

Actividad 12: Programando Árbol de Decisión en Python

Gerardo Enrique Torres Flores 2064063

30 de marzo de 2025

1. Introducción

Los *árboles de decisión* son representaciones gráficas de posibles soluciones a una decisión basadas en ciertas condiciones, es uno de los algoritmos de aprendizaje supervisado más utilizados en machine learning y pueden realizar tareas de clasificación o regresión.

Los árboles de decisión tienen un primer nodo llamado **raíz** (root) y luego se descomponen el resto de atributos de entrada en dos ramas planteando una condición que puede ser cierta o falsa. Se bifurca cada nodo en 2 y vuelven a subdividirse hasta llegar a las hojas que son los nodos finales y que equivalen a respuestas a la solución: **Si/No, Comprar/Vender, o lo que sea que estemos clasificando.**

En esta actividad se empleará un árbol de decisión para analizar un conjunto de datos musical (`artists_billboard_fix3.csv`), evaluar el desempeño del modelo y realizar predicciones sobre el rendimiento de artistas en el billboard.

2. Metodología

Para realizar esta actividad se siguieron los siguientes pasos:

1. Importación de librerías y configuración

Se importaron las librerías necesarias para el análisis, visualización y creación del modelo de Árbol de Decisión:

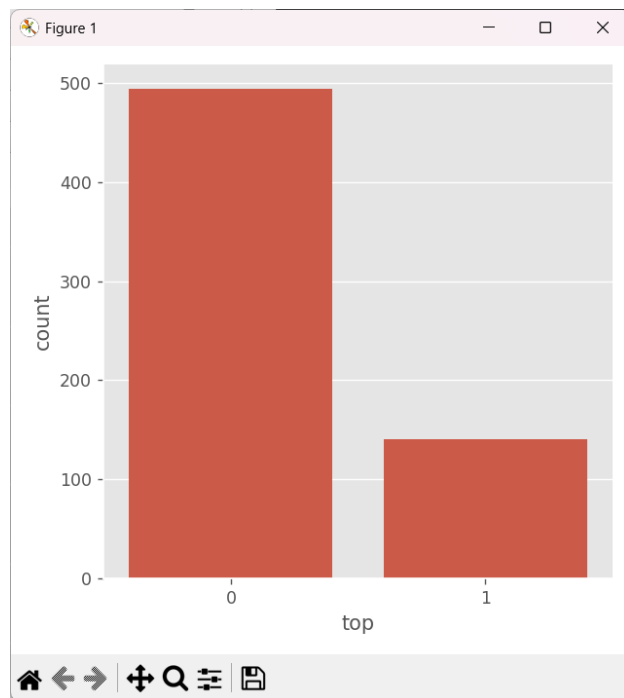
```
1 import numpy as np
2 import pandas as pd
3 import seaborn as sb
4 import matplotlib.pyplot as plt
5 plt.rcParams['figure.figsize'] = (16, 9)
6 plt.style.use('ggplot')
7 from sklearn import tree
8 from sklearn.metrics import accuracy_score
9 from sklearn.model_selection import KFold, cross_val_score
10 from IPython.display import Image as PImage
11 from subprocess import check_call
12 from PIL import Image, ImageDraw, ImageFont
```

2. Lectura y exploración de datos

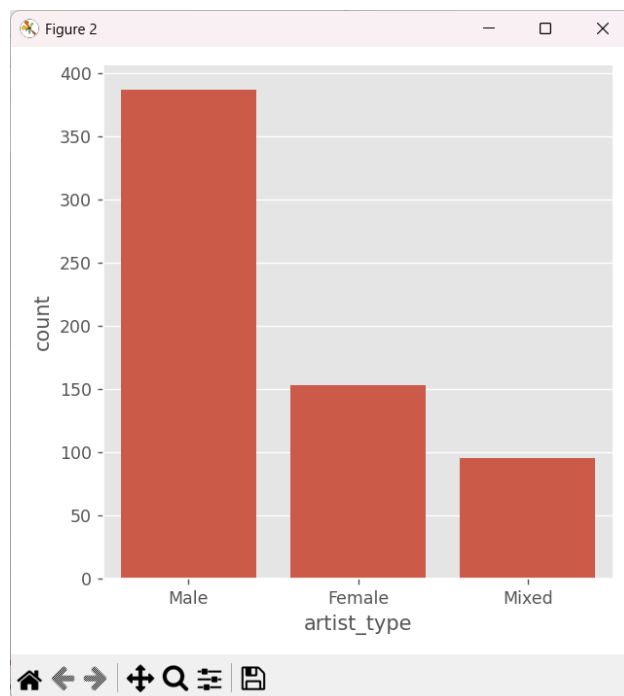
Se carga el conjunto de datos `artists_billboard_fix3.csv` y se exploran las distribuciones de las variables y las clases:

```
1 artists_billboard = pd.read_csv(r"artists_billboard_fix3.csv")
2 print(artists_billboard.groupby('top').size())
3
4 sb.catplot(x='top', data=artists_billboard, kind="count")
5 sb.catplot(x='artist_type', data=artists_billboard, kind="count")
6 sb.catplot(x='top', data=artists_billboard, hue='artist_type',
7           , kind="count")
8 sb.catplot(x='mood', data=artists_billboard, kind="count",
9           aspect=3)
10 sb.catplot(x='tempo', data=artists_billboard, hue='top', kind="count",
11           aspect=3)
12 sb.catplot(x='genre', data=artists_billboard, kind="count",
13           aspect=3)
14 sb.catplot(x='mood', data=artists_billboard, hue='top', kind="count",
15           aspect=3)
16 sb.catplot(x='anioNacimiento', data=artists_billboard, kind="count",
17           aspect=3)
```

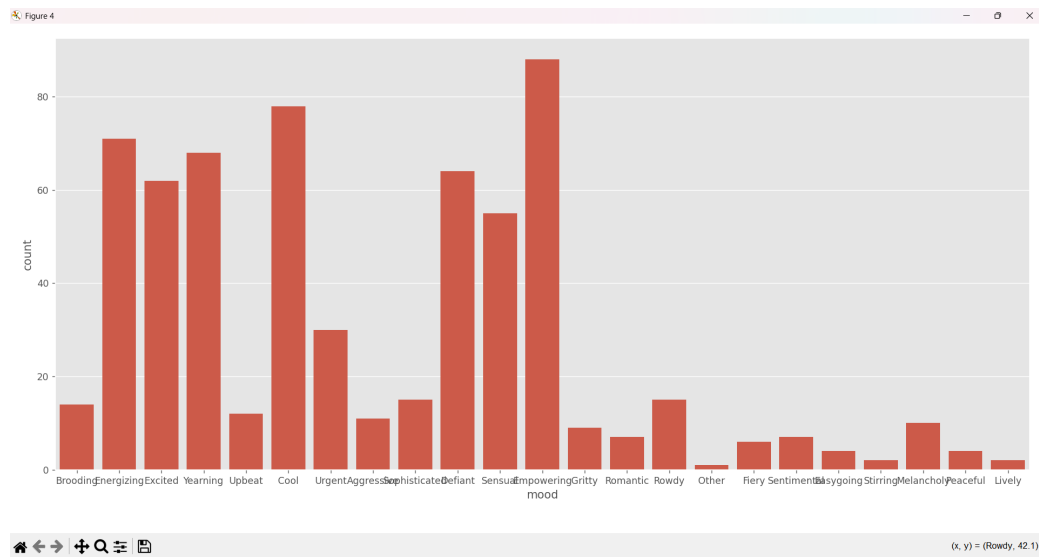
Visualización de cuántos estuvieron en Top 1 y no:



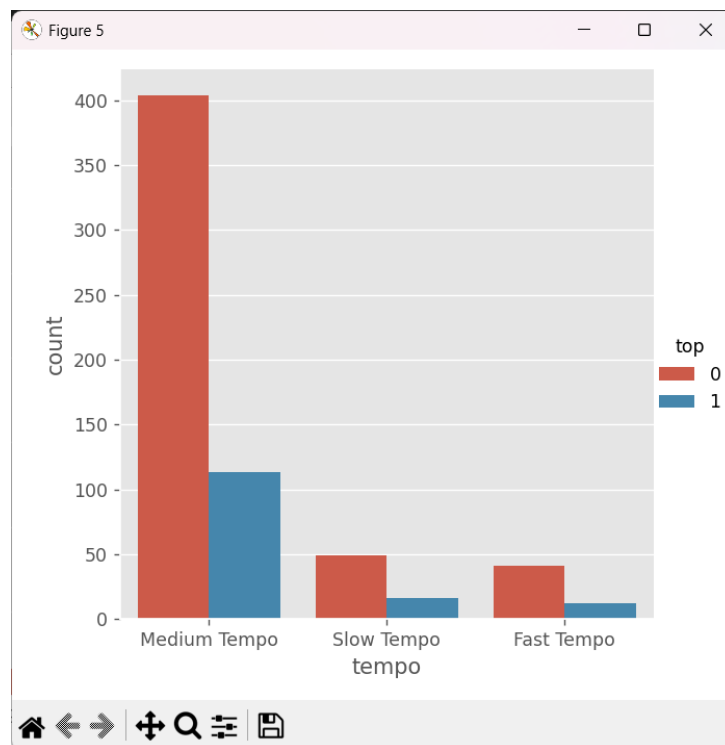
Visualización de Género del Artista(s):



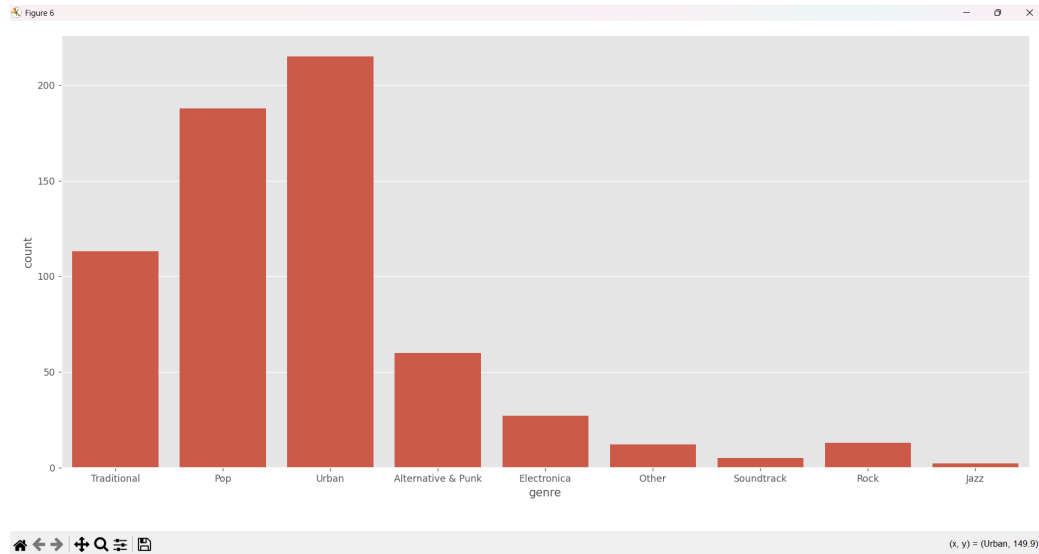
Visualización de Tipos de Mood:



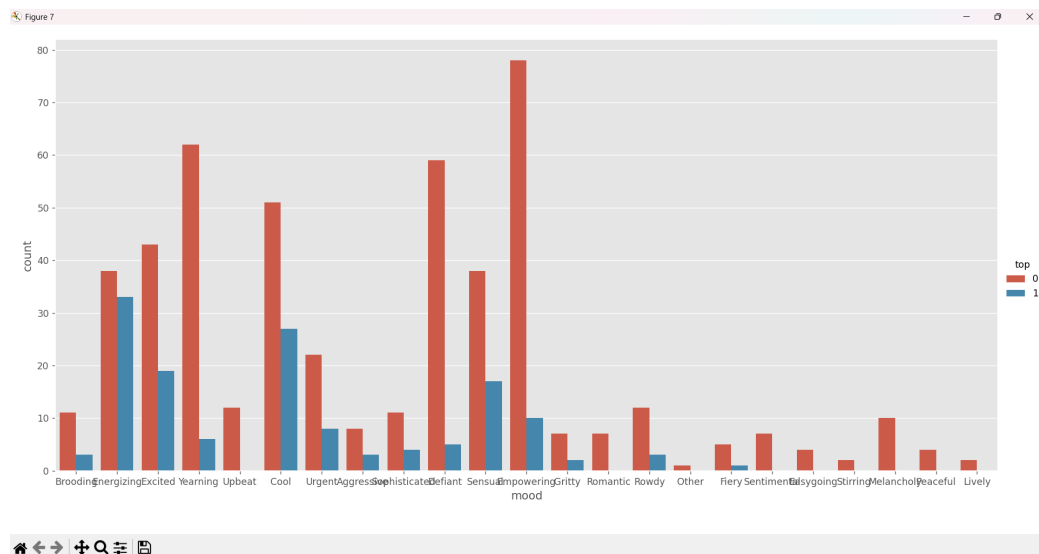
Visualización de Tempo de la Canción:



Visualización de Géneros Musicales:



Visualización de Géneros Musicales en Top 1 y no:



3. Visualización y preprocesamiento

Se realizaron varias visualizaciones para analizar relaciones entre variables (como duración de canción y año de nacimiento) y se realizó el preprocesamiento necesario (codificación de variables, manejo de datos faltantes, etc.):

```
1 # Visualizacion de relacion entre anio de nacimiento y
   duracion de cancion
2 colores = ['orange','blue']
3 tamanios = [60,40]
4 f1 = artists_billboard['anioNacimiento'].values
5 f2 = artists_billboard['durationSeg'].values
6
7 asignar = []
8 for index, row in artists_billboard.iterrows():
9     asignar.append(colores[row['top']])
10 plt.scatter(f1, f2, c=asignar, s=30)
11 plt.axis([1960,2005,0,600])
12
13 # Preprocesamiento: ajuste de la edad y codificacion de
   variables
14 def edad_fix(anio):
15     if anio == 0:
16         return None
17     return anio
18 artists_billboard['anioNacimiento'] = artists_billboard.apply
   (lambda x: edad_fix(x['anioNacimiento']), axis=1)
19 def calcula_edad(anio, cuando):
20     cad = str(cuando)
21     momento = cad[:4]
22     if anio == 0.0:
23         return None
24     return int(momento) - anio
25 artists_billboard['edad_en_billboard'] = artists_billboard.
   apply(lambda x: calcula_edad(x['anioNacimiento'], x['
   chart_date']), axis=1)
26
27 age_avg = artists_billboard['edad_en_billboard'].mean()
28 age_std = artists_billboard['edad_en_billboard'].std()
29 age_null_count = artists_billboard['edad_en_billboard'].
   isnull().sum()
30 age_null_random_list = np.random.randint(age_avg - age_std,
   age_avg + age_std, size=age_null_count)
31 artists_billboard.loc[np.isnan(artists_billboard['
   edad_en_billboard']), 'edad_en_billboard'] =
   age_null_random_list
32 artists_billboard['edad_en_billboard'] = artists_billboard['
   edad_en_billboard'].astype(int)
33 print("Edad Promedio: " + str(age_avg))
34 print("Desvio Std Edad: " + str(age_std))
35 print("Intervalo para asignar edad aleatoria: " + str(int(
```

```

age_avg - age_std)) + " a " + str(int(age_avg + age_std)))
36
37 # Codificación de variables (mood, tempo, genre, artist_type,
    edad, duracion)
38 artists_billboard['moodEncoded'] = artists_billboard['mood'].
    map({
39     'Energizing': 6, 'Empowering': 6, 'Cool': 5, 'Yearning':
        4,
40     'Excited': 5, 'Defiant': 3, 'Sensual': 2, 'Gritty': 3,
41     'Sophisticated': 4, 'Aggressive': 4, 'Fiery': 4, 'Urgent'
        : 3,
42     'Rowdy': 4, 'Sentimental': 4, 'Easygoing': 1, 'Melancholy'
        : 4,
43     'Romantic': 2, 'Peaceful': 1, 'Brooding': 4, 'Upbeat': 5,
44     'Stirring': 5, 'Lively': 5, 'Other': 0, '': 0
45 }).astype(int)
46 artists_billboard['tempoEncoded'] = artists_billboard['tempo'
    ].map({
47     'Fast Tempo': 0, 'Medium Tempo': 2, 'Slow Tempo': 1, '':
        0
48 }).astype(int)
49 artists_billboard['genreEncoded'] = artists_billboard['genre'
    ].map({
50     'Urban': 4, 'Pop': 3, 'Traditional': 2, 'Alternative &
        Punk': 1,
51     'Electronica': 1, 'Rock': 1, 'Soundtrack': 0, 'Jazz': 0,
52     'Other': 0, '': 0
53 }).astype(int)
54 artists_billboard['artist_typeEncoded'] = artists_billboard['
    artist_type'].map({
55     'Female': 2, 'Male': 3, 'Mixed': 1, '': 0
56 }).astype(int)
57 artists_billboard.loc[artists_billboard['edad_en_billboard']
    <= 21, 'edadEncoded'] = 0
58 artists_billboard.loc[(artists_billboard['edad_en_billboard']
    > 21) & (artists_billboard['edad_en_billboard'] <= 26), '
    edadEncoded'] = 1
59 artists_billboard.loc[(artists_billboard['edad_en_billboard']
    > 26) & (artists_billboard['edad_en_billboard'] <= 30), '
    edadEncoded'] = 2
60 artists_billboard.loc[(artists_billboard['edad_en_billboard']
    > 30) & (artists_billboard['edad_en_billboard'] <= 40), '
    edadEncoded'] = 3
61 artists_billboard.loc[artists_billboard['edad_en_billboard']
    > 40, 'edadEncoded'] = 4

```

```

62 artists_billboard.loc[artists_billboard['durationSeg'] <=
    150, 'durationEncoded'] = 0
63 artists_billboard.loc[(artists_billboard['durationSeg'] >
    150) & (artists_billboard['durationSeg'] <= 180), '
    durationEncoded'] = 1
64 artists_billboard.loc[(artists_billboard['durationSeg'] >
    180) & (artists_billboard['durationSeg'] <= 210), '
    durationEncoded'] = 2
65 artists_billboard.loc[(artists_billboard['durationSeg'] >
    210) & (artists_billboard['durationSeg'] <= 240), '
    durationEncoded'] = 3
66 artists_billboard.loc[(artists_billboard['durationSeg'] >
    240) & (artists_billboard['durationSeg'] <= 270), '
    durationEncoded'] = 4
67 artists_billboard.loc[(artists_billboard['durationSeg'] >
    270) & (artists_billboard['durationSeg'] <= 300), '
    durationEncoded'] = 5
68 artists_billboard.loc[artists_billboard['durationSeg'] > 300,
    'durationEncoded'] = 6
69
70 drop_elements = ['id', 'title', 'artist', 'mood', 'tempo', 'genre'
    , 'artist_type', 'chart_date', 'anioNacimiento', 'durationSeg'
    , 'edad_en_billboard']
71 artists_encoded = artists_billboard.drop(drop_elements, axis
    =1)
72 colormap = plt.cm.viridis
73 plt.figure(figsize=(12,12))
74 plt.title('Pearson Correlation of Features', y=1.05, size=15)
75 sb.heatmap(artists_encoded.astype(float).corr(), linewidths
    =0.1, vmax=1.0, square=True, cmap=colormap, linecolor='
    white', annot=True)
76 plt.show()

```

4. Árbol de Decisión

Se realiza la validación mediante KFold, se evalúa la precisión del modelo para diferentes profundidades y se crea el árbol de decisión definitivo para la clasificación:

```

1 cv = KFold(n_splits=10) # Numero de folds
2 accuracies = list()
3 max_attributes = len(list(artists_encoded))
4 depth_range = range(1, max_attributes + 1)
5

```



```

6 # Testeamos la profundidad de 1 a la cantidad de atributos +
  1
7 for depth in depth_range:
8     fold_accuracy = []
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         model = tree_model.fit(X=f_train.drop(['top'], axis
            =1),
18                                y=f_train["top"])
19         valid_acc = model.score(X=f_valid.drop(['top'], axis
            =1),
20                                y=f_valid["top"])
21         fold_accuracy.append(valid_acc)
22         avg = sum(fold_accuracy) / len(fold_accuracy)
23         accuracies.append(avg)
24
25 df = pd.DataFrame({"Max Depth": depth_range, "Average
        Accuracy": accuracies})
26 print(df.to_string(index=False))
27
28 # Entrenamos el arbol de decision con profundidad = 4
29 y_train = artists_encoded['top']
30 x_train = artists_encoded.drop(['top'], axis=1)
31 decision_tree = tree.DecisionTreeClassifier(criterion='
        entropy',
32                                                min_samples_split
            =20,
33                                                min_samples_leaf
            =5,
34                                                max_depth=4,
35                                                class_weight
            ={1:3.5})
36 decision_tree.fit(x_train, y_train)
37
38 # Exportamos el modelo a archivo .dot y generamos la imagen

```

```
del arbol
39 with open(r"tree1.dot", 'w') as f:
40     f = tree.export_graphviz(decision_tree,
41                               out_file=f,
42                               max_depth=7,
43                               impurity=True,
44                               feature_names=list(
45                                   artists_encoded.drop(['top'],
46                                                         axis=1)),
47                               class_names=['No', 'N1 Billboard'],
48                               rounded=True,
49                               filled=True)
50
51 check_call(['dot', '-Tpng', r'tree1.dot', '-o', r'tree1.png'])
52 PImage("tree1.png")
53
54 # Precision del arbol de decision
55 acc_decision_tree = round(decision_tree.score(x_train,
56                                                y_train) * 100, 2)
57 print(acc_decision_tree)
58
59 # Prediccion Camila Cabello
60 features = x_train.columns.tolist()
61 x_test = pd.DataFrame(columns=features)
62 x_test.loc[0] = [5, 2, 4, 1, 0, 3]
63 y_pred = decision_tree.predict(x_test)
64 print("Prediccion: " + str(y_pred))
65 y_proba = decision_tree.predict_proba(x_test)
66 print("Probabilidad de Acierto: " + str(round(y_proba[0][
67     y_pred][0] * 100, 2)) + "%")
68
69 # Prediccion Imagine Dragons
70 x_test = pd.DataFrame(columns=features)
71 x_test.loc[0] = [4, 2, 1, 3, 2, 3]
72 y_pred = decision_tree.predict(x_test)
73 print("Prediccion: " + str(y_pred))
74 y_proba = decision_tree.predict_proba(x_test)
75 print("Probabilidad de Acierto: " + str(round(y_proba[0][
76     y_pred][0] * 100, 2)) + "%")
```

3. Resultados

Algunos de los resultados que se obtuvieron fueron:

- **Evaluación de la profundidad del árbol:** Se generó una tabla con la precisión promedio para cada profundidad evaluada.
- **Visualización de correlación:** Se muestra el mapa de calor de correlación entre características.
- **Modelo final de Árbol de Decisión:** Se entrenó un árbol con profundidad 4, se exportó su estructura a un archivo .png y se calculó su precisión.
- **Predicciones:** Se realizaron predicciones para nuevos casos (ej. Camila Cabello e Imagine Dragons) junto con las probabilidades de acierto.

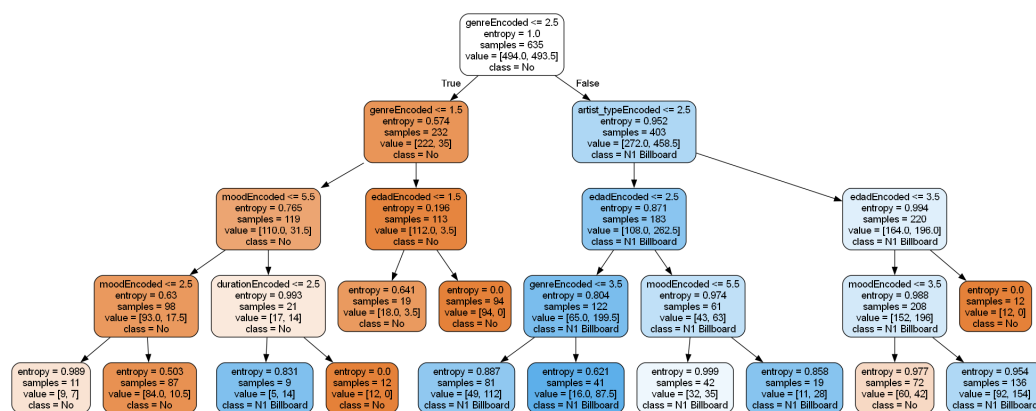
```

1 Prediccion: [1]
2 Probabilidad de Acierto: 84.54%
3 Prediccion: [0]
4 Probabilidad de Acierto: 88.89%

```

Como resultado de las predicciones, para Camila Cabello, **Havana** llegará al **top 1** con una **probabilidad del 84 %**. Por otra parte, para Imagine Dragons, la segunda predicción nos menciona que **Believer** **NO** llegará al **Top 1** con una **certeza del 88 %**.

Visualización del Árbol de Decisión de forma gráfica:



4. Conclusión

Se aplicó el **Árbol de Decisión** para clasificar y predecir el desempeño en el billboard a partir de diversas características de los artistas. Se exploraron y preprocesaron los datos, se evaluaron distintos parámetros del modelo mediante validación cruzada y se seleccionó una profundidad óptima para el árbol. Los resultados muestran que es una herramienta interpretativa y efectiva para la clasificación en escenarios complejos, permitiendo además identificar la importancia de cada característica en predicciones finales.

```

1 Edad Promedio: 30.10282258064516
2 Desvio Std Edad: 8.40078832861513
3 Intervalo para asignar edad aleatoria: 21 a 38
4 mood
5 Empowering      88
6 Cool            78
7 Energizing      71
8 Yearning        68
9 Defiant         64
10 Excited         62
11 Sensual         55
12 Urgent          30
13 Rowdy           15
14 Sophisticated   15
15 Brooding        14
16 Upbeat          12
17 Aggressive      11
18 Melancholy      10
19 Gritty           9
20 Sentimental     7
21 Romantic        7
22 Fiery           6
23 Easygoing       4
24 Peaceful        4
25 Lively          2
26 Stirring        2
27 Other           1
28 dtype: int64
29 ### ### ###
30 Tempos de Cancion: ['Medium Tempo' 'Slow Tempo' 'Fast Tempo']
31 ### ### ###
32 Tipos de Artista: ['Male' 'Female' 'Mixed']
33 ### ### ###
34 genre

```

```
35 Urban                215
36 Pop                  188
37 Traditional          113
38 Alternative & Punk   60
39 Electronica          27
40 Rock                 13
41 Other                12
42 Soundtrack           5
43 Jazz                 2
44 dtype: int64
45 Max Depth  Average Accuracy
46          1          0.556101
47          2          0.556126
48          3          0.564038
49          4          0.644122
50          5          0.601711
51          6          0.609400
52          7          0.645610
53 64.88
54 Prediccion: [1]
55 Probabilidad de Acierto: 84.54%
56 Prediccion: [0]
57 Probabilidad de Acierto: 88.89%
```