

Spotify Songs Analysis

Amy Seo November 2021 Data Analytics & Visualizations Bootcamp Final Project GitHub <u>Link</u>

Project Topic

Analysis of songs that have charted on the Spotify Top 200s Chart for the United States since January 1, 2021.

- Numerous songs are released every year in multiple countries
- Understanding insights into popularity of songs/artists in **2021**, focusing on **United States**
- Model to predict success of song based on auditory/musical features
 - Can be useful in anticipating what songs will be most popular and influential
 - Can provide insight into what types of songs are appealing to people at a given point in time

Project Questions

- What songs have the most streams?
- What songs have been #1 on the charts?
- Which artists have the most streams?
- Which artists have the most songs on the chart?
- Which artists have the most #1 hits?
- Which songs are most danceable? Most energetic? Most positive-sounding (high valence)? (and so on, for each audio feature)
- Can we predict if a song can break into the top 20 positions, based on the song's audio features?

Data Source

Spotify Top 200 Charts

- Top 200 Charts, United States
- Daily Data, from Jan. 1, 2021 to Nov. 17, 2021
- Source: Spotify Charts Site
- Scraped Data CSV file

Spotify Audio Features Data

- Spotify API Request for Tracks' audio features data
- Source: API Request Info Site
- Saved Audio Features <u>CSV file</u>

Data Extraction Process & Tools

- 1. Charts data was scraped using:
 - a. splinter
 - b. BeautifulSoup
 - c. ChromeDriverManager

```
28 # Define function to run scraping process
29 def scrape(date urls):
      # Initiate headless driver for deployment
      executable_path = {'executable_path': ChromeDriverManager().install()}
       browser - Browser('chrome', **executable_path, headless=False)
      # Run scraping function
       for url in date urls:
          # Visit URL
           browser.visit(url)
           # Set up HTML parser
           html = browser.html
           song soup = soup(html, 'html.parser')
           playlist = song soup.find('table', class ='chart-table')
           for i in playlist.find('tbody').findAll('tr'):
               song = i.find('td', class_='chart-table-track').find('strong').get_text()
               artist = i.find('td', class_='chart-table-track').find('span').get_text()
              artist = artist.replace("by ", "").strip()
               songurl = i.find('td', class ='chart-table-image').find('a').get('href')
               songid = songurl.split("track/")[1]
               position = i.find('td', class ='chart-table-position').get text()
               streams = i.find('td', class_='chart-table-streams').get_text()
               date = url.split("daily/")[1]
               data.append([songid, songurl, song, artist, date, position, streams])
      # Stop webdriver
      browser.quit()
```

- 2. Features data pulled using an API request:
 - a. requests

```
In [7]: # Import access tokens for Spotify API
        from config import client_id, client_secret
In [8]: # Pass the access token
        auth url = 'https://accounts.spotify.com/api/token
        auth response = requests.post(auth url, (
             'client id': client id.
             'client_secret': client_secret
        auth_response_data = auth_response.json()
access token = auth_response_data('access token')
In [9]: # Set up access token in header for GET request
         headers = {'Authorization': 'Bearer {token}'.format(token-access_token)}
        base_url = 'https://api.spotify.com/v1/audio-features/
In [11]: # GET request for song_id string lists
        data - []
         for i in song_idsl:
            req = requests.get(base_url + i, headers=headers)
             reg = reg.ison()
             data.append(reg)
         for i in song_ids2:
            req - requests.get(base_url + i, headers-headers)
             req = req.json()
             data.append(reg)
         for i in song ids3:
            req = requests.get(base_url + i, headers=headers)
             req = req.json()
            data.append(reg)
         for i in song ids4:
            req = requests.get(base_url + i, headers=headers)
             reg = reg.ison()
             data.append(reg)
            req = requests.get(base_url + i, headers=headers)
             req = req. ison()
         for i in song_ids6:
```

File Links:

"<u>scraping.py</u>" script <u>Jupyter Notebook</u> for data loading/cleaning/databasing

Data Cleaning Process & Tools

Scraped data table ("scrape_df") and API data table ("features_df") both cleaned using Python and Pandas.

- Unnecessary columns removed
- Column values relabeled as needed
- Data types fixed
- Data checked for null values

Additional data tables created for aggregates and merged for final table to be used for analysis:

- Total streams aggregated for each song
- Highest position in charts recorded for each song
- List of unique track names and artists
- The above three listed tables and "features_df" merged for final dataframe for analysis

```
64 SELECT DISTINCT ON (ta.song_id) ta.song_id,
        ta.song,
        ta.artist.
        ts.streams.
        hp.position,
        f.danceability,
        f.energy.
        f.loudness.
        f.mode.
        f.speechiness,
        f.acousticness.
        f.instrumentalness.
77
        f.liveness,
        f.valence,
        f.tempo,
        f.duration ms.
        f.time signature
    FROM track artist AS ta
    ON (ta.song_id = ts.song_id)
    ON (ta.song_id = hp.song_id)
89 ON (ta.song id = f.song id):
```

File Links:

<u>Jupyter Notebook</u> for data loading/cleaning/databasing SQL <u>script</u> for merging tables

Database

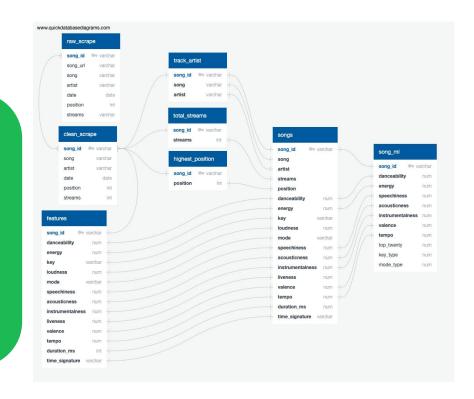
Database Type: PostgreSQL

Link: sqlalchemy

Number of Tables: 8

- Includes raw and cleaned tables
- "songs" table = final merged data used for analysis
- "song_ml" table = final data used for machine learning inputs

Primary Key: "song_id"



File Links:

<u>Jupyter Notebook</u> for data loading/cleaning/databasing SQL <u>script</u> for creating tables

Data Exploration Process & Tools

Tools Used: Python, Pandas, Plotly

Questions Answered

What songs have the most streams?

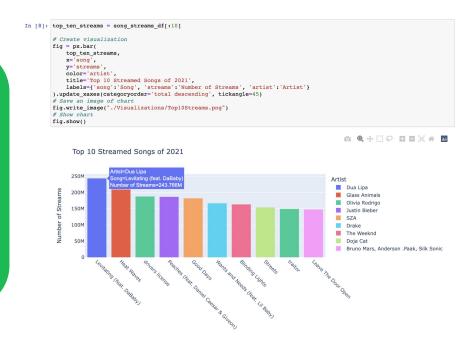
What songs have been #1 on the charts?

Which artists have the most streams?

Which artists have the most songs on the chart?

Which artists have the most #1 hits?

Which songs are most danceable? Most energetic? Most positive-sounding (high valence)? (and so on, for each audio feature)



File Links:

<u>Jupyter Notebook</u> for data exploration Plotly Chart Visualization <u>Image Files</u>

Data Exploration

Walk-through of Exploration Analysis & Visualizations
(Flask Dashboard)

Flask App/Dashboard Files for Compiled Chart Visualizations:

index.html

app.py

Machine Learning Model Question

Can we predict if a song can break into the top 20 positions, based on the song's audio features?

Tools Used:

- Pandas → to clean "songs" table for model ("song_ml" table)
- scikit-learn library
- StandardScaler
- LabelEncoder
- SMOTEENN
- RandomForestClassifier

File Links:

<u>Jupyter Notebook</u> for machine learning model

1) Additional column ("top_twenty") added to codify whether a song reached a rank within top 20 ranks or not.



2) Unnecessary columns removed from dataframe.

File Links:

Jupyter Notebook for machine learning model

3) Categorical variables encoded to numerical data. In [7]: # Generate categorical variable list song_cat = song_ml_df.dtypes[song_ml_df.dtypes == "object"].index.tolist() # Check number of unique values in each column song_ml_df[song_cat].nunique() Out[7]: key 12 dtype: int64 Encode Dataframe's non-numerical data In [8]: # Create instance for label encoder labelencoder = LabelEncoder() In [9]: # Encode categorical columns and store as another column song_ml_rf_df = song_ml_df.copy() song_ml_rf_df['key_type'] = labelencoder.fit_transform(song_ml_rf_df['key']) song_ml_rf_df['mode_type'] = labelencoder.fit_transform(song_ml_rf_df['mode']) # Drop the original columns song ml rf df = song ml rf df.drop(["key", "mode"], 1) song ml rf df.head() Out[9]: danceability energy speechiness acousticness instrumentalness valence tempo top_twenty key_type mode_type 0.0534 0.64300 0.003910 0.549 84.997 0.0282 0.24700 0.000073 0.442 139.971

4) Target (newly created "top_twenty" column) and features assigned.

```
Assign features and target variables

In [11]: # Assign preprocessed data into features and target arrays
y = song_ml_rf_df["top_twenty"].ravel()
X = song_ml_rf_df.drop(["top_twenty"], 1)
```

File Links:

Jupyter Notebook for machine learning model

01K4zKU104LyJ8gMb7227B

5) Data is split into training and test sets. Then, data is resampled using SMOTEENN technique (to help with class imbalance). Resampled data is then scaled given different scales of features.

```
Split, resample, and scale data for model
         Using SMOTEENN combination sampling to address disparate class sizes
In [13]: # Split preprocessed data into training and testing dataset
         X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=1)
In [14]: # Check dataset split of training set
         Counter(y train)
Out[14]: Counter({0: 1087, 1: 220})
In [15]: # Data resampled with SMOTEENN
         smote enn = SMOTEENN(random state=1)
         X resampled, y resampled = smote enn.fit resample(X train, y train)
         Counter(y_resampled)
Out[15]: Counter({0: 375, 1: 699})
In [16]: # Create StandardScaler instances
         scaler = StandardScaler()
In [17]: # Fit the StandardScaler
         X scaler = scaler.fit(X resampled)
In [18]: # Scale data
         X train scaled = X scaler.transform(X resampled)
         X test scaled = X scaler.transform(X test)
```

File Links:

Jupyter Notebook for machine learning model

6) Random Forest Model is run on the data. Model is trained using 100 estimators.

Random Forest Model

In [19]: # Create a random forest classifier
 rf_model = RandomForestClassifier(n_estimators=100, random_state=1)

In [20]: # Fitting the model
 rf_model = rf_model.fit(X_train_scaled, y_resampled)

File Links:

Jupyter Notebook for machine learning model

Data Analysis/Model Results

```
In [23]: # Training score & Testing score
         print(f'Training Data Score: {rf model.score(X train scaled, y resampled)}')
         print(f'Testing Data Score: {rf_model.score(X_test_scaled, y_test)}')
         Training Data Score: 1.0
         Testing Data Score: 0.5711009174311926
In [24]: # Confusion matrix
         cm = confusion matrix(y test, y pred)
         cm df = pd.DataFrame(cm, index=["Actual 0", "Actual 1"], columns = ["Predicted 0", "Predicted 1"])
         cm df
Out[24]:
                 Predicted 0 Predicted 1
          Actual 0
                       216
                                142
          Actual 1
                       45
                                 33
In [25]: # Classification report
         print(f'Classification Report:\n{classification report(y test, y pred)}')
         Classification Report:
                       precision
                                     recall f1-score
                                                        support
                             0.83
                                       0.60
                                                 0.70
                                                            358
                             0.19
                                       0.42
                                                 0.26
                                                             78
                                                 0.57
                                                            436
             accuracy
                             0.51
                                       0.51
                                                 0.48
                                                            436
            macro avg
         weighted avg
                            0.71
                                       0.57
                                                 0.62
                                                            436
```

File Links:

Jupyter Notebook for machine learning model

Data Analysis/Model

Walk-through of Interactive Machine Learning Model
(Flask Dashboard)

File Links:

<u>Jupyter Notebook</u> for machine learning model

Recommendations for Future Analyses

Recommendations

- Testing if audio features' have statistically different values between songs that do reach top 20 vs. don't.
- Testing other new features, such as...
 - Artist popularity/followers
 - Time of release (year/season)

Anything Done Different?

- Conduct in-depth, statistical analyses on each of audio features first to determine usability in model
- Test with wider dataset across multiple years and across different countries

Thank you!

