

Act12

Daniel Rojas

March 2025

1 Introducción

Un árbol de decisión es un algoritmo de aprendizaje supervisado perteneciente al aprendizaje profundo usado para realizar regresiones o clasificaciones.

En resumen, lo que hace es crear un nodo raíz el cuál se bifurca en dos opciones; dependiendo del valor del input, decidirá a cual de las dos hojas dirigirse. Cada nodo se dividará también en dos caminos posibles hasta llegar a los nodos finales, hojas, lo cuales darán una categoría o valor como salida dependiendo del input.

Los árboles de decisión son importantes ya que permiten realizar predicciones útiles en la gran mayoría de casos y con un bajo coste computacional. Además, el algoritmo encuentra la mejor configuración del árbol de decisión por si mismo, lo cual permite crear árboles que procesen muchos datos de entrada, sin necesidad de que un humano haga la optimización.

2 Metodología

Para realizar la práctica se tomo como referencia el artículo de Árbol de decisión del libro Aprende Machine Learning, adjuntado al final en las referencias.

La actividad consiste en predecir si la canción de un artista entrará al top de los billboard, dependiendo de diferentes factores como su tempo, edad del artista al momento del lanzamiento, género, etc. Para ello, primero se tuvo que crear una nueva columna con la edad del artista al llegar al billboard, usando las columnas de fecha de nacimiento y fecha de ingreso al billboard. Sin embargo, como las columnas de fecha de nacimiento estaban vacías, se tuvo que asignar valores aleatorios a dichos valores.

Todo el proceso se puede observar en los siguientes fragmentos de código:

```
1 #Librerías usadas
2 import numpy as np
```

```

3 import pandas as pd
4 from sklearn import tree
5 from sklearn.metrics import accuracy_score
6 from sklearn.model_selection import KFold
7 from sklearn.model_selection import cross_val_score
8 from IPython.display import Image as PImage
9 from subprocess import check_call
10 from PIL import Image, ImageDraw, ImageFont

1 #Cargar el dataset
2 artists_billboard = pd.read_csv('artists_billboard_fix3.csv')

1 #convertir valores 0 a None
2 def edad_fix(row):
3     if row['anioNacimiento']==0:
4         return None
5     return row['anioNacimiento']
6
7 artists_billboard['anioNacimiento']=artists_billboard.apply(
    ↪ edad_fix, axis=1)

1 #Nueva columna con edad del artista al llegar al billboard
2 def calcula_edad(row):
3     cad = str(row['chart_date'])
4     year = row['anioNacimiento']
5     momento = cad[:4]
6     if year==0.0:
7         return None
8     return int(momento) - year
9
10 artists_billboard['edad_en_billboard']=artists_billboard.apply(
    ↪ calcula_edad, axis=1)

1 #Agregar valores aleatorios a los valores None de
2 age_avg = artists_billboard['edad_en_billboard'].mean()
3 age_std = artists_billboard['edad_en_billboard'].std()
4 age_null_count = artists_billboard['edad_en_billboard'].isnull().
    ↪ sum()
5 age_null_random_list = np.random.randint(age_avg - age_std, age_avg
    ↪ + age_std, size=age_null_count)
6
7 conValoresNulos = np.isnan(artists_billboard['edad_en_billboard'])
8
9 artists_billboard.loc[np.isnan(artists_billboard['edad_en_billboard
    ↪ ']), 'edad_en_billboard'] = age_null_random_list
10 artists_billboard['edad_en_billboard'] = artists_billboard['
    ↪ edad_en_billboard'].astype(int)

```

Ahora, hay que crear distintas categorías de datos para pasarlas como argumentos a nuestro árbol. Esto se realiza en las siguientes líneas de código:

```

1 # Mood Mapping
2 artists_billboard['moodEncoded'] = artists_billboard['mood'].map( {
    ↪ 'Energizing': 6,
3                                     'Empowering': 6,

```

```

4         'Cool': 5,
5         'Yearning': 4, # anhelo,
6             ↳ deseo, ansia
7         'Excited': 5, #emocionado
8         'Defiant': 3,
9         'Sensual': 2,
10        'Gritty': 3, #coraje
11        'Sophisticated': 4,
12        'Aggressive': 4, #
13            ↳ provocativo
14        'Fiery': 4, #caracter
15            ↳ fuerte
16        'Urgent': 3,
17        'Rowdy': 4, #ruidoso
18            ↳ alboroto
19        'Sentimental': 4,
20        'Easygoing': 1, # sencillo
21        'Melancholy': 4,
22        'Romantic': 2,
23        'Peaceful': 1,
24        'Brooding': 4, #
25            ↳ melancolico
26        'Upbeat': 5, #optimista
27            ↳ alegre
28        'Stirring': 5, #emocionante
29        'Lively': 5, #animado
30        'Other': 0, '' : 0} ).astype(
31            ↳ int)
32
33 # Tempo Mapping
34 artists_billboard['tempoEncoded'] = artists_billboard['tempo'].map(
35     ↳ {'Fast Tempo': 0, 'Medium Tempo': 2, 'Slow Tempo': 1, '' : 0}
36     ↳ ).astype(int)
37
38 # Genre Mapping
39 artists_billboard['genreEncoded'] = artists_billboard['genre'].map(
40     ↳ {'Urban': 4,
41
42         'Pop': 3,
43         'Traditional': 2,
44         'Alternative & Punk': 1,
45         'Electronica': 1,
46         'Rock': 1,
47         'Soundtrack': 0,
48         'Jazz': 0,
49         'Other': 0, '' : 0}
50     ).astype(int)
51
52 # artist_type Mapping
53 artists_billboard['artist_typeEncoded'] = artists_billboard['
54     ↳ artist_type'].map( {'Female': 2, 'Male': 3, 'Mixed': 1, '' :
55     ↳ 0} ).astype(int)
56
57 # Mapping edad en la que llegaron al billboard
58 artists_billboard.loc[ artists_billboard['edad_en_billboard'] <=
59     ↳ 21, 'edadEncoded'] = 0
60 artists_billboard.loc[(artists_billboard['edad_en_billboard'] > 21)
61     ↳ & (artists_billboard['edad_en_billboard'] <= 26), '
62     ↳ edadEncoded'] = 1
63 artists_billboard.loc[(artists_billboard['edad_en_billboard'] > 26)

```

```

    ↪ & (artists_billboard['edad_en_billboard'] <= 30), '
    ↪ edadEncoded'] = 2
46 artists_billboard.loc[(artists_billboard['edad_en_billboard'] > 30)
    ↪ & (artists_billboard['edad_en_billboard'] <= 40), '
    ↪ edadEncoded'] = 3
47 artists_billboard.loc[ artists_billboard['edad_en_billboard'] > 40,
    ↪ 'edadEncoded'] = 4
48
49 # Mapping Song Duration
50 artists_billboard.loc[ artists_billboard['durationSeg'] <= 150, '
    ↪ durationEncoded'] = 0
51 artists_billboard.loc[(artists_billboard['durationSeg'] > 150) & (
    ↪ artists_billboard['durationSeg'] <= 180), 'durationEncoded']
    ↪ = 1
52 artists_billboard.loc[(artists_billboard['durationSeg'] > 180) & (
    ↪ artists_billboard['durationSeg'] <= 210), 'durationEncoded']
    ↪ = 2
53 artists_billboard.loc[(artists_billboard['durationSeg'] > 210) & (
    ↪ artists_billboard['durationSeg'] <= 240), 'durationEncoded']
    ↪ = 3
54 artists_billboard.loc[(artists_billboard['durationSeg'] > 240) & (
    ↪ artists_billboard['durationSeg'] <= 270), 'durationEncoded']
    ↪ = 4
55 artists_billboard.loc[(artists_billboard['durationSeg'] > 270) & (
    ↪ artists_billboard['durationSeg'] <= 300), 'durationEncoded']
    ↪ = 5
56 artists_billboard.loc[ artists_billboard['durationSeg'] > 300, '
    ↪ durationEncoded'] = 6

1 #Quitar las columnas no categorizadas o innecesarias
2 drop_elements = ['id','title','artist','mood','tempo','genre','
    ↪ artist_type','chart_date','anioNacimiento','durationSeg','
    ↪ edad_en_billboard']
3 artists_encoded = artists_billboard.drop(drop_elements, axis = 1)

```

Ahora, para obtener el mejor árbol en función de la máxima profundidad permitida, se realizó un diseño de experimento provando entre árboles con 1 y 7 de profundidad para encontrar el mejor. Además, se usó la técnica de kFolds para suplementar la escasez de datos. Por último, se ajustan los pesos de cada etiqueta top a 1:3.5 en class, ya que hay 3.5 veces menos valores con el valor top = 1 que con el valor top = 0. (O sea es más común que no se llegue al top del billboard):

```

1 cv = KFold(n_splits=10) #9 folds para entrenar y 1 para validar
2 accuracies = list() #Almacena la precisión de cada árbol creado
3 max_attributes = len(list(artists_encoded)) #Número máximo de
    ↪ columnas en el dataset usado
4 depth_range = range(1, max_attributes + 1) #Diferentes niveles de
    ↪ profundidad de los árboles
5
6 # Testearemos la profundidad de 1 a cantidad de atributos +1
7 for depth in depth_range:
8     fold_accuracy = [] #Precisión de cada fold
9     tree_model = tree.DecisionTreeClassifier(criterion='entropy',
10                                              min_samples_split=20,
11                                              min_samples_leaf=5,

```

```

12                                     max_depth = depth,
13                                     class_weight={1:3.5})
14     for train_fold, valid_fold in cv.split(artists_encoded):
15         f_train = artists_encoded.loc[train_fold]
16         f_valid = artists_encoded.loc[valid_fold]
17
18         model = tree_model.fit(X = f_train.drop(['top'], axis=1),
19                               y = f_train["top"])
20         valid_acc = model.score(X = f_valid.drop(['top'], axis=1),
21                                y = f_valid["top"]) # calculamos la
22                                                         ↳ precision con el segmento de
23                                                         ↳ validacion
24         fold_accuracy.append(valid_acc)
25
26     avg = sum(fold_accuracy)/len(fold_accuracy)
27     accuracies.append(avg)
28
29 # Mostramos los resultados obtenidos
30 df = pd.DataFrame({"Max Depth": depth_range, "Average Accuracy":
31                    ↳ accuracies})
32 df = df[["Max Depth", "Average Accuracy"]]
33 print(df.to_string(index=False))

```

Al ejecutar lo anterior se obtuvo que el mejor caso era de árbol con profundidad máxima de 4. Los valores obtenidos son los siguientes:

1. 0.556101
2. 0.556126
3. 0.564038
4. 0.650446
5. 0.603199
6. 0.620511
7. 0.636210

Ahora, para observar visualmente la estructura del mejor árbol, podemos ejecutar el siguiente código:

```

1 # Crear arrays de entrenamiento y las etiquetas que indican si
2   ↳ lleg a top o no
3 y_train = artists_encoded['top']
4 x_train = artists_encoded.drop(['top'], axis=1).values
5
6 # Crear Arbol de decision con profundidad = 4
7 decision_tree = tree.DecisionTreeClassifier(criterion='entropy',
8                                              min_samples_split=20,
9                                              min_samples_leaf=5,
10                                              max_depth = 4,
11                                              class_weight={1:3.5})
12 decision_tree.fit(x_train, y_train)

```

```

13 # exportar el modelo a archivo .dot
14 with open(r"tree1.dot", 'w') as f:
15     f = tree.export_graphviz(decision_tree,
16                             out_file=f,
17                             max_depth = 4,
18                             impurity = True,
19                             feature_names = list(artists_encoded.
20                                             ↪ drop(['top'], axis=1)),
21                             class_names = ['No', 'N1 Billboard'],
22                             rounded = True,
23                             filled= True )
24 # Convertir el archivo .dot a png para poder visualizarlo
25 check_call(['dot', '-Tpng', r'tree1.dot', '-o', r'tree1.png'])
26 PImage("tree1.png")

```

La imagen resultante es la siguiente:

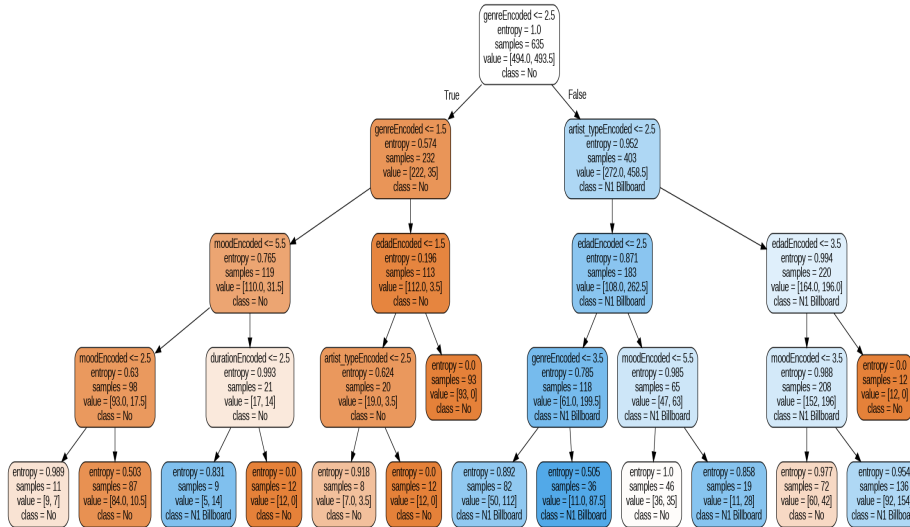


Figure 1: Árbol de 4 de profundidad

3 Resultados

Se puede observar que el mejor árbol es el que tiene una profundidad máxima de 4 niveles, con un porcentaje de precisión del 65.04%. Además, la imagen del árbol nos permite ver la forma en la que opera el árbol de decisión, observando que separa los posibles valores en función de si los valores de las features es mayor o menor a un cierto número.

4 Conclusiones

Con la actividad, aprendí de forma general el como se debe de utilizar apropiadamente los árboles de decisión, por ejemplo, el considerar que la proporción de de clases a predecir para cambiar los valores de los pesos. Por otro lado, pude observar de forma general el funcionamiento de los árboles de decisión por medio de la imagen generada. Por último, Aprendí algunos métodos de preprocesamiento útiles para casos específicos, como el uso de valores aleatorios alrededor de la media aritmética o convertir valores continuos a categóricos.

5 Referencias

Bagnato, J. I. (2020). Aprende Machine Learning.