

19-202-0604 Data Mining

Assignment Report

Aparna Nair

Predicting Song Popularity Based On Spotify
Classifying Features Using Logistic Regression

B.Tech Computer Science

CS-A Register Number: 20220021

July 30, 2023

Introduction

Predicting song popularity using spotify classifying features. These features are: Title, Artist, Top Genre, Year, Beats Per Minute (BPM), Energy, Danceability, Loudness (dB), Liveness, Valence, Length (Duration), Acousticness, Speechiness, Popularity.

Dataset

First, find the dataset required on Kaggle: Spotify - All Time Top 2000s Mega Dataset. This is a dataset containing audio features of Top 2000 songs on Spotify of all time. The data contains about 15 columns each describing the track and it's qualities. Songs released from 1956 to 2019 are included from some notable and famous artists like Queen, The Beatles, Guns N' Roses, etc.

Content

- Index: ID
- Title: Name of the Track
- Artist: Name of the Artist
- Top Genre: Genre of the track
- Year: Release Year of the track
- Beats per Minute(BPM): The tempo of the song
- Energy: The energy of a song - the higher the value, the more energtic. song
- Danceability: The higher the value, the easier it is to dance to this song.
- Loudness: The higher the value, the louder the song.
- Valence: The higher the value, the more positive mood for the song.
- Length: The duration of the song.
- Acoustic: The higher the value the more acoustic the song is.
- Speechiness: The higher the value the more spoken words the song contains.

These are the Spotify classifying features.

Implementation in Google Colab

```
1  #Mount the google drive
2  from google.colab import drive
3  drive.mount('/content/drive')
4
5  #Import nrequired libraries
6  # Visualizations
7  import pandas as pd
8  import numpy as np
9  import datetime
10 import plotly.express as px
11 import seaborn as sns
12 from matplotlib import pyplot as plt
13 %matplotlib inline
14 sns.set_style("whitegrid")
15 import matplotlib.pyplot as plt
16
17
18 # For transformations and predictions
19 from sklearn.preprocessing import StandardScaler
20 from sklearn.model_selection import train_test_split
21
22 # Models to be used, all from sklearn
23 from sklearn.ensemble import RandomForestClassifier
24 from sklearn.linear_model import LogisticRegression
25 from sklearn.neighbors import KNeighborsClassifier
26 from sklearn.tree import DecisionTreeClassifier
27
28 # Transformers
29 from scipy.optimize import curve_fit
30 from sklearn.preprocessing import FunctionTransformer
31 from sklearn.preprocessing import OneHotEncoder
32 from sklearn.metrics import pairwise_distances
33 from six import StringIO
34
35 # For comparing metrics
36 from sklearn.metrics import mean_squared_error as mse
37
38 # For validation
39 from sklearn.model_selection import train_test_split as split
40
41
42
43 # Step 1: Data Preprocessing
44 # Load the dataset
45 df = pd.read_csv("/content/drive/MyDrive/Spotify-2000.csv")
46
47 # Returns the first 1000 rows
```

```
48 df.head(1000)
49
50 df.isnull().sum()
51 # no nulls, all good
52
53 # print(len(df.index))
54
55 df = df.iloc[0:2000]
56
57 # There were some non-numeric values contained in the duration column,
58 # that should have been measured in milliseconds
59 df[df["Length (Duration)".str.contains(",")==True]
60
61 df["Length (Duration)"] = df["Length (Duration)".replace(["1,412", "
        1,121", "1,367", "1,292"], ['1412', '1121', '1367', '1292'])
62 #These were corrected manually
63
64 df["Length (Duration)"] = df["Length (Duration)".astype(np.int64)
65 #Casting all values to integers
66
67 columns = ['Index', 'Title']
68 # for col in columns:
69     # print(f'{col:<15}: {df[col].nunique()} unique values')
70
71 df = df.drop(labels=['Index', 'Title'], axis=1)
72 df.shape
73
74 #No duplicates
75 df.duplicated().sum()
76
77 #No null values
78 df.isnull().sum().sum()
79
80 #As danceability increases, popularity increases according to spotify
81 fig, ax = plt.subplots(1, figsize=(15, 6), sharey=True, sharex = True)
82 ax_data = df.groupby('Danceability')['Popularity'].mean().to_frame().
        reset_index()
83 ax = sns.scatterplot(x='Danceability', y='Popularity', data=ax_data,
        color='blue', ax=ax)
84 ax.set_title('Danceability')
85 ax.set_ylabel('Mean Popularity', fontsize=12)
86 plt.tight_layout()
87 plt.show()

88 #Results of initial numeric correlation, shows the features are not good
        indicators of target variable, popularity
89 corr = np.abs(numeric_df.corr())
90 series = np.abs(corr['Popularity']).sort_values(ascending=False)
91 print('The most linear correlated features to POPULARITY are:')
```

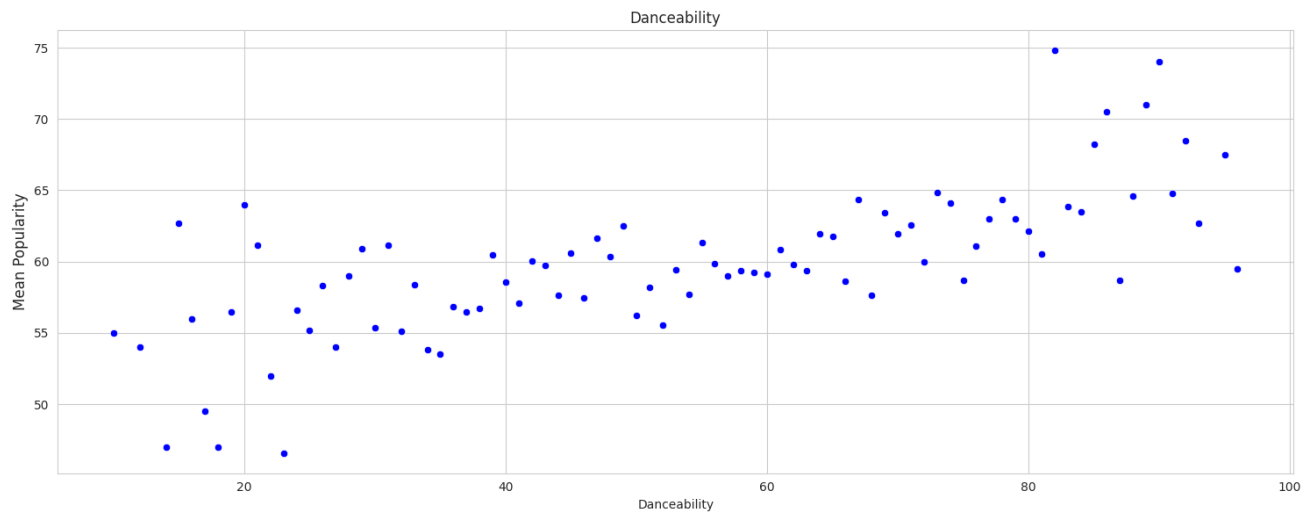


Figure 1: Scatter graph of Danceability vs Mean Popularity

```

92 for i, row in enumerate(series):
93     if 0.1 <= row < 1:
94         print(f'{series.index[i]:17} --> {row: .2f} (abs)')

```

The most linear correlated features to POPULARITY are:

Loudness (dB) -> 0.17 (abs)

Year -> 0.16 (abs)

Danceability -> 0.14 (abs)

Liveness -> 0.11 (abs)

Speechiness -> 0.11 (abs)

Energy -> 0.10 (abs)

```

95 #feature engineering, classify popular songs as 1, more unpopular songs
   as 0
96 df['Popularity'].describe()
97 df['Popularity'] = pd.qcut(df['Popularity'], q=2, labels=[0, 1])
98 df["Popularity"].value_counts()
99 df['Popularity']
100
101 #One-hot encoding to get dummies function for the non-numeric values.
102 def onehot_encode(df, column, prefix):
103     df = df.copy()
104     dummies = pd.get_dummies(df[column], prefix=prefix)
105     df = pd.concat([df, dummies], axis=1)
106     df = df.drop(column, axis=1)
107     return df
108
109 #one hot encoding the genre and artists
110 df = onehot_encode(df, 'Artist', 'artist')
111 df = onehot_encode(df, 'Top Genre', 'genre')

```

```
112 df
113
114 #splitting the data
115 y = df.loc[:, 'Popularity']
116 X = df.drop('Popularity', axis=1)
117 # X = X.drop('Top Genre', axis=1)
118
119 scaler = StandardScaler()
120 X = scaler.fit_transform(X)
121
122 #random = 369
123 X_train, X_test, y_train, y_test = train_test_split(X, y, train_size
    =0.7, random_state=369)
124
125 print("X_train: ",X_train.shape)
126 print("X_test: ",X_test.shape)
127 print("y_train: ",y_train.shape)
128 print("y_test: ",y_test.shape)
```

$X_{train} : (1395, 890)$

$X_{test} : (599, 890)$

$y_{train} : (1395,)$

$y_{test} : (599,)$

```
129 #Accuracy
130 log_model = LogisticRegression()
131 log_model.fit(X_train, y_train)
132 log_acc = log_model.score(X_test, y_test)
133 print("Logistic Regression Accuracy:", log_acc)
```

Logistic Regression Accuracy: 0.7161936560934892

Results and interpretation

From the accuracy count which is 0.71 using Logistic Regression, we can safely say that this model can be used to predict the popularity of a song against spotify classifying features.

My interpretation from this case study is that we can not in fact, successfully predict the popularity of a song from spotify features alone. While some features were largely useful, such as Danceability, a combination of BPM over duration, most were not, Valence, Energy, Loudness, etc.

Spotify Feature	Target Variable, (Popularity)
Loudness (dB)	0.17 (abs)
Year	0.16 (abs)
Danceability	0.14 (abs)
Liveness	0.11 (abs)
Speechiness	0.11 (abs)
Energy	0.10 (abs)

Table 1: Shows all the features that are not required for explaining song popularity