
Práctica 3

Inteligencia de Negocio

Amin Kasrou Aouam



**UNIVERSIDAD
DE GRANADA**

2021-01-12

Índice

Práctica 3	3
Introducción	3
Preprocesamiento de datos	3
Valores nulos	3
Valores no numéricos	3
Balanceo de clases	4
Elección de algoritmo	5
Resultados obtenidos	5
Análisis de resultados	5

Práctica 3

Introducción

En esta práctica, resolveremos un problema de clasificación multiclase, en concreto, trataremos de predecir la categoría de precio de una serie de coches

Preprocesamiento de datos

Valores nulos

Nuestro *dataset* contiene bastantes valores nulos, optamos por estrategias diferentes según las columnas:

- Eliminación: tipo marchas, descuento, ciudad
- Imputación: asientos, motor cc, potencia

El criterio que seleccionamos es el número de instancias nulas, en el caso de que sean muchas optamos por imputar, para mantener un número adecuado de datos.

La implementación se encuentra en la siguiente función:

```
1 def process_null_values(df_list):
2     drop_columns = ["tipo_marchas", "descuento", "ciudad"]
3     fill_columns = ["asientos", "motor_cc", "potencia"]
4     for df in df_list:
5         for column in fill_columns:
6             if column == "asientos":
7                 df[column].fillna(value=df[column].median(), inplace=
8                     True)
9             else:
10                df[column].fillna(
11                    value=df[column].str.extract("(\\d+)").mean(),
12                    inplace=True
13                )
14                df.drop(columns=drop_columns, inplace=True)
15                df.dropna(inplace=True)
16     return df_list
```

Valores no numéricos

Ciertas columnas contienen valores alfanúmericos, aunque se nos proporcionan distintos archivos CSV para realizar un *mapping*. En este caso, utilizamos un **LabelEncoder**, y como entrada le damos el CSV correspondiente.

Es primordial usar el mismo *LabelEncoder* para los datos de entrenamiento como de test. La implementación se encuentra en la siguiente función:

```
1 def encode_columns(df_list):
2     label_encoder = LabelEncoder()
3     files = [
4         "ao",
5         "asientos",
6         "combustible",
7         "consumo",
8         "kilometros",
9         "mano",
10        "motor_cc",
11        "nombre",
12        "potencia",
13    ]
14    for data in files:
15        for df in df_list:
16            label = label_encoder.fit(read_csv("data/" + data + ".csv",
17                                             squeeze=True))
18            if data == "ao":
19                df["año"] = label.transform(df["año"])
20            else:
21                df[data] = label.transform(df[data])
22    return df_list
```

Balanceo de clases

Observamos que la mayoría de coches son de la categoría de precio 3, lo cual no es idóneo para entrenar un modelo de inteligencia artificial.

Debemos realizar un balanceo de las clases, en este caso optamos por usar el modelo **SMOTEENN**, que combina un *over-sampling* mediante **SMOTE** y una limpieza gracias a *Edited Nearest Neighbours (ENN)*.

La implementación se encuentra en esta función:

```
1 def balance_training_data(df):
2     smote_enn = SMOTEENN(random_state=42)
3     data, target = split_data_target(df=df, dataset="data")
4     balanced_data, balanced_target = smote_enn.fit_resample(data,
5                                                            target)
6     balanced_data_df = DataFrame(
7         balanced_data, columns=df.columns.difference(["precio_cat"])
8     )
9     balanced_target_df = DataFrame(balanced_target, columns=["
10        precio_cat"])
11    return balanced_data_df, balanced_target_df
```

Elección de algoritmo

Elegimos el algoritmo **GradientBoostingClassifier**, que pertenece a los algoritmos de *ensemble*. Éstos combinan las predicciones de varios clasificadores, con el objetivo de mejorar la generalización y la robustez de las predicciones.

En particular, pertenece a la familia de *boosting methods*, cuya característica es que los clasificadores se crean de forma secuencial, y uno de ellos trata de reducir el sesgo de los demás.

Resultados obtenidos

Al ejecutar el programa en local obtenemos los siguientes resultados:

```
Accuracy Score: 0.9177745580991017
Cross validation score: 0.9167671441999843
P3 master
```

Figura 1: Resultados de ejecución

Desafortunadamente, en la plataforma Kaggle obtenemos unos resultados pésimos:


44	—	AminKasrouAouamCeuta		0.31578	3	11d
----	---	----------------------	---	---------	---	-----

Figura 2: Resultados de Kaggle

Análisis de resultados

Debido a la discrepancia entre los resultados de la ejecución en local, y de la plataforma Kaggle, intuimos que debe de haber un problema en el preprocesamiento de datos.

También es posible que el modelo no sea óptimo para la tarea, aunque no justificaría un rendimiento tan bajo, puede contribuir a ello.