

Act 12: Programando Árbol de Decisión en Python

Martín Alexis Martínez Andrade - 2049334

1. Introducción

Los árboles de decisión son algoritmos de aprendizaje supervisado que se utilizan tanto para tareas de clasificación como para regresión. Estos modelos representan gráficamente las decisiones mediante una estructura en forma de árbol, donde cada nodo interno corresponde a la evaluación de una condición y cada hoja representa una conclusión o respuesta final. En este ejercicio se usa scikit-learn para construir y visualizar un árbol de decisión que, mediante atributos de artistas y datos históricos del Billboard Hot 100, intenta predecir si un artista logrará alcanzar el puesto número 1 (top) en el ranking.

2. Metodología

La actividad se desarrolló siguiendo estos pasos:

1. Carga y preprocesamiento de datos:

- Se cargó el dataset `artists_billboard_fix3.csv` utilizando Pandas.
- Se revisaron los primeros registros y se imprimieron algunos valores estadísticos, así como el promedio y la desviación estándar de la edad de los artistas.
- Se trató el problema de valores nulos en la columna del año de nacimiento, sustituyéndolos por valores aleatorios dentro del intervalo (usando el promedio y el desvío).

2. Mapeo de variables:

- Se transformaron variables categóricas (como mood, tempo, género, tipo de artista) en valores numéricos.

- Se crearon nuevas columnas para `.edadEncoded` y `"durationEncoded"`.
- Se eliminó información innecesaria para formar el conjunto de datos final `artists_encoded`.

3. Construcción y validación del árbol de decisión:

- Se utilizó la función `KFold` para probar distintos niveles de profundidad (`max_depth` del árbol) y determinar el valor óptimo. Se realizaron pruebas para profundidades desde 1 hasta el número de atributos + 1.
- Se configuró el árbol de decisión y se usó el criterio de entropía para evaluar la calidad de las divisiones. Se configuró `class_weight` para compensar el desbalance en la variable de salida.
- Se ajustó el modelo sobre el subset de entrenamiento y se evaluó la precisión obtenida.

4. Visualización y predicción:

- Se visualizó el árbol con la función `plot_tree` de Matplotlib.
- Se realizaron pruebas de predicción con dos ejemplos: uno para Camila Cabello (alcanzó el puesto número 1) y otro para Imagine Dragons (no alcanzó el top).

Fragmento de Código en Python

```

1  # Author: Mart n Alexis Mart nez Andrade - 2049334
2
3  # Imports necesarios
4  import numpy as np
5  import pandas as pd
6  import seaborn as sb
7  import matplotlib.pyplot as plt
8  plt.rcParams['figure.figsize'] = (16, 9)
9  plt.style.use('ggplot')
10 from sklearn import tree
11 from sklearn.metrics import accuracy_score
12 from sklearn.model_selection import KFold
13 from sklearn.model_selection import cross_val_score
14 from IPython.display import Image as PImage
15 from subprocess import check_call
16 from PIL import Image, ImageDraw, ImageFont
17

```

```

18 # Datos de entrada
19 artists_billboard = pd.read_csv(r"artists_billboard_fix3.csv"
20 )
21 print(artists_billboard.head())
22
23 f1 = artists_billboard['chart_date'].values
24 f2 = artists_billboard['durationSeg'].values
25
26 colores = ['orange', 'blue'] # si no estaban declarados
27     previamente
28 tamanios = [60, 40] # si no estaban declarados previamente
29
30 asignar = []
31 asignar2 = []
32 for index, row in artists_billboard.iterrows():
33     asignar.append(colores[row['top']])
34     asignar2.append(tamanios[row['top']])
35
36 def edad_fix(anio):
37     if anio==0:
38         return None
39     return anio
40
41 artists_billboard['anioNacimiento']=artists_billboard.apply(
42     lambda x: edad_fix(x['anioNacimiento']), axis=1)
43
44 def calcula_edad(anio, cuando):
45     cad = str(cuando)
46     momento = cad[:4]
47     if anio==0.0:
48         return None
49     return int(momento) - anio
50
51 artists_billboard['edad_en_billboard']=artists_billboard.
52     apply(lambda x: calcula_edad(x['anioNacimiento'],x['
53     chart_date'])), axis=1);
54
55 age_avg = artists_billboard['edad_en_billboard'].mean()
56 age_std = artists_billboard['edad_en_billboard'].std()
57 age_null_count = artists_billboard['edad_en_billboard'].
58     isnull().sum()
59 age_null_random_list = np.random.randint(age_avg - age_std,
60     age_avg + age_std, size=age_null_count)
61 conValoresNulos = np.isnan(artists_billboard['
62     edad_en_billboard'])
63 artists_billboard.loc[np.isnan(artists_billboard['
64     edad_en_billboard']), 'edad_en_billboard'] =
65     age_null_random_list
66 artists_billboard['edad_en_billboard'] = artists_billboard['

```

```

    edad_en_billboard'].astype(int)
57 print("Edad Promedio: " + str(age_avg))
58 print("Desvi Std Edad: " + str(age_std))
59 print("Intervalo para asignar edad aleatoria: " + str(int(
    age_avg - age_std)) + " a " + str(int(age_avg + age_std)))
60
61
62 # Mood Mapping
63 artists_billboard['moodEncoded'] = artists_billboard['mood'].
    map( {'Energizing': 6,
64         'Empowering': 6,
65         'Cool': 5,
66         'Yearning': 4, # anhelo, deseo, ansia
67         'Excited': 5, #emocionado
68         'Defiant': 3,
69         'Sensual': 2,
70         'Gritty': 3, #coraje
71         'Sophisticated': 4,
72         'Aggressive': 4, # provocativo
73         'Fiery': 4, #caracter fuerte
74         'Urgent': 3,
75         'Rowdy': 4, #ruidoso alboroto
76         'Sentimental': 4,
77         'Easygoing': 1, # sencillo
78         'Melancholy': 4,
79         'Romantic': 2,
80         'Peaceful': 1,
81         'Brooding': 4, # melancolico
82         'Upbeat': 5, #optimista alegre
83         'Stirring': 5, #emocionante
84         'Lively': 5, #animado
85         'Other': 0,':':0} ).astype(int)
86
87 # Tempo Mapping
88 artists_billboard['tempoEncoded'] = artists_billboard['tempo'
    ].map( {'Fast Tempo': 0,
89

```

, Medium

Tempo
,

:

2,

,

Slow

Tempo

```

,
:
1,
,
,
:
0}
)
.
astype
(
int
)

90
91 # Genre Mapping
92 artists_billboard['genreEncoded'] = artists_billboard['genre'
93     ].map( {'Urban': 4,
94             'Pop': 3,
95             'Traditional': 2,
96             'Alternative & Punk': 1,
97             'Electronica': 1,
98             'Rock': 1,
99             'Soundtrack': 0,
100            'Jazz': 0,
101            'Other': 0, '' : 0}
102            ).astype(int)
103
104 # artist_type Mapping
105 artists_billboard['artist_typeEncoded'] = artists_billboard['
106     artist_type'].map(
107         {'Female': 2, 'Male': 3, 'Mixed': 1, '' : 0}).astype(int)
108
109 # Mapping edad en la que llegaron al billboard
110 artists_billboard.loc[ artists_billboard['edad_en_billboard']
111     <= 21, 'edadEncoded'] = 0
112 artists_billboard.loc[(artists_billboard['edad_en_billboard']
113     > 21) & (artists_billboard['edad_en_billboard'] <= 26), '
114     edadEncoded'] = 1
115 artists_billboard.loc[(artists_billboard['edad_en_billboard']
116     > 26) & (artists_billboard['edad_en_billboard'] <= 30), '
117     edadEncoded'] = 2
118 artists_billboard.loc[(artists_billboard['edad_en_billboard']
119     > 30) & (artists_billboard['edad_en_billboard'] <= 40), '
120     edadEncoded'] = 3

```

```

112 artists_billboard.loc[ artists_billboard['edad_en_billboard']
    > 40, 'edadEncoded'] = 4
113 # Mapping Song Duration
114 artists_billboard.loc[ artists_billboard['durationSeg'] <=
    150, 'durationEncoded'] = 0
115 artists_billboard.loc[(artists_billboard['durationSeg'] >
    150) & (artists_billboard[\
116 'durationSeg'] <= 180), 'durationEncoded'] = 1
117 artists_billboard.loc[(artists_billboard['durationSeg'] >
    180) & (artists_billboard[\
118 'durationSeg'] <= 210), 'durationEncoded'] = 2
119 artists_billboard.loc[(artists_billboard['durationSeg'] >
    210) & (artists_billboard[\
120 'durationSeg'] <= 240), 'durationEncoded'] = 3
121 artists_billboard.loc[(artists_billboard['durationSeg'] >
    240) & (artists_billboard[\
122 'durationSeg'] <= 270), 'durationEncoded'] = 4
123 artists_billboard.loc[(artists_billboard['durationSeg'] >
    270) & (artists_billboard[\
124 'durationSeg'] <= 300), 'durationEncoded'] = 5
125 artists_billboard.loc[ artists_billboard['durationSeg'] >
    300, 'durationEncoded'] = 6
126
127 drop_elements = ['id', 'title', 'artist', 'mood', 'tempo', 'genre',
    'artist_type', 'chart_date', 'anioNacimiento', 'durationSeg',
    'edad_en_billboard']
128 artists_encoded = artists_billboard.drop(drop_elements, axis
    = 1)
129
130 print(artists_encoded[['moodEncoded', 'top']].groupby(['
    moodEncoded'], as_index=False).agg(['mean', 'count', 'sum'
    ]))
131 print(artists_encoded[['artist_typeEncoded', 'top']].groupby(
    (['artist_typeEncoded'], as_index=False).agg(['mean', '
    count', 'sum'])))
132 print(artists_encoded[['genreEncoded', 'top']].groupby(['
    genreEncoded'], as_index=False).agg(['mean', 'count', 'sum'
    ]))
133 print(artists_encoded[['tempoEncoded', 'top']].groupby(['
    tempoEncoded'], as_index=False).agg(['mean', 'count', 'sum'
    ]))
134 print(artists_encoded[['durationEncoded', 'top']].groupby(['
    durationEncoded'], as_index=False).agg(['mean', 'count', '
    sum'])))
135 print(artists_encoded[['edadEncoded', 'top']].groupby(['
    edadEncoded'], as_index=False).agg(['mean', 'count', 'sum'
    ]))
136
137 cv = KFold(n_splits=10)

```

```

138 accuracies = list()
139 max_attributes = len(list(artists_encoded))
140 depth_range = range(1, max_attributes + 1)
141
142 for depth in depth_range:
143     fold_accuracy = []
144     tree_model = tree.DecisionTreeClassifier(criterion='
entropy',
145         min_samples_split=20,
146         min_samples_leaf=5,
147         max_depth = depth,
148         class_weight={1:3.5}
149     )
150
151     for train_fold, valid_fold in cv.split(artists_encoded):
152         f_train = artists_encoded.loc[train_fold]
153         f_valid = artists_encoded.loc[valid_fold]
154         model = tree_model.fit(X = f_train.drop(['top'], axis
=1),
155             y = f_train["top"])
156         valid_acc = model.score(X = f_valid.drop(['top'],
axis=1),
157             y = f_valid["top"])
158         fold_accuracy.append(valid_acc)
159
160     avg = sum(fold_accuracy)/len(fold_accuracy)
161     accuracies.append(avg)
162 # Mostramos los resultados obtenidos
163 df = pd.DataFrame({"Max Depth": depth_range, "Average
Accuracy": accuracies})
164 df = df[["Max Depth", "Average Accuracy"]]
165 print(df.to_string(index=False))
166
167 # Crear arrays de entrenamiento y las etiquetas que indican
si lleg a top o no
168 y_train = artists_encoded['top']
169 x_train = artists_encoded.drop(['top'], axis=1).values
170 # Crear Arbol de decision con profundidad = 4
171 decision_tree = tree.DecisionTreeClassifier(criterion='
entropy',
172     min_samples_split=20,
173     min_samples_leaf=5,
174     max_depth = 4,
175     class_weight={1:3.5})
176
177 decision_tree.fit(x_train, y_train)
178
179 # exportar el modelo a archivo .dot
180 with open(r"tree1.dot", 'w') as f:

```

```

181     f = tree.export_graphviz(decision_tree,
182                             out_file=f,
183                             max_depth = 7,
184                             impurity = True,
185                             feature_names = list(artists_encoded.drop(['top'],
186                             axis=1)),
187                             class_names = ['No', 'N1 Billboard'],
188                             rounded = True,
189                             filled= True,
190                             )
191 # Visualizar sin Graphviz
192 import matplotlib.pyplot as plt
193 plt.figure(figsize=(12, 8))
194 tree.plot_tree(decision_tree,
195               feature_names=list(artists_encoded.drop(['top'], axis=1).
196               columns),
197               class_names=['No', 'N1 Billboard'],
198               filled=True,
199               rounded=True,
200               impurity=True)
201 plt.show()
202
203 acc_decision_tree = round(decision_tree.score(x_train,
204               y_train) * 100, 2)
205 print(acc_decision_tree)
206
207 #predecir artista CAMILA CABELLO featuring YOUNG THUG
208 # con su canci n Havana llego a numero 1 Billboard US en
209 2017
210 x_test = pd.DataFrame(columns=('top', 'moodEncoded', '
211               tempoEncoded', 'genreEncoded', 'artist_typeEncoded', '
212               edadEncoded', 'durationEncoded'))
213 x_test.loc[0] = (1,5,2,4,1,0,3)
214 y_pred = decision_tree.predict(x_test.drop(['top'], axis = 1)
215 )
216 print("Prediccion: " + str(y_pred))
217 y_proba = decision_tree.predict_proba(x_test.drop(['top'],
218               axis = 1))
219 print("Probabilidad de Acierto: " + str(round(y_proba[0][
220               y_pred[0]]* 100, 2)) + "%")
221
222 #predecir artista Imagine Dragons
223 # con su canci n Believer llego al puesto 42 Billboard US en
224 2017
225 x_test = pd.DataFrame(columns=('top', 'moodEncoded', '
226               tempoEncoded', 'genreEncoded', 'artist_typeEncoded', '
227               edadEncoded', 'durationEncoded'))
228 x_test.loc[0] = (0,4,2,1,3,2,3)

```



```

218 y_pred = decision_tree.predict(x_test.drop(['top'], axis = 1)
    )
219 print("Prediccion: " + str(y_pred))
220 y_proba = decision_tree.predict_proba(x_test.drop(['top'],
    axis = 1))
221 print("Probabilidad de Acierto: " + str(round(y_proba[0][
    y_pred[0]]* 100, 2)) + "%")

```

Listing 1: Árbol de Decisión en Python

3. Resultados

- Tras optimizar la profundidad del árbol mediante cross validation, se obtuvo que una profundidad de 4 generaba la mejor precisión (alrededor del 64 % en entrenamiento).
- La visualización del árbol permitió observar que el nodo raíz realiza una división, por ejemplo, basada en el género, y que en niveles posteriores intervienen otras características como edad y duración.
- La precisión global del árbol, una vez ajustado, fue del 68.19 %.
- Para el caso de predicción de ejemplos:
 - Se predijo que Camila Cabello (con la canción "Havana") alcanzaría el top 1, con una probabilidad aproximada del 71.43 %.
 - Se predijo que Imagine Dragons (con "Believer") NO alcanzaría el top 1, con una probabilidad de acierto del 88.89 %.

4. Conclusión

En esta actividad se implementó un árbol de decisión utilizando Python y scikit-learn para predecir si un artista alcanzará el puesto número 1 en el Billboard Hot 100. Se procesaron y transformaron los datos originales, se realizó un mapeo de variables categóricas y se optimizó la profundidad del árbol usando cross validation. Aunque la precisión final del modelo no es muy alta, los experimentos con predicción de ejemplos nos permitieron evaluar el comportamiento del algoritmo frente a nuevos casos.

```
C:\AlexisAndradeDev\UANL\Inteligencia Artificial\Códigos\12>python main.py
```

```
   id      title  ... top anioNacimiento
0    0  Small Town Throwdown  ...  0      1975.0
1    1          Bang Bang  ...  0      1989.0
2    2          Timber  ...  1      1993.0
3    3      Sweater Weather  ...  0      1989.0
4    4      Automatic  ...  0      0.0
```

```
[5 rows x 11 columns]
```

```
Edad Promedio: 30.10282258064516
```

```
Desvió Std Edad: 8.40078832861513
```

```
Intervalo para asignar edad aleatoria: 21 a 38
```

```
   moodEncoded      top
```

```
              mean count sum
```

```
0          0  0.000000      1  0
```

```
1          1  0.000000      8  0
```

```
2          2  0.274194     62 17
```

```
3          3  0.145631    103 15
```

```
4          4  0.136986    146 20
```

```
5          5  0.294872    156 46
```

```
6          6  0.270440    159 43
```

```
   artist_typeEncoded      top
```

```
                  mean count sum
```

```
0                  1  0.305263    95 29
```

```
1                  2  0.320261   153 49
```

```
2                  3  0.162791   387 63
```

```
   genreEncoded      top
```

```
              mean count sum
```

```
0          0  0.105263     19  2
```

```
1          1  0.070000    100  7
```

```
2          2  0.008850    113  1
```

```
3          3  0.319149    188 60
```

```
4          4  0.330233    215 71
```

```
   tempoEncoded      top
```

```
              mean count sum
```

```
0          0  0.226415     53 12
```

```
1          1  0.246154     65 16
```

```
2          2  0.218569    517 113
```

```
   durationEncoded      top
```

```
              mean count sum
```

```
0          0.0  0.295775     71 21
```

```
1          1.0  0.333333     30 10
```

```
   edadEncoded      top
```

```
              mean count sum
```

```
0          0.0  0.246575     73 18
```

```
1          1.0  0.309677    155 48
```

```
2          2.0  0.256944    144 37
```

```
3          3.0  0.166607    216 36
```

```
4          4.0  0.042553     47  2
```

```
Max Depth Average Accuracy
```

```
1          0.556101
```

```
2          0.556126
```

```
3          0.564038
```

```
4          0.640923
```

```
5          0.600000
```

```
6          0.609524
```

```
7          0.643998
```

```
68.19
```

```
C:\Users\Alexi\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\utils\validation.py:2732: UserWarning: X has feature names, but DecisionTreeClassifier was fitted without fe
```

```
ature names
```

```
warnings.warn(
```

```
Prediccion: [1]
```

```
C:\Users\Alexi\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\utils\validation.py:2732: UserWarning: X has feature names, but DecisionTreeClassifier was fitted without fe
```

```
ature names
```

```
warnings.warn(
```

```
Probabilidad de Acierto: 71.43%
```

```
C:\Users\Alexi\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\utils\validation.py:2732: UserWarning: X has feature names, but DecisionTreeClassifier was fitted without fe
```

```
ature names
```

```
warnings.warn(
```

```
Prediccion: [0]
```

```
C:\Users\Alexi\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\utils\validation.py:2732: UserWarning: X has feature names, but DecisionTreeClassifier was fitted without fe
```

```
ature names
```

```
warnings.warn(
```

```
Probabilidad de Acierto: 88.89%
```

