

Emily You and Mia Inakage
Dr. Barbara Ericson
26 April 2022

SI 206 Final Project Report: Python Girls
<https://github.com/minakage/SI206PythonGirls>

Original Goals for the Project

For our SI 206 project, our goal was to retrieve data using a Spotify API (Spotipy) and Wikipedia using Beautiful Soup. For Spotify, we wanted to get playlist data from one project member's Spotify account, specifically the user's top 5 playlists of 2022 with their playlist name, song names in each playlist, and the artist of each song in the playlist. Once we knew the artist of each song in the playlist, we planned to gather the number of monthly listeners for each of the artists on Spotify. Given the artist names gathered from the Spotify data, we wanted to retrieve the ages of those artists from Wikipedia. We planned to calculate a correlation between the age and monthly listeners for each respective artist.

We wanted to create 3 visualizations using matplotlib:

1. One bar graph to show the number of monthly listeners (in millions) for each artist in the user's top 5 playlists on Spotify
 - a. X-axis would be artist name
 - b. Y-axis would be monthly listeners (in millions)
2. One bar graph to show which age group is most common among the artists in the user's top 5 playlists on Spotify
 - a. X-axis would be artist name
 - b. Y-axis would be age group (grouped by 10 years)
3. One scatter plot to identify a correlation between an artist's age and their monthly listeners
 - a. X-axis would be artist age
 - b. Y-axis would be monthly listeners (in millions)

Our goal was to find out more information about the most commonly played artists in 2022 like which artists have the most monthly listeners and the ages of artists in the user's top 5 playlists of 2022.

New Goals and Goals that were Achieved

After reviewing our project plan, we realized our goals were a little too ambitious and that we realized that the Spotipy API would not allow us to find the number of monthly listeners that we originally intended to do. We also had difficulty with accessing a user's Spotify account and finding their top 5 playlists of 2022 while also satisfying the requirement for extracting at

least 100 data points for each API. Thus, we decided to select 10 artists and extract the popularity values and duration of each song in their top 10 songs on Spotify.

10 Artists We Chose for Our Project

- Olivia Rodrigo
- Bazzi
- Billie Eilish
- Mac Miller
- Lady Gaga
- Harry Styles
- Post Malone
- Rihanna
- Elton John
- Madonna

The Spotipy API has a popularity index that rates songs on a scale of 0 to 100 and a song duration index that measures the length of each song in milliseconds. We also assigned each artist an artist id in order to avoid duplicate strings in our database. Our calculation resulted in the average popularity of the top 10 songs for each artist. Instead of using Beautiful Soup to scrape data from Wikipedia, we decided to use a Genius API. With the same 100 songs we used from the Spotify data, we could find the page view count and annotation count for each song within the Genius API. We kept with our original plan of creating visualizations using matplotlib but changed the data that would populate in our graphs.

We ended up creating 4 visualizations using matplotlib:

4. One bar graph to show the average popularity of the top 10 songs for 10 selected artists on Spotify.
 - a. X-axis is average popularity
 - b. Y-axis is artist name
5. One bar graph to show the duration of each of the top 10 songs for 10 selected artists on Spotify
 - a. X-axis is song title
 - b. Y-axis is duration (in ms)
6. One bar graph to show the page view count of each of the 100 songs from Genius
 - a. X-axis is song title
 - b. Y-axis is page views (in hundred thousands)
7. One bar graph to show the annotation count of each of the 100 songs from Genius
 - a. X-axis is song title
 - b. Y-axis is annotation count

We arrived at a few new goals. For Spotify, we wanted to identify the average popularity and variation of song durations for the top ten songs for 10 artists on Spotify. For Genius, our goal was to see if there was a pattern in the number of page views or annotations from the lyrics of each song that we took from Spotify. We successfully achieved our goals because we were familiar with working with APIs in previous homework and discussion assignments. We also were reminded of how to create calculations and databases. One key lesson that we learned was how to access the Spotify and Genius APIs and discover which indices we could find from them.

Problems Faced

We faced multiple problems while trying to achieve our new goals, one of which was having duplicate strings in our Spotify and Spotify 2 database tables. Specifically, since we gathered the top 10 songs from 10 artists, we would have the same artist for 10 rows in the 'artist' column of our Spotify and Spotify 2 database tables. For example, an artist (i.e. 'Olivia Rodrigo') would appear 10 times in the rows with a song name and popularity or song duration value. In order to eliminate the duplicate strings, we modified our tables so that each artist was given an artist_id (integer). Each table then had duplicate id values instead of duplicate strings. Another problem we faced was that the artist_id value could not be the primary key shared between the Spotify and Spotify 2 tables because artist_id had duplicate values. In order to ensure that both tables shared a primary key, we made the song title the primary key between the two tables and VARCHAR NOT NULL, which forces the data in the song title column to always contain a value. The song titles are unique to each row. We also had a problem with readjusting the size of the figure of our visualizations because the data was difficult to view. The titles and bars on the bar graph were often either too cramped together, or the titles were out of the frame and unable to be read. We found a solution to this problem by changing the figure size of the plot using `fig = plt.figure(figsize=(#, #))` and adjusting the spacing of the bottom using `fig.subplots_adjust(bottom=spacing)` to show the x-axis.

Git Commits

Both of us decided not to create a collaborative Github to avoid running into the problem of git merge conflicts. Since we almost always worked on the project together in-person, both of us knew which code we had to work on for each work session and took turns making changes. When meeting during Zoom, we would work on code together by enabling the Remote Control feature. This is why Mia made all the commits to Github.

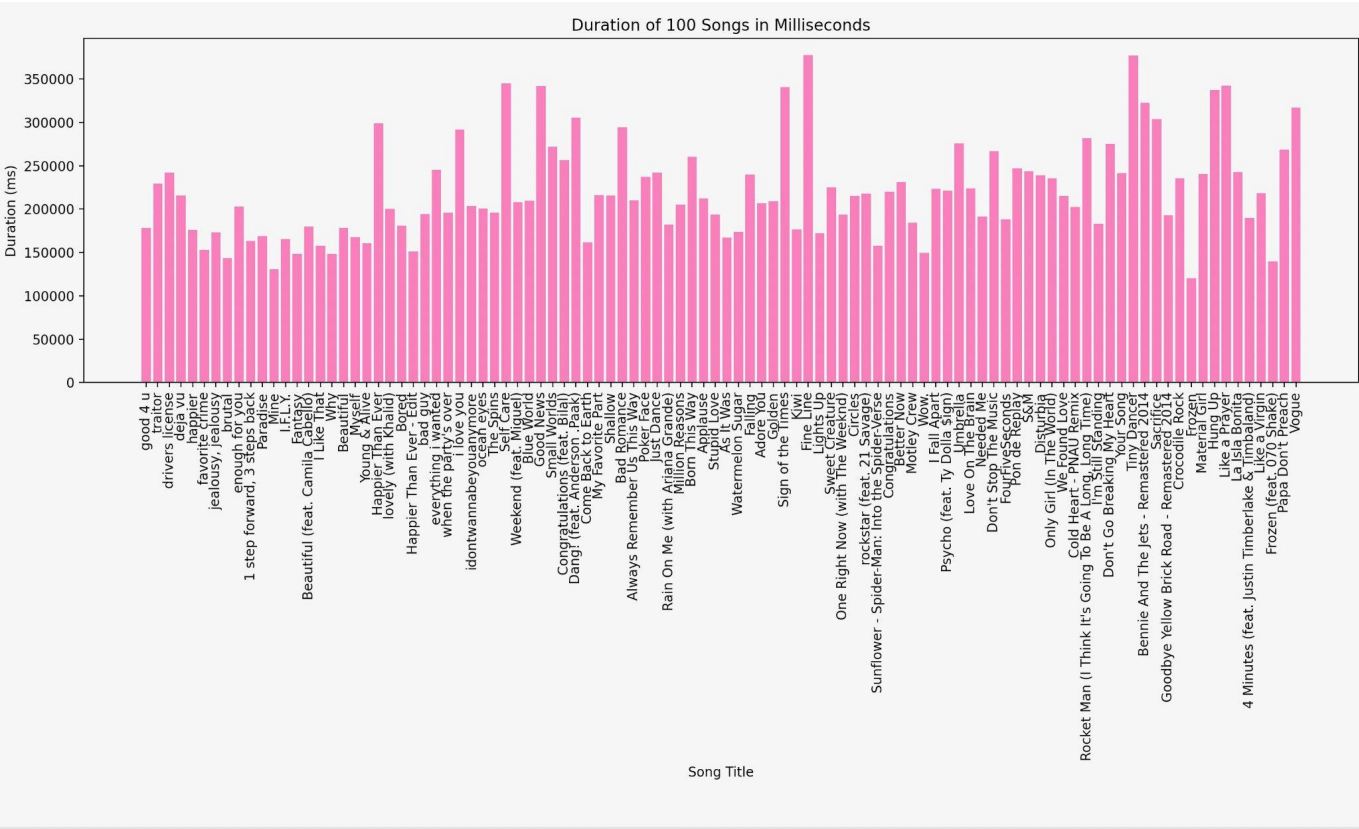
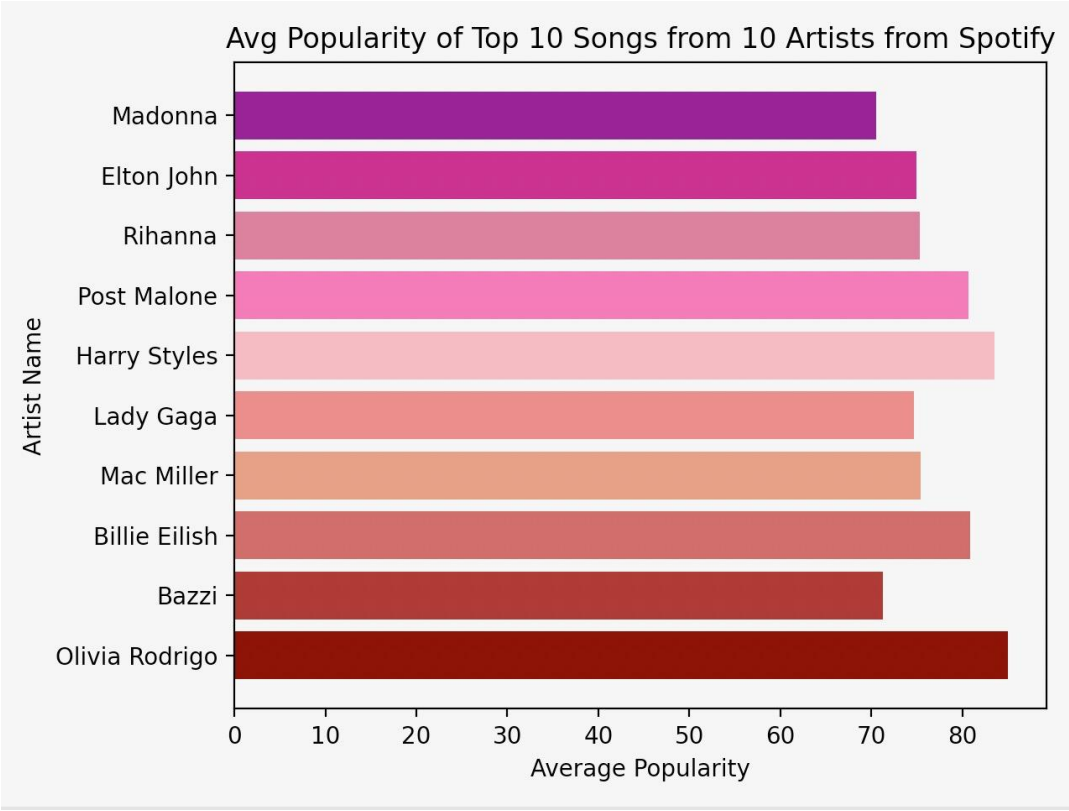
Files Containing Calculations from Data

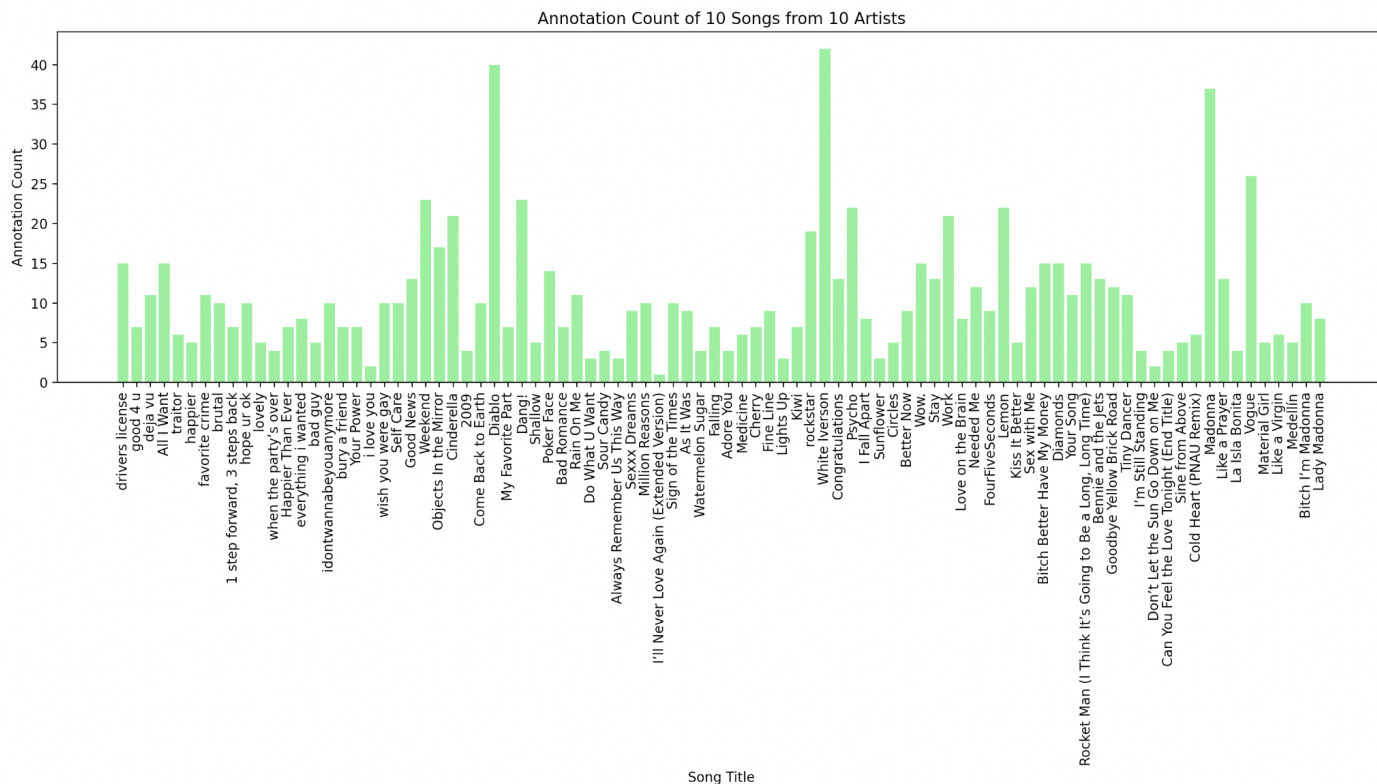
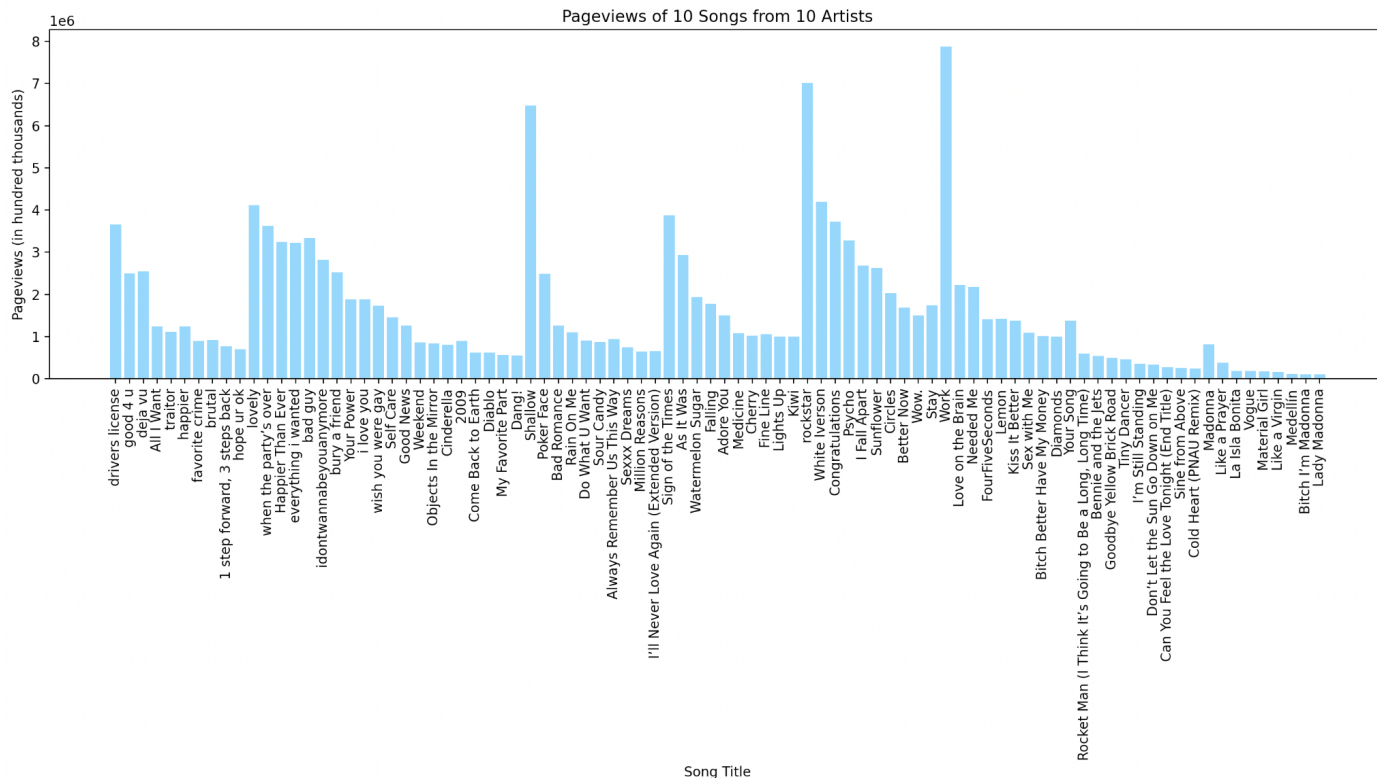
File containing our calculations from "spotifyPop.py" file: appears in popularityValues.txt

```
popularityValues.txt
```

```
1  Olivia Rodrigo has an average popularity of 85.0 from their top ten songs on Spotify
2  Bazzi has an average popularity of 71.3 from their top ten songs on Spotify
3  Billie Eilish has an average popularity of 80.9 from their top ten songs on Spotify
4  Mac Miller has an average popularity of 75.4 from their top ten songs on Spotify
5  Lady Gaga has an average popularity of 74.7 from their top ten songs on Spotify
6  Harry Styles has an average popularity of 83.5 from their top ten songs on Spotify
7  Post Malone has an average popularity of 80.7 from their top ten songs on Spotify
8  Rihanna has an average popularity of 75.3 from their top ten songs on Spotify
9  Elton John has an average popularity of 75.0 from their top ten songs on Spotify
10 Madonna has an average popularity of 70.5 from their top ten songs on Spotify
11
```

Visualizations (4)





Instructions for Running the Code

spotipyPop.py:

- Install spotipy by typing this in terminal [pip install spotipy --upgrade](#)
- Go to <https://developer.spotify.com/dashboard/applications>.
- Log in with your Spotify account.
- Click on 'Create an App'.
- Pick an 'App Name' and 'App Description' of your choice and mark the checkboxes.
- After creation, you see your 'Client Id', and you can click on 'Show Client Secret' to unhide your 'Client secret'.
- Use your 'Client id' and 'Client secret' to retrieve a token from the Spotify API. The function `get_spotify_api_token()` performs all necessary steps with your 'Client id' and 'Client secret' to retrieve a token.
- To get a redirected uri, go to the settings of the app you created and enter 'http://localhost:8888/callback'
- Go to spotipyPop.py and update your username, your client id, and your client secret
- Run the spotipyPop.py file
- It will give you the Python Girls database with 2 tables, Spotipy and Spotipy 2
 - Spotipy will have artist id, song, and popularity
 - Spotipy 2 will have artist id, song, and duration (in ms)
- It will also write out the calculations to a text file and generate 2 visualizations using matplotlib

genius.py:

- Go to https://genius.com/signup_or_login and sign up
- Once you're signed in, you will be taken to <https://genius.com/api-clients>, where you need to click the button that says "New API Client."
- After clicking "New API Client," you'll be prompted to fill out a short form about the "App" that you need the Genius API for. You only need to fill out "App Name" and "App Website URL."
- You can put "Song Lyrics Project" or anything you want for the "App Name" and the URL as "<https://melaniewalsh.github.io/Intro-Cultural-Analytics/>" for the "App Website URL."
- When you click "Save," you'll be given a series of API Keys: a "Client ID" and a "Client Secret." To generate your "Client Access Token," which is the API key that we'll be using in this file, click "Generate Access Token".
- Copy and paste your "Client Access Token" into genius.py
- Run the genius.py file after the spotipyPop.py file
- It will give you the Python Girls database with 2 more populated tables, GeniusPageViews and GeniusAnnotationCount
 - GeniusPageViews will have song title and page views

- GeniusAnnotationCount will have song title and annotation count
- It will also generate 2 visualizations using matplotlib

Documentation: more clearly explained in terms of lines rather than the input/output of functions because we only have a few functions)

spotifyPop.py (*all explanations of code is also commented in spotifyPop.py file*):

- Lines 5-15: import all necessary files and applications to successfully run spotifyPop.py file
- Lines 19-22: define username of the user we are retrieving spotify data from and finish authorizing use of the data with client id, access token, and redirect uri
- Lines 22-255 **def total_list_songs_popularity:** This function creates the database table by collecting data (top 10 song title and its popularity value) from each of the 10 artist's unique uri. From the database table, it calculates average of the top 10 songs popularity values for each artist and outputs it to a txt file. Then, it creates a list of tuples of artist id, song name and popularity by accessing each artist's top tracks from Spotify. It generates a bar chart that shows the average popularity of each artist's top ten songs.
 - Lines 29-111: assign each artist and artist id value, and create list of tuples of artist id, song name and popularity by accessing each artist's top tracks from Spotify
 - Lines 114-115: create a connection to the PythonGirls database
 - Lines 118-130: create a table called Spotipy in the Python Girls database and input the artist_id, song, and popularity
 - Lines 123-130: if else statement ensures that only 25 items are inserted into Spotipy each time the code is run
 - Lines 136-207: use select statements to get the popularity index for each artist and calculate the average popularity
 - Lines 212-232: open and write to text file to show all of the averages of the top 10 songs from each artist
 - Lines 237-248: create a bar chart that shows the artist name and average popularity for their top 10 songs from Spotify
- Lines 260-392 **def total_list_songs_duration_ms:** This function creates another database table for describing the song duration (in ms) for each song from the top 10 songs of each artist we selected. It then generates a bar chart showing the duration of each song with their respective song title.
 - Lines 261-352: create list of tuples of artist id, song name and duration (in ms) by accessing each artist's top tracks from Spotify
 - Lines 355-356: create a connection to the PythonGirls database
 - Lines 359-371: create a table called Spotipy2 in the Python Girls database and input the artist_id, song, and duration (in ms)

- Lines 364-371: if else statement ensures that only 25 items are inserted into Spotipy2 each time the code is run
- Lines 375-379: append song title and duration for each song into two lists, song_title_list and duration_list
- Lines 382-392: create a bar chart that shows the song title and duration (in ms) for their top 10 songs from Spotify

genius.py (all explanations of code is also commented in genius.py file):

- Lines 5-12: import all necessary files and applications to successfully run genius.py file
- Line 15: generate a client access token to authorize the use of the Genius API
- Lines 20-88: access song and page views for each artist by accessing json data from a search url, appending song title and pageviews to a list to make a list of tuples
- Lines 94-95: create a connection to the PythonGirls database
- Lines 98-110: create a table called GeniusPageViews in the Python Girls database and input the song title and page views
- Lines 103-110: if else statement ensures that only 25 items are inserted into GeniusPageViewstable each time the code is run
- Lines 114-118: append all song titles and page views into two lists, song_title and pageviewsVal_list
- Lines 121-131: create a bar chart that shows the song title and page views for each of the 10 songs from 10 artists we chose from Spotify
- Lines 141-145: create a table called GeniusAnnotationCount in the Python Girls database
- Lines 149-183: gather json data by accessing indices and input the song title and annotation count for each artist
- Lines 187-227: append song title and annotation count for each artist from the json data into two lists, song_title_list and annotation_list
- Lines 232-242: create a bar chart that shows the song title and annotation count for each of the 10 songs from 10 artists we chose from Spotify

Resources

Date	Issue Description	Location of Resource	Result
04/15/22	Trouble with the Spotipy API	https://spotipy.readthedocs.io/en/2.19.0/	We were able to successfully use the Spotipy API and create functions to access top songs of artists we chose
04/16/22	Find Spotipy Client ID and Client Secret	https://cran.r-project.org/web/packages/spotidy/vignettes/Connecting-with-the-Spotify-API.html	We were able to sign up for Spotify Developer and input the Client Secret
04/17/22	Trouble with the Genius API	https://lyricsgenius.readthedocs.io/en/master/ https://melaniewalsh.github.io/Intro-Cultural-Analytics/04-Data-Collection/07-Genius-API.html https://docs.genius.com/#/getting-started-h1	We were able to successfully use the Genius API and get help with accessing the Client Access Token
04/18/22	Search for Spotify feature to scrape from Spotipy API	https://developer.spotify.com/documentation/web-api/reference/#/operations/get-several-tracks https://towardsdatascience.com/extracting-song-data-from-the-spotify-api-using-python-b1e79388d50	We were able to find two indices to scrape from Spotipy API: song popularity and song duration (in ms)
04/19-21/22	How to create databases	file:///Users/emilyyou/Downloads/Databases-v7%20(1).pdf file:///Users/emilyyou/Downloads/DatabaseNormAndJoin-v4.pdf	We were able to follow examples from slides in order to get information into our databases for Spotify and Genius data
04/23/22	How to visualize data	file:///Users/emilyyou/Downlo	We successfully

	using Matplotlib	ads/Matplotlib-v7.pdf	plotted visualizations using matplotlib in Python files
04/25/22	How to add and change colors of matplotlib bar chart	https://matplotlib.org/3.5.0/gallery/color/named_colors.html	Customized colors using this website for each bar in all bar graphs