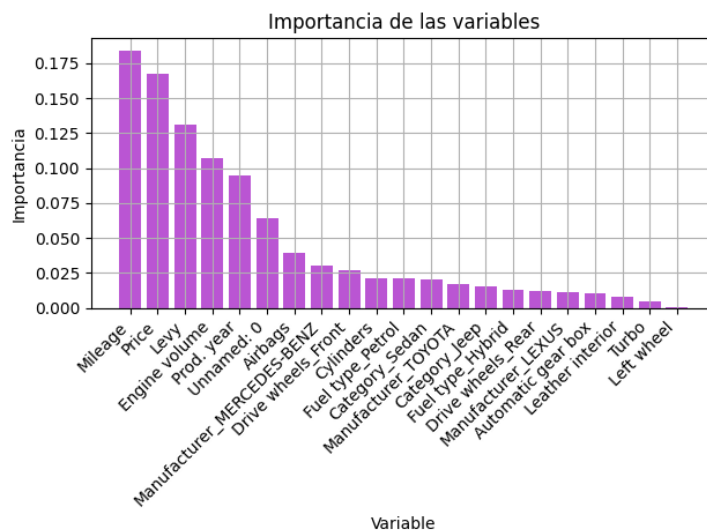


Figura 16: Importancia de las variables en el modelo ganador de RF.

5. Modelo de *XGBoost*.

El último modelo emplea un algoritmo de *Extreme Gradient Boosting* o *XGBoost*. Este consiste en una variante del método de *Gradient Boosting* más rápido y eficiente. Este proceso de ensamblado construye árboles que se emplean como una mejora del árbol anterior. Tienen mayor riesgo de sobreajuste que los modelos anteriores, por lo que es más difícil encontrar los parámetros óptimos.

Un modelo inicial es el siguiente. Los parámetros que dan lugar al clasificador inicial se han elegido de forma arbitraria para reducir el sobreajuste de las predicciones lo máximo posible.

```
1 xgb_classifier = XGBClassifier(booster = 'gbtree', n_estimators = 200,
2                               eta = 0.05, gamma = 3.5, random_state=123, max_depth = 25,
3                               tree_method = 'hist')
4 xgb_classifier.fit(X_train, y_train)
5 y_pred_base = xgb_classifier.predict(X_test)
6 accuracy_xgb = accuracy_score(y_test, y_pred_base)
7 y_train_pred = xgb_classifier.predict(X_train)
8 y_test_pred = xgb_classifier.predict(X_test)
9 print(f'Se tiene un accuracy para train de: {accuracy_score(y_train,y_train_pred)}')
10 print(f'Se tiene un accuracy para test de: {accuracy_score(y_test,y_test_pred)}')
11 print()
12 print(f'Diferencia de precisión entre train y test:
13       {accuracy_score(y_train,y_train_pred) - accuracy_score(y_test,y_test_pred)}')
```

Listing 21: Clasificador inicial del modelo XGB.

Sobre este clasificador inicial se construye una red de árboles sobre la cual se busca la mejor combinación de parámetros para maximizar la precisión del modelo. Se ha comprobado que el número de estimadores (árboles) y la profundidad de los mismos no contribuyen al sobreajuste cuando el resto de parámetros tienen valores parecidos a los siguientes:

```
1 params = {
2     'n_estimators': [150, 200, 250],
3     'eta' : [0.01, 0.05, 0.1, 0.15, 0.2],
4     'gamma' : [0.5, 1, 1.5, 2],
5     'max_depth': [3, 5, 10, 15]}
6
7 scoring_metrics = ['accuracy', 'precision_macro', 'recall_macro', 'f1_macro']
8 grid_search_XGB = GridSearchCV(estimator=xgb_classifier,
9                                param_grid=params,
10                                cv=4, scoring = scoring_metrics, refit='accuracy')
11 grid_search_XGB.fit(X_train, y_train)
12 modelo_XGB = grid_search_XGB.best_estimator_
13 y_train_pred_xgb = modelo_XGB.predict(X_train)
14 y_test_pred_xgb = modelo_XGB.predict(X_test)
15 accuracy_xgb_gscv = accuracy_score(y_test,y_test_pred_xgb)
16 print(f'Se tiene un accuracy para train de: {accuracy_score(y_train,y_train_pred_xgb)}')
17 print(f'Se tiene un accuracy para test de: {accuracy_score(y_test,y_test_pred_xgb)}')
18 print()
19 print('Resultados para Modelo')
20 print(classification_report(y_test, y_test_pred_xgb))
21 print()
22 print('Parámetros del mejor estimador:')
23 params_interes = ['n_estimators', 'max_depth', 'eta', 'gamma']
24 for param in params_interes:
25     valor = modelo_XGB.get_params()[param]
26     print(f'- {param}: {valor}')
```

Listing 22: Clasificador inicial del modelo XGB y búsqueda *GridSearchCV*.

Se puede comprobar que a pesar de intentar paliar el sobreajuste entre las predicciones de entrenamiento y de prueba, el modelo da lugar a mucha diferencia en la precisión de estas. Si se genera un diagrama de cajas con los mejores resultados de la búsqueda se obtiene lo siguiente:

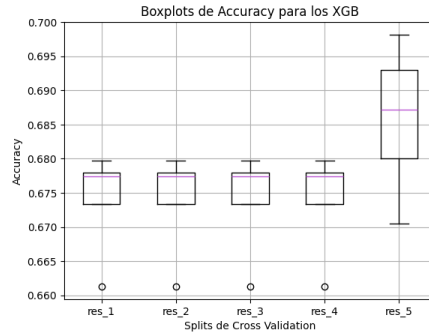


Figura 17: Diagrama de cajas para accuracy del modelo XGB.

Por otro lado, la evolución de la precisión del modelo en función de estos parámetros es la siguiente:

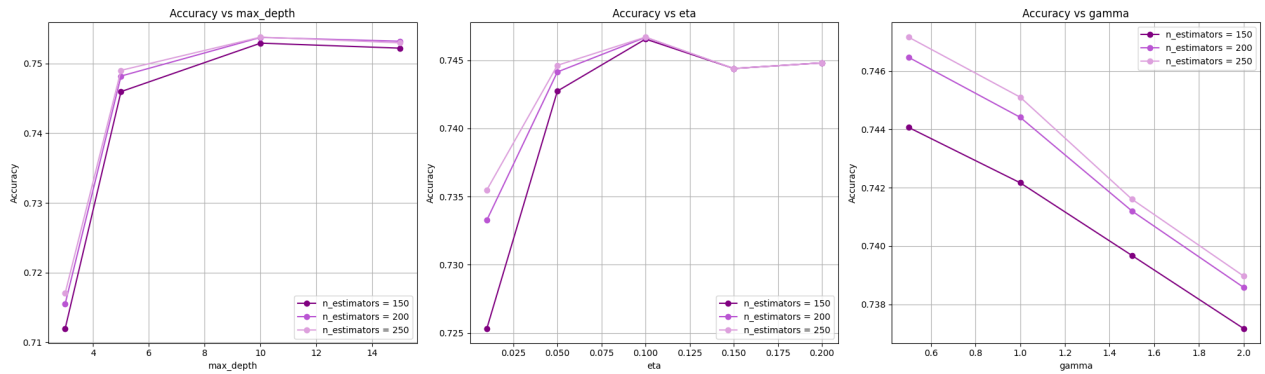


Figura 18: Evolución de la precisión del modelo XGB en función de los parámetros.

El número de árboles no afecta de manera significativa a la precisión del modelo, por lo que es mejor escoger el valor que minimice el tiempo de computación (el menor de todos). Sin embargo, la elección de otros parámetros se puede optimizar. Por ejemplo, los mejores parámetros para **accuracy** son **max_depth** = 10, **eta** = 0,1 y **gamma** = 0,5. La búsqueda manual proporciona los siguientes resultados, donde claramente el mejor modelo es aquel con índice 203 debido a su poca dispersión y bajo rango intercuartílico:

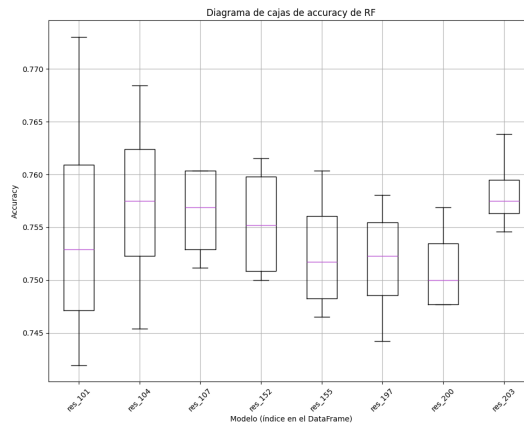


Figura 19: Diagrama de cajas de **accuracy** para XGB. Se muestran modelos con los parámetros determinados por la figura 18.

5.1. Resultados del modelo de *XGBoost*.

Los parámetros de todos los modelos medidos en esta sección se encuentran en la siguiente tabla:

Búsqueda	n_estimators	max_depth	eta	gamma
Base	200	25	3,5	0,05
<i>GridSearchCV</i>	150	10	0,20	1,00
<i>Boxplot</i>	200	3	0,01	1,50
Manual	250	15	0,20	0,50

Tabla 4: Búsqueda de parámetros para XGB.

Estos modelos se pueden ver comparados en cuanto a su precisión en el siguiente gráfico de barras:

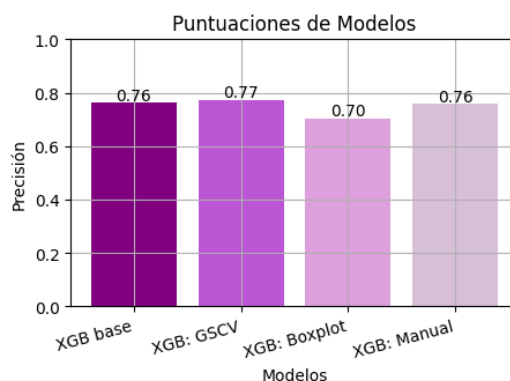


Figura 20: Puntuaciones de los modelos de XGB.

Por lo que el modelo ganador vuelve a ser aquel encontrado por la búsqueda con *GridSearchCV*.

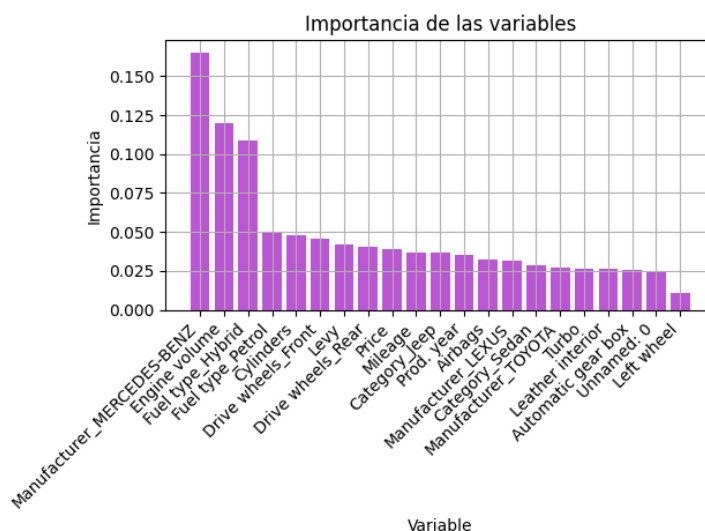


Figura 21: importancia de las variables para el modelo ganador de XGB.

Se puede comprobar que en este caso, las variables **Price** y **Mileage** han perdido importancia (tal y como se esperaba según el método de *Extreme Gradient Boosting*). Las variables más importantes en el desarrollo de este modelo son, en orden de importancia: **Manufacturer_MERCEDES-BENZ** (variable *dummie*), **Engine volume** y **Fuel type_Hybrid** (*dummie*).

6. Comparación de modelos y conclusiones.

6.1. Comparación.

A continuación se encuentra una comparación en forma de gráfico de barras que contiene las precisiones de los distintos modelos ganadores de cada sección:

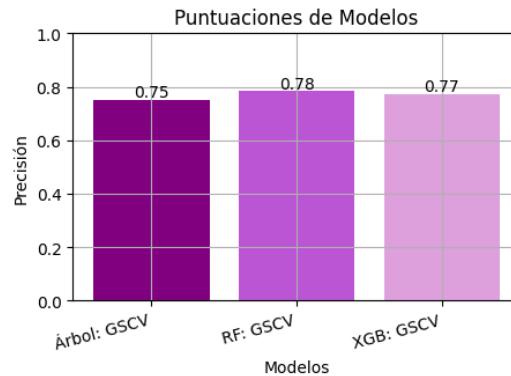


Figura 22: Precisión de los modelos ganadores.

Tal y como se observa, en los tres casos el modelo ganador ha sido aquel obtenido mediante el método *best_estimator* de la función *GridSearchCV*. El modelo con mayor precisión es el de *Random Forest*, aunque las diferencias con los otros dos no son significativas. Esto se debe, probablemente, a la búsqueda manual de parámetros en los estimadores o clasificadores base, que se ha llevado a cabo para reducir el sobreajuste en las primeras predicciones hechas en cada sección. La curva ROC de estos tres modelos ganadores son las siguientes:

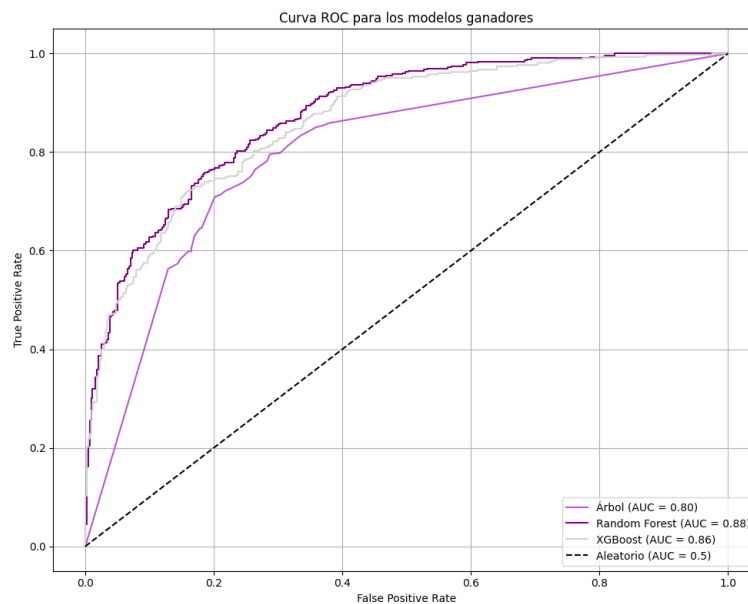


Figura 23: Curvas ROC de los modelos ganadores.

Donde el modelo de *Random Forest* presenta la mayor superficie bajo la curva ROC.

6.2. Conclusiones.

En este trabajo se han comparado tres algoritmos de clasificación: árboles de decisión y *Random Forest* y *Extreme Gradient Boosting* (*XGBoost* o XGB). Cada uno presenta ventajas y limitaciones que se han evaluado en función del coste computacional, precisión del modelo, importancia de variables y análisis de las curvas ROC.

Coste computacional: el árbol de decisiones es el modelo más ligero en términos de tiempo de entrenamiento y consumo de recursos, lo que lo convierte en una buena opción cuando se requiere rapidez o se trabaja con entornos con limitaciones computacionales. En cambio, el "bosque" obtenido por RF supone un mayor coste, especialmente al aumentar la profundidad, ya que implica muchas particiones recursivas sobre el conjunto de datos. Sin embargo, el modelo más costoso es el de XGB, debido a la necesidad de entrenar múltiples árboles (hasta 250 en este caso), lo que implica una mayor demanda computacional tanto en tiempo como en memoria.

Precisión de los modelos: En términos de rendimiento, RF ha obtenido la mayor precisión en el conjunto de prueba, alcanzando un 78.45 %, seguido por XGB con un 77,30 %, y el árbol de decisión con un 75,23 %. Aunque la diferencia entre modelos no es abismal, sí se observa una mejora progresiva al incorporar modelos más complejos y robustos frente al sobreajuste. *Random Forest* destaca especialmente por su capacidad de generalización y su resistencia al sobreajuste frente a árboles individuales.

Importancia de las variables: En la regresión logística, la interpretación de los coeficientes permite identificar qué variables tienen mayor peso en la predicción, destacando por ejemplo las variables relacionadas con la antigüedad/uso y el precio. En el árbol de decisión, la importancia se deriva del orden de las divisiones, y variables como **Price** y **Mileage** aparecen como claves. Para *Random Forest*, al combinar múltiples árboles, se ofrece una medida más estable y agregada de importancia, donde nuevamente variables como **Price**, **Mileage** o **Levy** parecen ser determinantes. En el caso del modelo ganador por XGB, Las variables más importantes son más generales y no tienen tanto que ver con los datos propios de cada vehículo. Estas son **Manufacturer**, **Engine volume** y **Fuel Type**.

Curvas ROC: Las curvas ROC de los tres modelos muestran cómo se comportan frente a distintas tasas de falsos positivos. El árbol de decisiones tiene $AUC = 0,80$, XGB alcanza $AUC = 0,86$ y RF ha logrado $AUC = 0,88$, el más alto de todos. Esto refuerza su capacidad para separar correctamente las clases, con una alta sensibilidad y especificidad. Las curvas indican que *Random Forest* no solo tiene mayor precisión sino también mayor robustez en la clasificación.

En resumen, aunque el árbol simple es adecuado por su simplicidad y bajo coste, el algoritmo de *Random Forest* ofrece el mejor equilibrio entre rendimiento y fiabilidad, a cambio de un coste computacional superior. La elección del modelo dependerá del contexto y de los recursos disponibles.

7. Anexo.

En esta sección se pueden encontrar listados de códigos, tablas e imágenes que aportan información adicional a la práctica pero no son de gran importancia. **NO HACE FALTA CORREGIR ESTA SECCIÓN.**

7.1. Librerías.

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import sklearn
5 from sklearn.tree import plot_tree
6 import seaborn as sns
7 from sklearn.tree import DecisionTreeClassifier, export_text, DecisionTreeRegressor
8 from sklearn.model_selection import train_test_split, GridSearchCV
9 from sklearn.metrics import confusion_matrix, classification_report, accuracy_score,
10 precision_score, recall_score, f1_score, roc_auc_score, roc_curve, auc
11 from sklearn.metrics import make_scorer, mean_absolute_error, mean_squared_error,
12 r2_score
13 from sklearn.ensemble import BaggingClassifier, RandomForestClassifier, BaggingRegressor
14 from sklearn.ensemble import GradientBoostingClassifier, GradientBoostingRegressor
15 from xgboost import XGBClassifier, XGBRegressor
16 from sklearn.model_selection import learning_curve
17
18 import warnings
19 warnings.filterwarnings('ignore')
```

Listing 23: Librerías usadas en la práctica.

7.2. Funciones.

```
1 def ImputacionCuant(var, tipo):
2     """
3     Esta función realiza la imputación de valores faltantes en una variable cuantitativa
4     .
5     Datos de entrada:
6     - var: Serie de datos cuantitativos con valores faltantes a imputar.
7     - tipo: Tipo de imputación ('media', 'mediana' o 'aleatorio').
8     Datos de salida:
9     - Una nueva serie con valores faltantes imputados.
10    """
11    # Realiza una copia de la variable para evitar modificar la original
12    vv = var.copy()
13    if tipo == 'media':
14        # Imputa los valores faltantes con la media de la variable
15        vv[np.isnan(vv)] = round(np.nanmean(vv), 4)
16    elif tipo == 'mediana':
17        # Imputa los valores faltantes con la mediana de la variable
18        vv[np.isnan(vv)] = round(np.nanmedian(vv), 4)
19    elif tipo == 'aleatorio':
20        # Imputa los valores faltantes de manera aleatoria basada en la distribución de
21        # valores existentes
22        x = vv[~np.isnan(vv)]
23        frec = x.value_counts(normalize=True).reset_index()
24        frec.columns = ['Valor', 'Frec']
25        frec = frec.sort_values(by='Valor')
26        frec['FrecAcum'] = frec['Frec'].cumsum()
27        random_values = np.random.uniform(min(frec['FrecAcum']), 1, np.sum(np.isnan(vv)))
28    )
29    imputed_values = list(map(lambda x: list(frec['Valor'][frec['FrecAcum'] <= x])
30    [-1], random_values))
31    vv[np.isnan(vv)] = [round(x, 4) for x in imputed_values]
32
33    return vv
```

Listing 24: Funciones importadas en la práctica.

7.3. Imágenes.

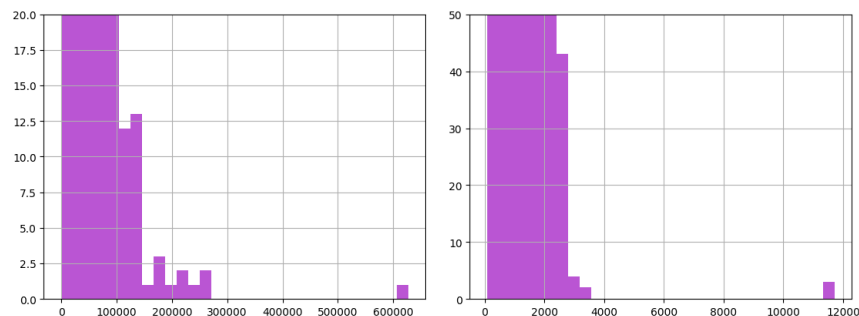


Figura 24: Histograma de las variables **Price** y **Levy**.

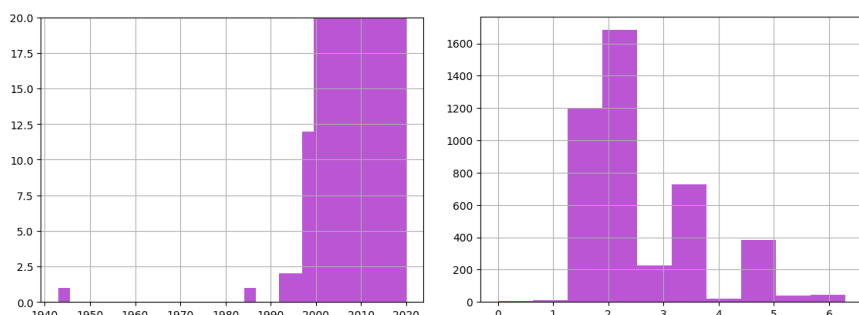


Figura 25: Histograma de las variables **Prod. year** y **Engine volume**.

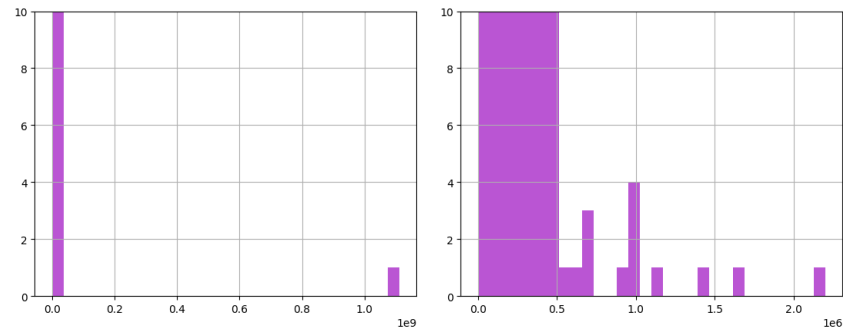


Figura 26: Histograma de la variable **Mileage**.

7.4. Códigos para generar imágenes.

```

1  grid_search = GridSearchCV(estimator=arbol1, param_grid=params, cv=4,
2  scoring = scoring_metrics, refit='accuracy')
3  grid_search.fit(X_train, y_train)
4  results = pd.DataFrame(grid_search.cv_results_)
5  scoring_metrics = ['accuracy', 'precision_macro', 'recall_macro', 'f1_macro']
6  fig, axes = plt.subplots(2, 2, figsize = (16, 12))
7  axes = axes.flatten()
8  palette = ['thistle', 'plum', 'pink', 'violet', 'mediumorchid', 'darkorchid', 'purple']
9
10 for i, metric in enumerate(scoring_metrics):
11     ax = axes[i]
12     mean_col = f'mean_test_{metric}'
13     std_col = f'std_test_{metric}'
14     grouped = results.groupby(['param_max_depth', 'param_min_samples_split'])[[mean_col,

```



```

15         std_col]].mean().reset_index()
16 min_split_vals = sorted(grouped['param_min_samples_split'].unique())
17
18 for j, val in enumerate(min_split_vals):
19     subset = grouped[grouped['param_min_samples_split'] == val]
20     x = subset['param_max_depth']
21     y = subset[mean_col]
22     ax.plot(x, y, marker = 'o', label = f'{val}', color = palette[j])
23
24     ax.set_title(f'Métrica {metric.capitalize()} en función de la profundidad
25                 y el número de splits')
26     ax.set_xlabel('Max Depth')
27     ax.set_ylabel(metric.capitalize())
28     ax.grid(True)
29     ax.legend(title = 'min_sample_split', loc = 'upper right', fontsize = 'small')
30 plt.tight_layout()
31 plt.show()

```

Listing 25: Código para generar la figura 3.

```

1 filtro = (
2     (results['param_max_depth'] >= 10) &
3     (results['param_max_depth'] <= 20) &
4     (results['param_min_samples_split'] == 10) &
5     (results['param_criterion'] == 'entropy')
6 )
7 subset = results[filtro]
8 split_cols = ['split0_test_accuracy', 'split1_test_accuracy', 'split2_test_accuracy',
9              'split3_test_accuracy']
10
11 data = [row[split_cols].values for _, row in subset.iterrows()]
12 labels = [f'res_{i}' for i in subset.index]
13 plt.figure(figsize = (10, 8))
14 plt.boxplot(data, labels = labels, medianprops = {'color': 'mediumorchid'})
15 plt.title('Diagrama de cajas de accuracy por modelo (max_depth = 20, min_samples_split =
16           15,
17           criterion = "entropy")')
18 plt.ylabel('Accuracy')
19 plt.xlabel('Modelo (índice en el DataFrame)')
20 plt.xticks(rotation = 45)
21 plt.grid(True)
22 plt.tight_layout()
23 plt.show()

```

Listing 26: Código para generar la figura 4.

```

1 sorted_results = results.sort_values(by='rank_test_accuracy', ascending=True).head(5)
2 res_1 = results[['split0_test_accuracy', 'split1_test_accuracy', 'split2_test_accuracy',
3                 'split3_test_accuracy']].iloc[0]
4 res_2 = results[['split0_test_accuracy', 'split1_test_accuracy',
5                 'split2_test_accuracy', 'split3_test_accuracy']].iloc[1]
6 res_3 = results[['split0_test_accuracy', 'split1_test_accuracy',
7                 'split2_test_accuracy', 'split3_test_accuracy']].iloc[2]
8 res_4 = results[['split0_test_accuracy', 'split1_test_accuracy',
9                 'split2_test_accuracy', 'split3_test_accuracy']].iloc[3]
10 res_5 = results[['split0_test_accuracy', 'split1_test_accuracy',
11                 'split2_test_accuracy', 'split3_test_accuracy']].iloc[4]
12
13 plt.boxplot(
14     [res_1.values, res_2.values, res_3.values, res_4.values, res_5.values],
15     labels = ['res_1', 'res_2', 'res_3', 'res_4', 'res_5'],
16     medianprops = {'color': 'mediumorchid'})
17 plt.title('Boxplots de Accuracy para los 4 Splits')
18 plt.xlabel('Splits de Cross Validation')
19 plt.ylabel('Accuracy')
20 plt.grid()
21 plt.show()
22 sorted_results['params'].iloc[1]

```

Listing 27: Código para generar la figura 5.

```

1     model_names = ['Árbol base', 'Árbol:\nGridSearchCV', 'Árbol:\nBoxplot']
2     scores = [accuracy_arbol_base, accuracy_arbol_gscv, accuracy_arbol_bmanual]
3     plt.figure(figsize=(5, 3))
4     bars = plt.bar(np.arange(len(model_names)), scores, color=['purple', 'mediumorchid', '
5     plum'])
6     plt.xticks(np.arange(len(model_names)), model_names, rotation=45, ha='right')
7     plt.ylim(0, 1.0)
8     plt.title('Puntuaciones de Modelos')
9     plt.xlabel('Modelos')
10    plt.ylabel('Precisión')
11    plt.grid()
12    for bar, score in zip(bars, scores):
13        plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height() + 0.01, f'{score:.2f}',
14                ha='center', color='black')
15    plt.show()

```

Listing 28: Código para generar la figura 6.

```

1     y_train_auc = pd.get_dummies(y_train, drop_first=True)
2     y_prob_train = best_model.predict_proba(X_train)[: , 1]
3     fpr, tpr, thresholds = roc_curve(y_train_auc, y_prob_train)
4     roc_auc = auc(fpr, tpr)
5     print(f"\nÁrea bajo la curva ROC (AUC): {roc_auc:.2f}")
6     plt.figure(figsize=(8, 6))
7     plt.plot(fpr, tpr, color='mediumorchid', lw=2, label=f'AUC = {roc_auc:.2f}')
8     plt.plot([0, 1], [0, 1], color='lightgrey', lw=2, linestyle='--')
9     plt.xlabel('Tasa de Falsos Positivos (FPR)')
10    plt.ylabel('Tasa de Verdaderos Positivos (TPR)')
11    plt.title('Curva ROC para el conjunto de entrenamiento')
12    plt.legend(loc="lower right")
13    plt.grid()
14    plt.show()
15
16    y_test_auc = pd.get_dummies(y_test, drop_first=True)
17    y_prob_test = best_model.predict_proba(X_test)[: , 1]
18    fpr, tpr, thresholds = roc_curve(y_test_auc, y_prob_test)
19    roc_auc_test = auc(fpr, tpr)
20    print(f"\nÁrea bajo la curva ROC (AUC): {roc_auc:.2f}")
21    plt.figure(figsize=(8, 6))
22    plt.plot(fpr, tpr, color='mediumorchid', lw=2, label=f'AUC = {roc_auc:.2f}')
23    plt.plot([0, 1], [0, 1], color='lightgrey', lw=2, linestyle='--')
24    plt.xlabel('Tasa de Falsos Positivos (FPR)')
25    plt.ylabel('Tasa de Verdaderos Positivos (TPR)')
26    plt.title('Curva ROC para el conjunto de prueba')
27    plt.legend(loc="lower right")
28    plt.grid()
29    plt.show()

```

Listing 29: Código para generar la figura 8.

```

1     df_importancia = pd.DataFrame({'Variable': best_model.feature_names_in_,
2     'Importancia': best_model.feature_importances_}).sort_values(by='Importancia',
3     ascending=False)
4     plt.bar(df_importancia['Variable'], df_importancia['Importancia'], color='mediumorchid')
5     plt.xlabel('Variable')
6     plt.ylabel('Importancia')
7     plt.title('Importancia de las variables')
8     plt.xticks(rotation=45, ha='right')
9     plt.tight_layout()
10    plt.grid()
11    plt.show()

```

Listing 30: Código para generar la figura 9.

```

1     results = pd.DataFrame(grid_search_RF.cv_results_)
2     sorted_results = results.sort_values(by='mean_test_accuracy', ascending=True).head(5)
3
4     res_1 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',
5     'split2_test_accuracy', 'split3_test_accuracy']].iloc[0]
6     res_2 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',

```

```

7         'split2_test_accuracy', 'split3_test_accuracy']].iloc[1]
8     res_3 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',
9         'split2_test_accuracy', 'split3_test_accuracy']].iloc[2]
10    res_4 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',
11        'split2_test_accuracy', 'split3_test_accuracy']].iloc[3]
12    res_5 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',
13        'split2_test_accuracy', 'split3_test_accuracy']].iloc[4]
14
15    plt.boxplot([res_1.values, res_2.values, res_3.values, res_4.values, res_5.values],
16        labels = ['res_1', 'res_2', 'res_3', 'res_4', 'res_5'],
17        medianprops = {'color': 'mediumorchid'})
18    plt.title('Boxplots de Accuracy para los 4 Splits')
19    plt.xlabel('Splits de Cross Validation')
20    plt.ylabel('Accuracy')
21    plt.grid()
22    plt.show()
23
24    best_model_rf_boxplot = RandomForestClassifier(**sorted_results['params'].iloc[4],
25        random_state=123)
26    best_model_rf_boxplot.fit(X_train, y_train)
27    print('Parámetros del mejor modelo de RF por boxplot:')
28    print(best_model_rf_boxplot.get_params())
29    print()
30    y_train_pred = best_model_rf_boxplot.predict(X_train)
31    y_test_pred = best_model_rf_boxplot.predict(X_test)
32    print(f'Se tiene un accuracy para train de: {accuracy_score(y_train, y_train_pred)}')
33    print(f'Se tiene un accuracy para test de: {accuracy_score(y_test, y_test_pred)}')
34    print(f'Diferencia de precisión entre train y test: {accuracy_score(y_train, y_train_pred)
35        - accuracy_score(y_test, y_test_pred)}')
36    print()
37    accuracy_rf_boxplot = accuracy_score(y_test, y_test_pred)
38    print(f'Precisión del modelo: {accuracy_rf_boxplot}')
39    print(f'Resultados para el modelo: {classification_report(y_test, y_test_pred)}')

```

Listing 31: Código para generar la figura 11.

```

1    scoring_metrics = ['accuracy', 'precision_macro', 'recall_macro', 'f1_macro']
2    palette = ['purple', 'mediumorchid', 'plum', 'orchid', 'hotpink', 'violet', 'darkorchid']
3
4    titles = ['con reemplazamiento', 'sin reemplazamiento']
5    results['param_bootstrap'] = results['param_bootstrap'].astype(bool)
6    results['param_criterion'] = results['param_criterion'].astype(str)
7    index = 0
8    for boot in [True, False]:
9        fig, axes = plt.subplots(2, 2, figsize=(16, 12))
10        axes = axes.flatten()
11
12        for crit in ['gini', 'entropy']:
13            filtered = results[
14                (results['param_bootstrap'] == boot) &
15                (results['param_criterion'] == crit)
16            ]
17
18            for i, metric in enumerate(scoring_metrics):
19                ax = axes[i]
20                mean_col = f'mean_test_{metric}'
21
22                grouped = filtered.groupby(
23                    ['param_max_depth', 'param_n_estimators']
24                )[mean_col].mean().reset_index()
25
26                for j, n_est in enumerate(sorted(grouped['param_n_estimators'].unique())):
27                    subset = grouped[grouped['param_n_estimators'] == n_est]
28                    ax.plot(subset['param_max_depth'], subset[mean_col],
29                        marker='o', label=f'{crit}, n={n_est}',
30                        color=palette[j % len(palette)])
31
32                ax.set_title(f'{metric}')
33                ax.set_xlabel('max_depth')
34                ax.set_ylabel(metric.capitalize())

```

```

34         ax.grid(True)
35         ax.legend(title='criterion + n_estimators', fontsize='small')
36
37     plt.suptitle(f'Evolución de las métricas {titles[index]}', fontsize=16)
38     index += 1
39     plt.tight_layout(rect=[0, 0.03, 1, 0.95])
40     plt.show()

```

Listing 32: Código para generar las figuras 12 y 13.

```

1     filtro = (
2         (results['param_max_depth'] == 15) &
3         (results['param_n_estimators'] == 50) &
4         (results['param_criterion'] == 'entropy')
5     )
6     subset = results[filtro]
7     split_cols = ['split0_test_accuracy', 'split1_test_accuracy', 'split2_test_accuracy', 'split3_test_accuracy']
8     subset['mean_accuracy'] = subset[split_cols].mean(axis = 1)
9     subset = subset[subset['mean_accuracy'] > 0.76]
10    data = [row[split_cols].values for _, row in subset.iterrows()]
11    labels = [f'res_{i}' for i in subset.index]
12    plt.figure(figsize = (10, 8))
13    plt.boxplot(data, labels = labels, medianprops = {'color': 'mediumorchid'})
14    plt.title('Diagrama de cajas de accuracy por modelo (media > 0.7)')
15    plt.ylabel('Accuracy')
16    plt.xlabel('Modelo (índice en el DataFrame)')
17    plt.xticks(rotation = 45)
18    plt.grid(True)
19    plt.tight_layout()
20    plt.show()

```

Listing 33: Código para generar la figura 14.

```

1    model_names = ['RF base', 'RF: GSCV', 'RF: Boxplot', 'RF: Manual']
2    scores = [accuracy_rf_base, accuracy_rf_gscv, accuracy_rf_boxplot, accuracy_modelo_1254]
3    plt.figure(figsize=(5, 3))
4    bars = plt.bar(np.arange(len(model_names)), scores, color=['purple', 'mediumorchid', 'plum', 'thistle'])
5    plt.xticks(np.arange(len(model_names)), model_names, rotation=15, ha='right')
6    plt.ylim(0, 1.0)
7    plt.title('Puntuaciones de Modelos')
8    plt.xlabel('Modelos')
9    plt.ylabel('Precisión')
10   plt.grid()
11   for bar, score in zip(bars, scores):
12       plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height() + 0.01, f'{score:.2f}',
13               ha='center', color='black')
14   plt.show()

```

Listing 34: Código para generar la figura 15.

```

1    best_model = grid_search_RF.best_estimator_
2    df_importancia = pd.DataFrame({'Variable': best_model.feature_names_in_, 'Importancia':
3    best_model.feature_importances_}).sort_values(by='Importancia', ascending=False)
4    plt.bar(df_importancia['Variable'], df_importancia['Importancia'], color='mediumorchid')
5    plt.xlabel('Variable')
6    plt.ylabel('Importancia')
7    plt.title('Importancia de las variables')
8    plt.xticks(rotation=45, ha='right')
9    plt.tight_layout()
10   plt.grid()
11   plt.show()

```

Listing 35: Código para generar la figura 16.

```

1    results = pd.DataFrame(grid_search_XGB.cv_results_)
2
3    sorted_results = results.sort_values(by='mean_test_accuracy', ascending=True).head(5)
4    res_1 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',

```

```

5         'split2_test_accuracy', 'split3_test_accuracy']].iloc[0]
6     res_2 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',
7                             'split2_test_accuracy', 'split3_test_accuracy']].iloc[1]
8     res_3 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',
9                             'split2_test_accuracy', 'split3_test_accuracy']].iloc[2]
10    res_4 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',
11                            'split2_test_accuracy', 'split3_test_accuracy']].iloc[3]
12    res_5 = sorted_results[['split0_test_accuracy', 'split1_test_accuracy',
13                            'split2_test_accuracy', 'split3_test_accuracy']].iloc[4]
14
15    plt.boxplot([res_1.values, res_2.values, res_3.values, res_4.values, res_5.values], labels =
16                ['res_1', 'res_2', 'res_3', 'res_4', 'res_5'], medianprops = {'color': 'mediumorchid'})
17    plt.title('Boxplots de Accuracy para los 4 Splits')
18    plt.xlabel('Splits de Cross Validation')
19    plt.ylabel('Accuracy')
20    plt.grid()
21    plt.show()
22    print('Parámetros del mejor estimador:')
23    print(sorted_results['params'].iloc[4])

```

Listing 36: Código para generar la figura 17.

```

1     results = pd.DataFrame(grid_search_XGB.cv_results_)
2     fig, axes = plt.subplots(1, 3, figsize = (20, 6))
3     axes = axes.flatten()
4     palette = ['purple', 'mediumorchid', 'plum']
5
6     for i, n in enumerate(sorted(results['param_n_estimators'].unique())):
7         subset = results[results['param_n_estimators'] == n]
8         grouped = subset.groupby('param_max_depth')['mean_test_accuracy'].mean().reset_index()
9         axes[0].plot(grouped['param_max_depth'], grouped['mean_test_accuracy'],
10                      label = f'n_estimators = {n}', marker = 'o', color = palette[i])
11    axes[0].set_title('Accuracy vs max_depth')
12    axes[0].set_xlabel('max_depth')
13    axes[0].set_ylabel('Accuracy')
14    axes[0].legend()
15    axes[0].grid(True)
16
17    for i, n in enumerate(sorted(results['param_n_estimators'].unique())):
18        subset = results[results['param_n_estimators'] == n]
19        grouped = subset.groupby('param_eta')['mean_test_accuracy'].mean().reset_index()
20        axes[1].plot(grouped['param_eta'], grouped['mean_test_accuracy'],
21                     label = f'n_estimators = {n}', marker = 'o', color = palette[i])
22    axes[1].set_title('Accuracy vs eta')
23    axes[1].set_xlabel('eta')
24    axes[1].set_ylabel('Accuracy')
25    axes[1].legend()
26    axes[1].grid(True)
27
28    for i, n in enumerate(sorted(results['param_n_estimators'].unique())):
29        subset = results[results['param_n_estimators'] == n]
30        grouped = subset.groupby('param_gamma')['mean_test_accuracy'].mean().reset_index()
31        axes[2].plot(grouped['param_gamma'], grouped['mean_test_accuracy'],
32                     label = f'n_estimators = {n}', marker = 'o', color = palette[i])
33    axes[2].set_title('Accuracy vs gamma')
34    axes[2].set_xlabel('gamma')
35    axes[2].set_ylabel('Accuracy')
36    axes[2].legend()
37    axes[2].grid(True)
38    plt.tight_layout()
39    plt.show()

```

Listing 37: Código para generar la figura 18.

```

1     filtro = (
2         (results['param_max_depth'] >= 5) &
3         (results['param_max_depth'] <= 15) &
4         (results['param_eta'] >= 0.1) &
5         (results['param_gamma'] == 0.5) &
6         (results['param_n_estimators'] == 250))
7     subset = results[filtro]

```

```

8 split_cols = ['split0_test_accuracy', 'split1_test_accuracy', 'split2_test_accuracy', '
split3_test_accuracy']
9 subset['mean_accuracy'] = subset[split_cols].mean(axis = 1)
10 subset = subset[subset['mean_accuracy'] > 0.75]
11 data = [row[split_cols].values for _, row in subset.iterrows()]
12 labels = [f'res_{i}' for i in subset.index]
13 plt.figure(figsize = (10, 8))
14 plt.boxplot(data, labels = labels, medianprops = {'color': 'mediumorchid'})
15 plt.title('Diagrama de cajas de accuracy de RF')
16 plt.ylabel('Accuracy')
17 plt.xlabel('Modelo (índice en el DataFrame)')
18 plt.xticks(rotation = 45)
19 plt.grid(True)
20 plt.tight_layout()
21 plt.show()

```

Listing 38: Código para generar la figura 19.

```

1 model_names = ['XGB base', 'XGB: GSCV', 'XGB: Boxplot', 'XGB: Manual']
2 scores = [accuracy_xgb_base, accuracy_xgb_gscv, accuracy_xgb_boxplot,
accuracy_modelo_203]
3 plt.figure(figsize=(5, 3))
4 bars = plt.bar(np.arange(len(model_names)), scores, color=['purple', 'mediumorchid', '
plum', 'thistle'])
5 plt.xticks(np.arange(len(model_names)), model_names, rotation=15, ha='right')
6 plt.ylim(0, 1.0)
7 plt.title('Puntuaciones de Modelos')
8 plt.xlabel('Modelos')
9 plt.ylabel('Precisión')
10 plt.grid()
11 for bar, score in zip(bars, scores):
12     plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height() + 0.01, f'{score:.2f}',
ha='center', color='black')
13 plt.show()

```

Listing 39: Código para generar la figura 20.

```

1 best_model = grid_search_XGB.best_estimator_
2 df_importancia = pd.DataFrame({'Variable': best_model.feature_names_in_,
3 'Importancia': best_model.feature_importances_}).sort_values(by='Importancia',
ascending=False)
4 plt.bar(df_importancia['Variable'], df_importancia['Importancia'], color='mediumorchid')
5 plt.xlabel('Variable')
6 plt.ylabel('Importancia')
7 plt.title('Importancia de las variables')
8 plt.xticks(rotation=45, ha='right')
9 plt.tight_layout()
10 plt.grid()
11 plt.show()

```

Listing 40: Código para generar la figura 21.

```

1 model_names = ['Árbol: GSCV', 'RF: GSCV', 'XGB: GSCV']
2 scores = [accuracy_arbol_gscv, accuracy_rf_gscv, accuracy_xgb_gscv]
3 plt.figure(figsize=(5, 3))
4 bars = plt.bar(np.arange(len(model_names)), scores, color=['purple', 'mediumorchid', '
plum', 'thistle'])
5 plt.xticks(np.arange(len(model_names)), model_names, rotation=15, ha='right')
6 plt.ylim(0, 1.0)
7 plt.title('Puntuaciones de Modelos')
8 plt.xlabel('Modelos')
9 plt.ylabel('Precisión')
10 plt.grid()
11 for bar, score in zip(bars, scores):
12     plt.text(bar.get_x() + bar.get_width() / 2, bar.get_height() + 0.01, f'{score:.2f}',
ha='center', color='black')

```

Listing 41: Código para generar la figura 22.

```

1 # Modelo 1: árbol
2 modelo_1 = grid_search.best_estimator_
3 modelo_1.fit(X_train, y_train)
4 y_prob1 = modelo_1.predict_proba(X_test)[:, 1]
5 fpr1, tpr1, _ = roc_curve(y_test, y_prob1)
6 roc_auc1 = auc(fpr1, tpr1)
7
8 # Modelo 2: RF
9 modelo_2 = grid_search_RF.best_estimator_
10 modelo_2.fit(X_train, y_train)
11 y_prob2 = modelo_2.predict_proba(X_test)[:, 1]
12 fpr2, tpr2, _ = roc_curve(y_test, y_prob2)
13 roc_auc2 = auc(fpr2, tpr2)
14
15 # Modelo 3: XGB
16 modelo_3 = grid_search_XGB.best_estimator_
17 modelo_3.fit(X_train, y_train)
18 y_prob3 = modelo_3.predict_proba(X_test)[:, 1]
19 fpr3, tpr3, _ = roc_curve(y_test, y_prob3)
20 roc_auc3 = auc(fpr3, tpr3)
21
22 # Gráfica ROC
23 plt.figure(figsize=(10, 8))
24 plt.plot(fpr1, tpr1, label = f'Árbol (AUC = {roc_auc1:.2f})', color='mediumorchid')
25 plt.plot(fpr2, tpr2, label = f'Random Forest (AUC = {roc_auc2:.2f})', color='purple')
26 plt.plot(fpr3, tpr3, label = f'XGBoost (AUC = {roc_auc3:.2f})', color='lightgrey')
27 plt.plot([0, 1], [0, 1], 'k--', label = 'Aleatorio (AUC = 0.5)')
28
29 plt.title('Curva ROC para los modelos ganadores')
30 plt.xlabel('False Positive Rate')
31 plt.ylabel('True Positive Rate')
32 plt.legend(loc = 'lower right')
33 plt.grid(True)
34 plt.tight_la

```

Listing 42: Código para generar la figura 23.