
Práctica 1

Inteligencia de Negocio

Amin Kasrou Aouam



**UNIVERSIDAD
DE GRANADA**

2020-11-10

Índice

Práctica 1	3
Introducción	3
Procesado de datos	3
Valores nulos	5
Valores no numéricos	6
Separación de datos	7

Práctica 1

Introducción

En esta práctica, usaremos distintos algoritmos de aprendizaje automático para resolver un problema de clasificación.

Procesado de datos

Antes de proceder con el entrenamiento de los distintos modelos, debemos realizar un preprocesado de los datos, para asegurarnos que nuestros modelos aprenden de un *dataset* congruente.

La integridad de la lógica del preprocesado se encuentra en el archivo *preprocessing.py*, cuyo contenido mostramos aquí:

```
1 from pandas import read_csv
2 from sklearn.preprocessing import LabelEncoder
3 from sklearn.model_selection import KFold
4
5
6 def replace_values(df):
7     columns = ["BI-RADS", "Margin", "Density", "Age"]
8     for column in columns:
9         df[column].fillna(value=df[column].mean(), inplace=True)
10    return df
11
12
13 def process_na(df, action):
14     if action == "drop":
15         return df.dropna()
16     elif action == "fill":
17         return replace_values(df)
18     else:
19         print("Unknown action selected. The choices are: ")
20         print("fill: fills the na values with the mean")
21         print("drop: drops the na values")
22         exit()
23
24
25 def encode_columns(df):
26     label_encoder = LabelEncoder()
27     encoded_df = df.copy()
28     encoded_df["Shape"] = label_encoder.fit_transform(df["Shape"])
29     encoded_df["Severity"] = label_encoder.fit_transform(df["Severity"])
30    return encoded_df
31
32
33 def split_train_target(df):
34     train_data = df.drop(columns=["Severity"])
35     target_data = df["Severity"]
36    return train_data, target_data
37
38
39 def split_k_sets(df):
40     k_fold = KFold(shuffle=True, random_state=42)
41    return k_fold.split(df)
42
43
44 def parse_data(source, action):
45     df = read_csv(filepath_or_buffer=source, na_values="?")
46     processed_df = process_na(df=df, action=action)
47     encoded_df = encode_columns(df=processed_df)
48     test_data, target_data = split_train_target(df=encoded_df)
49    return test_data, target_data
```

A continuación, mostraremos cada uno de los pasos que realizamos para obtener el *dataset* final:

Valores nulos

Nuestro *dataset* contiene valores nulos, representados mediante un signo de interrogación (?). Optamos por evaluar 2 estrategias:

1. Eliminar los valores nulos

```
1 df = read_csv(filepath_or_buffer="../data/mamografia.csv",
2               na_values="?")
3 processed_df = process_na(df=df, action="drop")
4 print("DataFrame sin preprocesamiento: ")
5 print(df.describe())
6 print("DataFrame sin preprocesamiento: ")
7 print(processed_df.describe())
```

```
1 DataFrame sin preprocesamiento:
2      BI-RADS      Age      Margin      Density
3 count  959.000000  956.000000  913.000000  885.000000
4 mean    4.296142   55.487448   2.796276   2.910734
5 std     0.706291   14.480131   1.566546   0.380444
6 min     0.000000   18.000000   1.000000   1.000000
7 25%     4.000000   45.000000   1.000000   3.000000
8 50%     4.000000   57.000000   3.000000   3.000000
9 75%     5.000000   66.000000   4.000000   3.000000
10 max    6.000000   96.000000   5.000000   4.000000
11 DataFrame sin preprocesamiento:
12     BI-RADS      Age      Margin      Density
13 count  847.000000  847.000000  847.000000  847.000000
14 mean    4.322314   55.842975   2.833530   2.909091
15 std     0.703762   14.603754   1.564049   0.370292
16 min     0.000000   18.000000   1.000000   1.000000
17 25%     4.000000   46.000000   1.000000   3.000000
18 50%     4.000000   57.000000   3.000000   3.000000
19 75%     5.000000   66.000000   4.000000   3.000000
20 max    6.000000   96.000000   5.000000   4.000000
```

Observamos que el número de instancias disminuye considerablemente, hasta un máximo de 112, en el caso del *BI-RADS*. Aún así, los valores de la media y desviación estándar no se ven afectados de forma considerable.

2. Imputar su valor con la media

```
1 df = read_csv(filepath_or_buffer="../data/mamografia.csv",
  na_values="?")
2 processed_df = process_na(df=df, action="fill")
3 print("DataFrame sin preprocesamiento: ")
4 print(df.describe())
5 print("DataFrame sin preprocesamiento: ")
6 print(processed_df.describe())
```

```
1 DataFrame sin preprocesamiento:
2      BI-RADS      Age      Margin      Density
3 count  961.000000  961.000000  961.000000  961.000000
4 mean    4.296142   55.487448   2.796276   2.910734
5 std     0.705555   14.442373   1.526880   0.365074
6 min     0.000000   18.000000   1.000000   1.000000
7 25%     4.000000   45.000000   1.000000   3.000000
8 50%     4.000000   57.000000   3.000000   3.000000
9 75%     5.000000   66.000000   4.000000   3.000000
10 max    6.000000   96.000000   5.000000   4.000000
11 DataFrame sin preprocesamiento:
12      BI-RADS      Age      Margin      Density
13 count  961.000000  961.000000  961.000000  961.000000
14 mean    4.296142   55.487448   2.796276   2.910734
15 std     0.705555   14.442373   1.526880   0.365074
16 min     0.000000   18.000000   1.000000   1.000000
17 25%     4.000000   45.000000   1.000000   3.000000
18 50%     4.000000   57.000000   3.000000   3.000000
19 75%     5.000000   66.000000   4.000000   3.000000
20 max    6.000000   96.000000   5.000000   4.000000
```

Esta alternativa nos permite mantener el número de instancias en todas las columnas, sin alterar la media ni la desviación típica.

Valores no numéricos

La mayoría de algoritmos de aprendizaje automática trabaja con datos numéricos, desafortunadamente nuestro *dataset* contiene dos columnas con datos descriptivos.

Procedemos a convertirlos en valores numéricos mediante un *LabelEncoder*:

```
1 encoded_df = encode_columns(df=processed_df)
2 print(encoded_df.head())
```

	BI-RADS	Age	Shape	Margin	Density	Severity	
1	0	5.0	67.0	1	5.0	3.000000	1
2	1	4.0	43.0	4	1.0	2.910734	1
3	2	5.0	58.0	0	5.0	3.000000	1
4	3	4.0	28.0	4	1.0	3.000000	0
5	4	5.0	74.0	4	5.0	2.910734	1

Vemos como las columnas **Shape** y **Severity** se componen ahora únicamente de valores numéricos.

Separación de datos

Como último paso, separamos la columna objetivo de los demás datos.

```
1 test_data, target_data = split_train_target(df=encoded_df)
2 print("Datos de entrenamiento: ")
3 print(test_data.head())
4 print("Datos objetivo: ")
5 print(target_data.head())
```

```
1 Datos de entrenamiento:
2   BI-RADS  Age  Shape  Margin  Density
3  0     5.0  67.0     1     5.0  3.000000
4  1     4.0  43.0     4     1.0  2.910734
5  2     5.0  58.0     0     5.0  3.000000
6  3     4.0  28.0     4     1.0  3.000000
7  4     5.0  74.0     4     5.0  2.910734
8 Datos objetivo:
9  0     1
10 1     1
11 2     1
12 3     0
13 4     1
14 Name: Severity, dtype: int64
```