

SANTA CLARA UNIVERSITY

Department of Computer Science and Engineering

I HEREBY RECOMMEND THAT THE THESIS PREPARED
UNDER MY SUPERVISION BY

Matthew Seminatore, Peter Hay

ENTITLED

Finder: A Diverse Music Recommendation System

BE ACCEPTED IN PARTIAL FULFILLMENT OF THE REQUIREMENTS
FOR THE DEGREE OF

BACHELOR OF SCIENCE
IN
COMPUTER SCIENCE AND ENGINEERING

Thesis Advisor(s) (use separate line for each advisor) date

Department Chair(s) (use separate line for each chair) date

FINDER: A DIVERSE MUSIC RECOMMENDATION SYSTEM

By

Matthew Seminatore, Peter Hay

SENIOR DESIGN PROJECT REPORT

Submitted to

the Department of Computer Science and Engineering

of

SANTA CLARA UNIVERSITY

in Partial Fulfillment of the Requirements

for the degree of

Bachelor of Science in Computer Science and Engineering

Santa Clara, California

2023

FINDER: A DIVERSE MUSIC RECOMMENDATION SYSTEM

Matthew Seminatore

Peter Hay

Department of Computer Science and Engineering
Santa Clara University

Abstract

Many popular media recommendation algorithms suffer from over recommending the most popular results. We wanted to create a music recommendation system that relied completely on song data and no demographic data to provide less biased song results. We created a Python and Flask web application that takes data from Spotify about songs and targets similar songs using the K-Nearest Neighbors algorithm. Our user interface allows users to login to Spotify, select songs as seeds for our algorithm, and then examine data about the recommended results. Our algorithm provided recommendations that were rated on average half as popular as Spotify's algorithm, and provided almost twice as much coverage of long-tail items. The algorithm shows promising results, and with further testing, computing power, and time, it could be improved even more.

Keywords: Long-Tail, K-Nearest Neighbors, Recommendation Algorithms, Music

Table of Contents

Abstract.....	2
Chapter 1 - Introduction.....	5
1.1 Motivation.....	5
1.2 Problem.....	5
1.3 Solution.....	6
Chapter 2 - Requirements.....	8
2.1 Functional Requirements.....	8
2.2 Non-Functional Requirements.....	8
2.3 Design Constraints.....	9
Chapter 3 - Use Cases.....	11
3.1 Use Case Diagram.....	11
3.2 Use Case Description.....	11
Chapter 4 - Technologies Used.....	12
4.1 Functional Technologies.....	12
Backend.....	12
Libraries/API.....	12
Source Control.....	12
4.2 User Interface Technologies.....	13
Chapter 6 - Application Implementation.....	14
6.1 Preprocessing.....	14
Parameter Selection.....	14
6.2 Methods.....	15
Dimensionality Reduction.....	15
K-Nearest Neighbors.....	16
6.3 User Interface.....	17
Chapter 7 - Development Timeline.....	18
7.1 Projected Timeline.....	18
7.2 Problems Encountered.....	18
7.3 Final Timeline.....	19
Chapter 8 - Ethical Analysis.....	20
Chapter 9 - Results.....	21
9.1 Metrics.....	21
Average Recommendation Popularity (ARP).....	21
Average Percentage of Long Tail items (APLT).....	21
9.2 Data.....	21
Table 1. Spotify vs Finder Results.....	22
9.3 Analysis.....	22

Chapter 10 - Conclusion.....	23
10.1 Summary.....	23
10.2 Problems Encountered.....	23
10.3 Lessons Learned.....	24
10.4 Future Work.....	24
Bibliography.....	25
Appendix A - Source Code.....	26
A.1 Server Code.....	26
A.2 Recommender Code.....	28
A.3 Front End.....	31
A.4 Correlation Matrix.....	34
A.5 PCA and TSNE Visualizations.....	35
Appendix B - Correlation Matrix.....	37

Chapter 1 - Introduction

1.1 Motivation

Independent musicians now have access to unparalleled chances due to the music industry's digital transformation. Without relying on record companies, online platforms, social media, and streaming services have given musicians the ability to make, market, and distribute their music worldwide. With their increased freedom, artists may continue to exercise creative control and establish direct relationships with their audience. Independent performers, however, struggle to stand out in the congested music industry. By prioritizing artists signed to established record labels, current music recommendation algorithms aggravate this problem. These methods put commercial concerns and popularity measurements first. Independent musicians thus struggle to gain recognition and establish connections with potential audiences who may value their skill. By creating a music recommendation system that actively boosts the exposure of independent musicians, our research seeks to correct this imbalance. We want to promote an inclusive and democratic music ecosystem by developing a platform that features all musicians. Through the creation of meaningful relationships and the provision of new musical experiences, this system will empower independent musicians and listeners.

1.2 Problem

More and more song recommendations are being generated by algorithms instead of radio DJs. A common problem in recommendation systems is a bias towards popular songs. We call this tendency to recommend popular music popularity bias. The algorithm is more likely to recommend popular music because many people enjoy popular songs. This creates a feedback loop that makes popular music more and more popular. However not everyone likes the most popular songs. Many people have

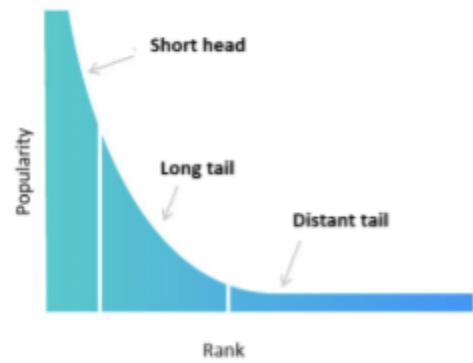


Figure 1. Example of Popularity Bias [1]

diverse and unique music tastes. Songs that are less recommended are called long tail items. From the creator perspective, small, diverse, independent artists with more unique sounds will struggle to get recommended. This creates an environment where it's hard for diverse artists to break out and become more popular. One popular method called collaborative filtering works by comparing people's similar playlists. For example if two people had identical playlists, but one person had one additional song, this song would be recommended to the other person. In Figure 2 below, this would mean that the red books would be recommended to the respective users. This method substantially increases popularity bias, as people are more likely to have popular songs in their playlist.

COLLABORATIVE FILTERING PROBLEM SETTING					
User	Item	<i>1984</i> by George Orwell	<i>Brave New World</i> by Aldous Huxley	<i>On the Road</i> by Jack Kerouac	<i>Of Mice and Men</i> by John Steinbeck
User X		✓	✓	?	✓
User Y		✓	?	✓	✓

altexsoft
software & engineering

Figure 2. Collaborative Filtering Example [2]

1.3 Solution

We wanted to create a recommender system that is less biased towards short head items, and gives more attention to long-tail songs. What we are going to do differently is look at purely the data and provide recommendations to users based off of their data only to decrease as much bias as possible. Using k-nearest neighbors (k-NN) for song data variables can offer a less biased approach compared to collaborative filtering. Collaborative filtering relies on user preferences and recommendations from similar users, which can introduce biases based on demographics, cultural backgrounds, or personal preferences. In contrast, k-NN leverages the similarity of song data variables such as genre, tempo, mood, and instrumentation to determine song recommendations. By considering these objective features, k-NN can provide a more impartial and data-driven approach to recommend songs. This approach ensures that recommendations are

based on similarities in the musical attributes rather than subjective user preferences, ultimately leading to a less biased and potentially more diverse set of song recommendations.

Chapter 2 - Requirements

2.1 Functional Requirements

The goal of the music recommendation system is to deliver a quick and dependable user interface that will improve the user experience and guarantee a seamless experience. The following functional requirements have been identified:

1. Fast and Responsive User Interface: The system should provide a smooth and responsive user interface that provides real-time recommendations while minimizing latency. Users should be able to move about the application, choose music from their favorite songs collection, or add songs one at a time to build a custom playlist to base the recommendations off of. The experience should be user-friendly that allows for easy interaction and intuitive controls.
2. Accurate and Relevant Recommendations: The system should produce recommendations that are precise, pertinent, and customized to the user's preferences. By analyzing the selected songs from the user's liked songs library or the individually added songs, the system should be able to suggest music that aligns with the user's taste. The recommendations should consider factors such as tempo, danceability, energy and many more song attributes.
3. Personalization and Customization: The system should be able to successfully adapt to users changing preferences. It shouldn't save any prior information about users that would affect recommendations. The system should also not base its recommendations off of any demographic data or identifying information that would result in bias.

2.2 Non-Functional Requirements

The music recommendation system has several non-functional needs in addition to its functional requirements, which affect its overall performance, scalability, user experience. The following non-functional requirements have been identified:

1. Performance and Scalability: The system needs to be able to manage a sizable user base and scale efficiently to meet rising demand. In order to quickly produce suggestions, it

should effectively process and evaluate the music that the user has chosen. In order to provide rapid reaction times while adding music or producing suggestions, the interface should be built to reduce latency.

2. Reliability and Availability: The system should be very reliable to guarantee users' continuous availability. It should be robust enough to mistakes and recover quickly if any unanticipated downtime or problems occur. The interface should be designed so that the user can never cause a user generated error to occur.

2.3 Design Constraints

There are several limitations on the design of the music recommendation system that must be taken into account during the development process. The design, features, and general performance of the system are significantly shaped by these constraints. The following design constraints have been identified:

1. Limited Computing Power: The small amount of computer power that is available to operate the project is one of the biggest obstacles. There is a limit on the number of computations that can be done while the system remains responsive because it is being implemented on a standard personal computer rather than a robust server architecture. Due to this restriction, effective algorithms and optimizations are required to keep the suggestion generating process workable within the current computational limits. The algorithms should be created to balance computational complexity and accuracy, allowing the system to provide recommendations quickly without using up all of the available processing power.
2. Data Availability and Quality: The quality and quantity of the data are essential for the music recommendation system to function well. The system needs a large and varied collection of songs, containing song properties. A difficult task is ensuring the availability and continual upkeep of a sizable and high-quality dataset. The system should also be able to gracefully handle missing or partial data, using suitable methods such as data imputation or clever algorithms to lessen the effect of missing data on the recommendation process.

3. User Engagement and Adoption: An important design restriction is how to promote user acceptance of the music recommendation system. Designing an intuitive and user-friendly interface that encourages users to interact with the system and explore the suggested music is crucial for the system's success. The user interface has to be visually appealing, simple to use, and explicit about how to add music and personalize suggestions.

Chapter 3 - Use Cases

3.1 Use Case Diagram

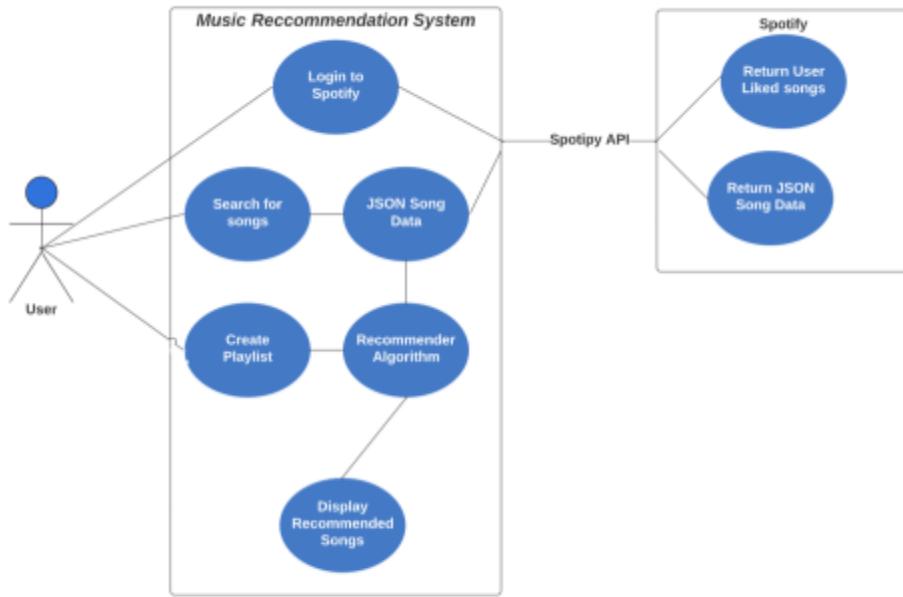


Figure 3. Use Case Diagram

3.2 Use Case Description

When a user opens our web app, the first thing they will see is three columns, the first of which shows a login button. Once a user logs in using their spotify credentials, we use the Spotify API to return the users liked songs in the first column. We also retrieve a JSON that contains the data on song features. Once that has been done the user has the option of selecting seed songs from their liked songs, or searching the entire Spotify catalog for another song of their choice. They select up to 10 songs which will show up in the middle section. Once they are satisfied with their selections they hit the “Create Playlist” button. Our algorithm will retrieve the features of all the seed songs, create a mean vector, and locate the 10 most similar songs in the Spotify catalog. If a user wants to see more information about a song they can simply click on a song and the song features will pop up.

Chapter 4 - Technologies Used

4.1 Functional Technologies

Backend

Python: A high-level, general-purpose programming language that is widely used for backend development. It has an easily readable syntax and allows for quick development and prototyping.

Flask: Flask is a lightweight web framework for Python. It allows developers to build web applications by providing tools and libraries for handling HTTP requests, routing, and rendering HTML templates. Flask enables Python code to dynamically update HTML pages and respond to user interactions.

Libraries/API

TensorFlow: TensorFlow is an open-source machine learning library. It provides a comprehensive ecosystem of tools, libraries, and resources for building and deploying machine learning models.

Spotify: Spotify is a Python library that provides a simple and convenient way to interact with the Spotify Web API. It allows developers to authenticate with Spotify, access user data, search for songs and playlists, and perform various operations on the music streaming platform.

PyPlot: PyPlot is a graphing framework for Python that provides a simple and intuitive interface to create and customize various types of plots, charts, and graphs.

Source Control

GIT: Git is a distributed version control system. It enables multiple developers to collaborate on a project by tracking changes, managing different versions of code, and merging modifications. Git allows developers to work on their local copies of a project, making it easier to manage and share code across teams.

4.2 User Interface Technologies

HTML: HTML is the standard markup language for creating web pages. It provides the structure and content of a webpage, defining elements such as headings, paragraphs, links, images, and tables.

CSS: CSS is a stylesheet language used to describe the presentation and formatting of HTML documents. It allows developers to control the appearance of web pages, defining properties such as colors, fonts, layouts, and animations. CSS enhances the visual design and layout of user interfaces.

JavaScript: JavaScript is a high-level programming language primarily used for adding interactivity and dynamic behavior to web pages. It runs on the client-side, enabling developers to manipulate HTML elements, handle events, and create animations.

Chapter 6 - Application Implementation

6.1 Preprocessing

Parameter Selection

Spotify provides data about each song's musical qualities. Using the Spotify API we can access this data. Shown in Figure (ADD FIGURE NUMBER), for the song "Levitating" by Dua Lipa, the following qualities are given. In doing our data analysis, the first thing we want to remove is non numerical values. They won't show the types of relationships we are looking for. Once parameters were narrowed down to numbers we still had a high number of parameters. With our limited processing power and time we decided to narrow down our parameters more. To try and preserve relationships despite removing variables we used a correlation matrix to determine which variables were most related to each other. A high positive correlation indicates the variables are extremely similar, while a negative correlation indicates they are inverses. A score of 0 is perfectly uncorrelated. We examined the correlations and took the least correlated variables to try and preserve the most relationships possible. The full matrix can be seen in the appendix. The final columns that we ended up selecting were danceability, energy, loudness, speechiness, acousticness, instrumentalness, liveness, valence, and tempo.

Once we have the features selected, we scale them using the standard scalar so that each variable is weighted equally.

#	Column	Non-Null Count	Dtype
0	name	1 non-null	object
1	year	1 non-null	int64
2	explicit	1 non-null	int64
3	duration_ms	1 non-null	int64
4	popularity	1 non-null	int64
5	danceability	1 non-null	float64
6	energy	1 non-null	float64
7	key	1 non-null	int64
8	loudness	1 non-null	float64
9	mode	1 non-null	int64
10	speechiness	1 non-null	float64
11	acousticness	1 non-null	float64
12	instrumentalness	1 non-null	float64
13	liveness	1 non-null	float64
14	valence	1 non-null	float64
15	tempo	1 non-null	float64
16	type	1 non-null	object
17	id	1 non-null	object
18	uri	1 non-null	object
19	track_href	1 non-null	object
20	analysis_url	1 non-null	object
21	time_signature	1 non-null	int64

Figure 4. Example Song Data

6.2 Methods

Dimensionality Reduction

The curse of dimensionality refers to the phenomenon where the increase in the number of dimensions in a dataset leads to various challenges and problems. As the number of dimensions increases, the data becomes increasingly sparse and scattered, resulting in several issues. This makes it difficult to find meaningful patterns or relationships in the data. The amount of processing time also increases exponentially with high dimensional data. Finally, it also becomes difficult to tell distances and separation apart as a human who can only visualize in 3 dimensions. For these reasons we decided that it was necessary to reduce the dimensionality of our data as a part of preprocessing. PCA (Principal Component Analysis) and t-SNE (t-Distributed Stochastic Neighbor Embedding) are both dimensionality reduction techniques commonly used in data analysis and visualization. However, they have distinct differences in their underlying algorithms and outputs.

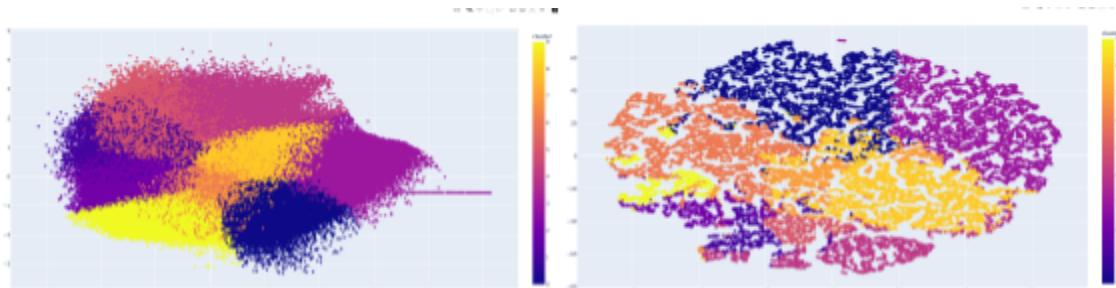


Figure 5. PCA Reduction with K-Means Grouping

Figure 6. TSNE Reduction with K-Means Grouping

PCA is a linear technique that aims to find the directions (principal components) along which the data varies the most. It projects the data onto a lower-dimensional subspace while preserving the maximum amount of variance. PCA is computationally efficient and scales well with the number of features. It is particularly useful when dealing with high-dimensional data and can provide a global overview of the data structure. However, PCA may struggle to capture non-linear relationships or complex patterns in the data.

On the other hand, t-SNE is a non-linear technique that focuses on preserving the local structure of the data. It attempts to map high-dimensional data into a lower-dimensional space by optimizing the similarity between data points. t-SNE is effective at revealing clusters and identifying intricate patterns, but it is computationally intensive and may be more time-consuming, especially with larger datasets. We chose to use PCA due to our computing restraints.

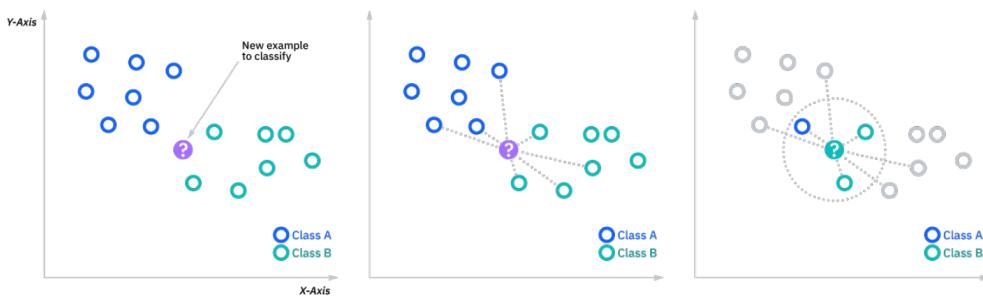


Figure 6. K-Nearest Neighbors Example[3]

K-Nearest Neighbors

Once our data is flattened out, we take the seed songs selected by the user and calculate a mean vector of all of the features for each song. We use a technique called k-Nearest Neighbors to select songs to recommend. This is essentially a method to select the k amount of songs nearest to our mean vector. In our case we use k = 10 to select the 10 most accurate songs. We use a technique called cosine similarity to tell the difference between songs. Cosine similarity finds the other vectors with angles closest to our mean vector. This works better at higher dimensions.

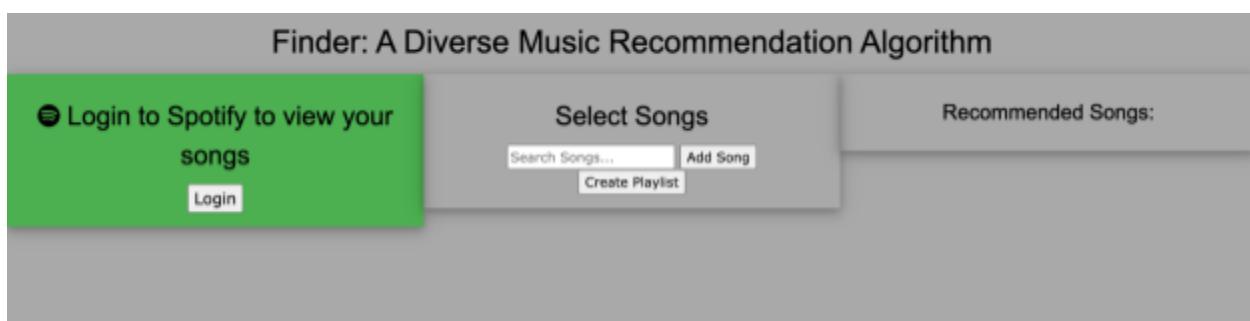


Figure 7. User Interface Layout

6.3 User Interface

When designing the user interface the most important in our minds was simplicity and ease of use. Neither of us had much experience with web technologies. Neither of us are very artistically talented either, so we figured the UI would look better if we kept it simple and avoided any unnecessary designs. We settled on a three column layout. The first is where a user can log in and view their saved Spotify songs. The second is where they can choose the songs they want for their seed list. The third and final column is where recommendations show up after a user has submitted their seeds. Our UI is made using HTML and CSS. We use JavaScript and Python to handle dynamic requests. The Flask framework makes it possible to do all of the backend work in Python and then inject it into the HTML.

Chapter 7 - Development Timeline

7.1 Projected Timeline

	October	November	December	January	February	March	April	May
Enable Spotify Oauth								
Import Song Data								
Create Base Algorithm								
Visualize Song Data and Feature Importance								
Work on UI								
Tuning Algorithm								
Integrate Frontend and Backend								
Testing								
Deployment								

Figure 8. Projected Timeline

7.2 Problems Encountered

Web Design: The team has a background more focused on backend development, which posed challenges in developing the frontend components of the system. Learning new frontend technologies and integrating backend and frontend effectively required overcoming a learning curve.

Learning New Technologies: With limited experience in frontend development and combining backend and frontend components, the team had to acquire new skills and adapt to the challenges of working with different technologies in real time.

Compatibility: Spotify's API did not seamlessly integrate with Jupyter Notebook, the initial choice for development. As a result, the team had to switch to using Git as our source control and development environment to ensure compatibility.

Music Recommendations are Subjective: Music preferences are subjective, and users may have varying preferences. Some users may not find the recommendations generated by the

system aligned with their tastes. Fine-tuning the recommendation algorithms and personalization approaches are necessary to address this challenge.

Large Scale Testing: Processing a significant number of seeds and computations on a personal computer led to longer processing times. In real-world scenarios, such as in a data center, more efficient computational resources would be required for large-scale testing.

Despite the challenges encountered in web design, frontend development, compatibility, subjective music preferences, and large-scale testing, we adapted, learned, refined, and optimized our music recommendation system. These obstacles caused delays, but our team's determination and problem-solving skills allowed us to overcome them. The final timeline demonstrates our adjustments and the successful progress achieved.

7.3 Final Timeline



Figure 9. Actual timeline

Chapter 8 - Ethical Analysis

One of the things that inspired us to pursue this project was the realization that many of the songs we were hearing/being recommended were songs we had already heard before. They were either songs that were currently on the top charts or previously on the charts. We decided that we wanted to create a system that could recommend a more diverse range of songs than we were getting. When we say diverse, we don't mean diverse in the sense of the artist. We mean diverse in the song selections. The way our system promotes diversity of sound is by cutting out many of the factors that could lead to biased recommendations. Our system doesn't look at any demographic or personal information about the user. Instead we trust the user to choose songs that they like as seeds, and then recommend new songs based off of the characteristics of the seeds.

While our system doesn't look at personal information, there are still ethical concerns in our project. Data privacy is becoming a larger and larger issue. Our system does not save any information from a user, and each session starts new with a new login to Spotify. We retrieve all song data from spotify where it is publicly accessible. Artists consent to the song data being published when they publish the songs to Spotify.

Another issue to keep in mind is whether our system is unintentionally introducing bias into recommendations. Since our system works based on mean vectors, recommendations will be in the middle of a user's seed songs. Being in the middle means that recommendations will trend away from outlier songs, and recommendations will be less diverse than desired. In order to avoid this, users should choose seed songs that are similar to each other. If they choose songs that are all different, the mean vector will point towards extremely average featured songs. Another way we could help avoid this is by not recommending mean songs, and instead choosing a song near each of the seeds. This and other algorithm changes are something that we will go into greater detail in the future work section.

Chapter 9 - Results

9.1 Metrics

To help measure the results of our recommendation algorithm, we wanted to collect data about the results. Since our goal is specifically to target better allocations to long-tail items, we chose two data driven metrics.

Average Recommendation Popularity (ARP)

ARP is a simple measurement that looks at the popularity of the items in a song list. It is calculated using the following formula. Given that n is the number of songs, and $Popularity_i$ is the popularity of the i -th song:

$$ARP = \frac{\sum_{i=0}^{numberOfSongs-1} Popularity_i}{numberOfSongs}$$

Average Percentage of Long Tail items (APLT)

APLT is a measurement that looks at the number of results that belong to the long tail. After plotting the distribution of song popularities we found that the long tail fall off was around 60. So we defined the long tail as popularities less than 60. With this in mind,

$$APLT = \frac{songsUnder60}{totalNumberOfSongs}$$

9.2 Data

We based our data collection on 1000 sets of 10 randomly chosen songs. We then used these sets as seeds for the baseline Spotify recommend() function and for the Finder algorithm. We compared the features of the songs recommended by both algorithms. Over the course of 1000 sets we received the following stats:

Table 1. Spotify vs Finder Results

Spotify Baseline	Finder Algorithm
ARP: 68	ARP: 32
APLT: 44%	APLT: 75%

9.3 Analysis

Looking at the data collected, we were happy with the tuning of our algorithm. The ARP for Spotify was twice the ARP of Finder. Looking at the APLT, a little under half of Spotify recommendations were under a popularity of 60. Finder had 3 out 4 songs as long tail items.

We would have liked to use a larger set of testing data, but were limited by computing time. Each set took 5-25 seconds to complete, which ended up being a limiting factor in testing. With a more powerful computer or a more optimized algorithm, we may have been able to do more extensive testing.

When interpreting the results, it is important to remember that music is a personal and subjective topic. It is near impossible for data to perfectly mimic peoples taste in music. Instead, we think that the best way to look at Finder is simply as another alternative to other recommendation algorithms. Spotify does a fantastic job recommending music for many people. Finder just recommends a more diverse set of songs.

Chapter 10 - Conclusion

10.1 Summary

Our group was able to successfully implement a full music recommendation system, including front end user interface. The system allows users to log into their Spotify account, select seed songs of their choice, or search the whole database. After selecting song seeds, our system extracts songs features, creates a mean vector, flattens the dataset, and then calculates the closest songs to the seed. The recommended songs are then displayed on the UI, where the user can click on recommendations to view the songs features. Overall we were very pleased with the final user experience. Our data supported our goal to recommend more long tail popularity songs. We think that our final algorithm worked effectively and can recommend songs that users might enjoy.

10.2 Problems Encountered

We encountered several problems while working on this project. One of the first was limited computing power. We are running our program on a personal laptop as opposed to a commercial server, so there is not as much dedicated computing time. As a result we had to limit the amount of data features that we processed. We tried to select features without much overlapping information, but we still may have lost data relationships. Another result of limited computing was using PCA as opposed to TSNE. We tested TSNE on our dataset and found that it was capable of creating much more separation between songs, however it was too slow for our purposes. A final result of limited computing power is in testing the results of our recommendations. Our system can take anywhere from 5-25 seconds to compute results. When testing different song sets to calculate ARP and APLT we found that testing became extremely long, and for that reason we had to limit our number of random sets to 1000. With a larger server, we would be able to improve our performance and make a more powerful recommender system.

Another problem encountered was limited experience with front end design. It ended up taking much longer than expected to create the UI, as we were trying to learn frameworks while simultaneously prototyping the design. This resulted in a UI that is not quite as polished as we

would have liked. It also delayed many of the other final step tasks such as data collection while we waited on a fully finished design.

10.3 Lessons Learned

One of the most important lessons we learned is that while planning, follow the beginning steps of a plan aggressively. When harder tasks pop up and end up leading to delays, the tasks start to pile up on top of each other. We encountered this while working on the UI. The UI took longer than we expected to get working, which ended up in delaying other important steps of our development such as testing.

10.4 Future Work

The music recommendation system may be improved in a number of ways going forward. Accuracy, scalability, and efficacy across user segments will be improved by ongoing refinement based on user feedback and testing at a greater scale. Users will be more satisfied with the system if we explore new recommendation algorithms, incorporate more data sources, and maybe take into account cutting-edge techniques such as serendipity-based recommendations. A more complete and interesting music discovery platform will result from integration with features like personalized radio stations and curated playlists. The music recommendation system can develop into a more sophisticated and individualized platform, accommodating a variety of musical tastes and promoting music exploration, by taking these future approaches into account.

Bibliography

- [1] Himan Abdollahpouri, Robin Burke, and Bamshad Mobasher. 2019.
- [2] Editor. (2021, July 27). *Recommender systems: Behind the scenes of machine learning-based personalization*. AltexSoft.
- [3] What is the K-nearest neighbors algorithm?. IBM. (n.d.). <https://www.ibm.com/topics/knn>
- Web API | Spotify for Developers. (n.d.). Web API | Spotify for Developers.
<https://developer.spotify.com/documentation/web-api>

Appendix A - Source Code

A.1 Server Code

```
1 import spotipy
2 import json
3 from recommender import *
4 import pandas as pd
5 import urllib.parse
6 import ast
7 #import libraries
8
9
10 #import flask tools and spotipy API
11 from flask import Flask, render_template, redirect, request, jsonify
12 from spotipy import util
13
14 app = Flask(__name__)
15 token = None
16 sp = None
17
18 #read in downloaded track info - presaved for responsiveness
19 data = pd.read_csv("tracks.csv")
20
21 #defined track data class
22 class Track:
23     def __init__(self, name, artist, songID):
24         self.name = name
25         self.artist = artist
26         self.songID = songID
27
28 #root page backend
29 @app.route('/')
30 def index():
31     #determine if logged in
32     login = request.values.get('login', False)
33     savedTracks = request.values.get('savedTracks',None)
34
35     #replace quotes
36     if savedTracks is not None:
37         savedTracks.replace("'", "")
38     #print(savedTracks)
39     savedTracks = ast.literal_eval(savedTracks)
40
41     #return the page with login data
42     return render_template('index.html', login = login, savedTracks = savedTracks)
43
44
45 #page shown once logged in
46 @app.route('/login')
47 def login():
48     redirect_uri = "http://localhost:5001/callback"
49     #username = "12128523243"
50     #spotify login credentials
51     global sp
52     token = util.prompt_for_user_token(scope="user-library-read", client_id="6c37a9c7679445a48df19f5f9692e07a",
53                                         client_secret="8974c498ee1b4648887ce59c5559a3ed", redirect_uri = redirect_uri)
54     sp = spotipy.Spotify(auth = token)
55
56
57     #return redirect('/?login=True')
58
59
60 #@app.route('/savedOrTop')
61 #def savedOrTop():
62
63     selection = request.args.get('savedOrTop')
64     print(selection)
65     #if(selection == "top"):
66     #    #savedTracks = sp.current_user_top_tracks(limit = 10)
67
68     #elif(selection == 'saved'):
69
70     #grab users saved tracks
71     savedTracks = sp.current_user_saved_tracks(limit = 15)
72
73     #can remove - saves tracks
74     f = open("savedTracks.json", "w")
75     str = json.dumps(savedTracks)
76     f.write(str)
77     f.close()
78
79
```

```

80     #format track list
81     print(savedTracks)
82     tracks = {}
83     tracks["items"] = []
84     for item in savedTracks["items"]:
85         trackDict = {}
86         trackDict["name"] = item["track"]["name"].replace('&', '%26')
87         trackDict["artist"] = item["track"]["artists"][0]["name"].replace('&', '%26')
88         trackDict["songID"] = item["track"]["id"]
89         tracks["items"].append(trackDict)
90
91
92     #need at least song name, artist, and id
93     print("\n\n\n\n\n")
94     print("Tracks")
95     print(tracks)
96
97
98
99     return redirect('/?login=True&savedTracks=%s' % tracks)

```

```

102    #generates recommendations based on seedList
103    @app.route('/createPlaylist', methods= ['POST'])
104    def createPlaylist():
105        #python goes here
106        seedList = request.get_json()['seedList']
107
108        """
109        baseline = sp.recommendations(seed_tracks = seedList, limit = 10)
110        print("baseline:")
111
112        sumPop = 0
113        longTail = 0
114        for track in baseline:
115            popularity = track["popularity"]
116            if popularity < 60:
117                longTail +=1
118
119            sumPop += popularity
120
121
122        print("\n\n\n")
123        print(sumPop/10)
124        print(longTail/10)
125        """
126
127
128        seedTracks = []
129        print("\n\n\n\n\n")
130        print("Seedlist")
131        for item in seedList:
132            print(item)
133            seedTracks.append(sp.track(item))
134
135        print("\n\n\n\n\n")
136        print("SeedTracks")
137        print(seedTracks)
138
139        #get seed list, look up song data, pass to recommend songs
140        #format to recommend input dict [{name:'name',year:'year'}]
141

```

```

142        seedListDict = []
143        for track in seedTracks:
144            seedDict = {}
145            seedDict['name'] = track['name']
146            seedDict['year'] = int(track['album']['release_date'][4:])
147            seedListDict.append(seedDict)
148
149        print("\n\n\n\n\n")
150        print("SeedListDict")
151        print(seedListDict)
152
153
154        #call recommender.py here
155        recommendedList = recommend_songs(seedListDict,data)
156
157        print("\n\n\n\n\n")
158        print("recommendedList")
159        print(recommendedList)
160
161
162        #look up recommended list features
163
164        recommendedTracks = []
165        for song in recommendedList:
166            recommendedTracks.append(sp.track(song['id']))
167
168
169
170
171        #print("\n\n\n\n\n")
172        #print("recommendedTracks")
173        #print(recommendedTracks)
174
175
176        #return jsonify(recommendedList)
177        return jsonify(recommendedList)
178
179

```

```

181 #get song data for given songID
182 @app.route('/retrieveSongData', methods= ['POST'])
183 def retrieveSongData():
184     data = request.get_json()
185     songID = data['id']
186
187     songData = sp.audio_features(songID)
188     songData.append(sp.track(songID)[ "name"])
189     songData.append(sp.track(songID)[ "artists"][0][ "name"])
190
191     return jsonify(songData)
192
193
194
195
196
197 #backend for search bar
198 @app.route('/search/', methods = ['GET'])
199 def search():
200     searchTerm = request.values.get('search')
201     result = sp.search(q="track:" +searchTerm, type='track', limit = 1)
202
203     print(result)
204
205     track = result['tracks']['items'][0]
206
207     print("\n \n \n \n \n")
208     print(track['id'])
209
210
211
212
213     return jsonify(track)
214
215 if __name__ == '__main__':
216     app.run(debug=True)

```

A.2 Recommender Code

```

61 # method to find the song data using spotify api using name and year
62 def find_song(name, year):
63     song_data = defaultdict()
64     results = sp.search(q="track: {} year: {}".format(name,year), limit=1)
65     if results['tracks']['items'] == []:
66         return None
67
68     results = results['tracks']['items'][0]
69     track_id = results['id']
70
71     # get all the features of the song using track id
72     audio_features = sp.audio_features(track_id)[0]
73
74     song_data['name'] = [name]
75     song_data['year'] = [year]
76     song_data['explicit'] = [int(results['explicit'])]
77     song_data['duration_ms'] = [results['duration_ms']]
78     song_data['popularity'] = [results['popularity']]
79
80     # create song data dictionary with features and their values
81     for key, value in audio_features.items():
82         song_data[key] = value
83
84     return pd.DataFrame(song_data)
85
86
87
88 # getting the song data from song
89 def get_song_data(song, spotify_data):
90
91     try:
92         song_data = spotify_data[spotify_data['name'] == song['name']].iloc[0]
93         return song_data
94
95     except IndexError:
96         return find_song(song['name'], song['year'])

```

```

99 |     def get_mean_vector(song_list, spotify_data,relevant_cols):
100|         song_vectors = []
101|         # for each song in list
102|         for song in song_list:
103|             # get song data from song
104|             song_data = get_song_data(song, spotify_data)
105|             if song_data is None:
106|                 print('Warning: {} does not exist in Spotify or in database'.format(song['name']))
107|                 continue
108|
109|             # taking relevant columns and forming a vector
110|             song_vector = song_data[relevant_cols].values
111|             # creating a list of song vectors
112|             song_vectors.append(song_vector)
113|
114|         # making a matrix from song vectors
115|         song_matrix = np.array(list(song_vectors))
116|
117|         # returning the mean of the song matrix
118|         return np.mean(song_matrix, axis=0)
119|
120|
121|     """method to flatten out the list of songs into a single dictionary with keys as 'name' and 'year'
122|     and their values as list """
123|     def flatten_dict_list(dict_list):
124|
125|         # creating empty dict with keys as 'name' and 'year'
126|         flattened_dict = defaultdict()
127|         for key in dict_list[0].keys():
128|             flattened_dict[key] = []
129|
130|         # creating lists of values for their keys
131|         for dictionary in dict_list:
132|             for key, value in dictionary.items():
133|                 flattened_dict[key].append(value)
134|
135|         return flattened_dict
136|
137|
138|
139|     def recommend_songs( song_list, spotify_data, n_songs=10):
140|
141|         print("\n\n\n\n ")
142|         print("Song_list")
143|         print(song_list)
144|
145|         scaler = standard_scaler
146|
147|         # flatten the given song list to single dictionary
148|         metadata_cols = ['name', 'artists','id']
149|         song_dict = flatten_dict_list(song_list)
150|
151|         #relevant columns
152|         relevant_cols = ['danceability', 'energy', 'loudness', 'speechiness', 'acousticness',
153|                          'instrumentalness', 'liveness', 'valence', 'tempo']
154|
155|         pca = PCA(n_components=2)
156|
157|         spotify_data_pca = pca.fit_transform(spotify_data)
158|
159|         song_list_pca = pca.transform(song_list)
160|
161|         # scale our whole spotify songs training data
162|         scaled_fit = scaler.fit_transform(spotify_data_pca[relevant_cols])
163|
164|
165|         # Get the mean song vector of the given song list
166|         song_center = get_mean_vector(song_list_pca, spotify_data_pca,relevant_cols)
167|
168|         # scale the mean song vector and reshape it
169|         scaled_song_center = scaler.transform(song_center.reshape(1, -1))
170|
171|         # Calculate cosine distances between the mean song vector and the spotify database
172|         # cosine distances is a measure to check how similar the vectors are
173|         distances = cosdist(scaled_song_center, scaled_fit, 'cosine')
174|
175|         # sort the vectors in ascending format
176|         index = list(np.argsort(distances)[0])
177|
178|         print(index)

```

```

185     # sort the vectors in ascending format
186     index = list(np.argsort(distances)[0])
187
188     print(index)
189
190     # find songs in spotify database from the indices
191     rec_songs = spotify_data.iloc[index]
192     rec_songs = rec_songs[rec_songs['name'].isin(set(song_dict['name']))]
193
194     # filter non-English songs
195     rec_songs = rec_songs[rec_songs['name'].apply(lambda x: filter_non_english(x))]
196
197     calculateMetrics(rec_songs[metadata_cols].to_dict(orient='records')[:10])
198
199     return rec_songs[metadata_cols].to_dict(orient='records')[:n_songs]
200
201
202
203 def calculateMetrics(rec_songs):
204
205     print(rec_songs)
206     sumPop = 0
207     longTail = 0
208     for track in rec_songs:
209         song = sp.track(track["id"])
210         popularity = song["popularity"]
211         if popularity < 60:
212             longTail += 1
213
214         sumPop += popularity
215
216     print("\n\n\n")
217     print(sumPop/10)
218     print(longTail/10)
219
220
221     return

```

A.3 Front End

```

1  <!DOCTYPE html>
2
3
4
5  <head >
6      <title>Music Recommender</title>
7      <meta charset="utf-8">
8      <meta name="viewport" content="width=device-width,initial-scale=1">
9      <link rel="stylesheet" href="https://www.w3schools.com/w3css/4/w3.css">
10     <link rel="stylesheet" href="https://cdnjs.cloudflare.com/ajax/libs/font-awesome/4.7.0/css/font-awesome.min.css">
11     <link rel="stylesheet" href="//code.jquery.com/ui/1.13.2/themes/base/jquery-ui.css">
12     <link rel="stylesheet" href="/resources/demos/style.css">
13     <script src="https://code.jquery.com/jquery-3.6.0.js"></script>
14     <script src="https://code.jquery.com/ui/1.13.2/jquery-ui.js"></script>
15
16
17     <style>
18         h1 {text-align:center;}
19         h2 {text-align:center}
20         form {text-align:center; font-size: 18px;}
21     </style>
22
23 </head>
24
25 <body style="background-color: #darkgray;">
26     <h1 > Finder: A Diverse Music Recommendation System </h1>
27
28     <!--First Third of layout-->
29     <div class = "w3-cell-row">
30         {%if not login%}
31             <div class = "w3-card-4 w3-green w3-cell w3-third w3-text-black">
32                 <div class = "w3-container w3-center w3-padding-16">
33                     <i><2><i class="fa fa-brands fa-spotify"></i> Login to Spotify to view your songs</i></h2>
34
35             <form action = "login">
36                 <input type="submit" value = "Login"/>
37             </form>
38

```

```

39     </div>
40     </div>
41     {%else%}
42         <div class = "w3-third w3-container w3-green w3-cell w3-text-black">
43             <div class = "w3-container w3-center w3-padding-16">
44
45                 <h2><i class="fa fa-brands fa-spotify"></i> Your Saved Tracks:</h2>
46                 <ul class = "w3-ul w3-card-4 w3-hoverable">
47
48                     <!--      <button id = "savedButton" onclick="savedOrTop('saved')">Your Saved Tracks</button>
49                     <button id = "topButton" onclick="savedOrTop('top')">Your Top Tracks</button>
50
51 -->
52             {%for item in savedTracks["items"] %}
53
54                 <li class = "w3-display-container" onclick="addSongFromSpot({{item.songID}},{{item.name}}',
55                         |{{item.artist}})">{{item.name}} - {{item.artist}}</li>
56             {%endfor%}
57
58         </ul>
59     </div>
60     </div>
61     {%endif%}
62
63
64     <!--Second Third of layout-->
65     <div class = "w3-card-4 w3-cell w3-third">
66         <div class = "w3-container w3-center w3-padding-16">
67             <h2>Select Seed Songs</h2>
68
69             <input type = "text" placeholder="Search Songs..." id= "searchbar">
70             <button type = "submit" onclick = "addSong()">Add Song</button>
71
72             <ul id = "list" class = "w3-ul w3-card-4 w3-hoverable">
73
74             </ul>
75
76             <button type="submit" onclick = "submitSeeds()" >Create Custom Playlist</button>
77
78         </div>
79     </div>
80

```

```

89 }      </div>
90 }    </div>
91 }
92 }
93 }
94 }  </div>
95 }
96 <!-- Song data popup -->
97 <div id="popup" class="w3-modal" >
98   <div class="w3-modal-content w3-animate-top" style="width: 600px">
99     <div class="w3-container w3-center">
100       <span onclick="document.getElementById('popup').style.display='none'">&times;</span>
101       <h2></h2>
102       <div id = "song-data" style = "padding-bottom: 16px;" ></div>
103     </div>
104   </div>
105 </div>
106
107
108
109
110 </body>

```

```

141  function addSong(){
142
143   var value = $('#searchbar').val();
144
145
146   $.ajax({
147     url: "/search",
148     method: "GET",
149     data: {search : value},
150     success: function(response){
151       console.log(response)
152
153       var track = response;
154
155       console.log(track)
156
157       var id = track['id']
158       console.log(id)
159
160       var li = document.createElement("li")
161       li.classList.add("w3-container")
162       li.id = id
163       li.textContent = track['name'] + " - " + track['artists'][0]['name']
164       li.addEventListener("click", createEventListener(id));
165
166       var button = document.createElement("button")
167       button.classList.add("w3-button")
168       button.classList.add("w3-right")
169
170       button.innerHTML = "&times;"
171
172       button.addEventListener("click",function(event){
173         event.stopPropagation();
174         li.remove()
175       })
176       li.appendChild(button)
177
178
179       document.getElementById("list").appendChild(li)
180     },
181     error: function(error){
182       console.log(error)
183     }
184   });

```

```

188     function savedOrTop(savedOrTop){
189         fetch('/savedOrTop?savedOrTop=' + savedOrTop, {
190             method: 'GET',
191         })
192     }
193
194
195     function addSongfromSpot(songID,songName,songArtist){
196
197         var id = songID
198         var text = songName + " - " + songArtist
199
200
201         var li = document.createElement("li")
202         li.classList.add("w3-container")
203         li.textContent = text
204         li.id = id
205         li.addEventListener("click", createEventListener(id));
206
207         var button = document.createElement("button")
208         button.classList.add("w3-button")
209         button.classList.add("w3-right")
210
211         button.innerHTML = "&times;"
212
213         button.addEventListener("click",function(event){
214             event.stopPropagation();
215             li.remove()
216         })
217         li.appendChild(button)
218
219         document.getElementById("list").appendChild(li)
220
221     }
222

```

```

223     function submitSeeds(){
224         var seedList = []
225
226         var liElements = document.getElementById("list").getElementsByTagName("li")
227
228         console.log(liElements)
229         console.log(liElements[0].id)
230         for(var i=0;i<liElements.length;i++){
231             seedList.push(liElements[i].id);
232         }
233
234         console.log("hi")
235
236         console.log(seedList)
237
238         fetch('/createPlaylist', {
239             method: 'POST',
240             headers:{
241                 'Content-Type': 'application/json'
242             },
243             body: JSON.stringify({ seedList: seedList })
244         })
245         .then(response => {
246             if (response.ok) {
247                 console.log('SeedList sent successfully!');
248                 return response.json()
249             } else {
250                 console.error('Error sending seedList.');
251             }
252         })

```

```

253         .then(recommendedList => {
254             for(var i = 0; i < recommendedList.length; i++) {
255
256                 var li = document.createElement("li");
257                 li.classList.add("w3-container");
258
259                 var id = recommendedList[i].id;
260                 li.textContent = recommendedList[i].name + " - " + recommendedList[i].artists.replace('[','').replace(']','');
261
262                 li.addEventListener("click", createEventListener(id));
263
264                 document.getElementById("recommended").appendChild(li)
265
266             };
267         });
268
269     })
270     .catch(error => {
271         console.error('Error sending seedList:', error);
272     });
273 }

```

```

277 |     function createEventListener(id){
278 |       return function() {
279 |         fetch('/retrieveSongData', {
280 |           method: 'POST',
281 |           headers: {
282 |             'Content-Type': 'application/json'
283 |           },
284 |           body: JSON.stringify({id})
285 |         })
286 |       .then(response => {
287 |         if (response.ok) {
288 |           console.log('SongID sent');
289 |           return response.json();
290 |         } else {
291 |           console.error('Error retrieving song features.');
292 |         }
293 |       })
294 |       .then(songData => {
295 |         console.log(songData);
296 |         var popup = document.getElementById('popup');
297 |         var title = popup.querySelector('h2');
298 |         title.textContent = songData[1] + ' - ' + songData[2];
299 |
300 |         /*content = JSON.stringify(songData[0], null, 4);*/
301 |         content = songData[0];
302 |
303 |         var relevant_cols = ['danceability', 'energy', 'loudness', 'speechiness', 'acousticness', 'instrumentalness', 'liveness', 'tempo'];
304 |         let html = '';
305 |
306 |         for (const key of relevant_cols) {
307 |           html += `<strong>${key}</strong> ${songData[0][key]}<br>`;
308 |         }
309 |         document.getElementById('song-data').innerHTML = html;
310 |
311 |         popup.style.display='block';
312 |       })
313 |       .catch(error => {
314 |         console.error('Error retrieving song features:', error);
315 |       });
316 |     );
317 |   );
318 | }
319 |
320 | }

```

A.4 Correlation Matrix

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3 import numpy as np
4
5 import pandas as pd
6
7 # Read the CSV file into a DataFrame
8 data = pd.read_csv('tracks.csv')
9
10 # Compute the correlation matrix
11 correlation_matrix = data.corr()
12
13 # Print the correlation matrix
14 print(correlation_matrix)
15
16 sns.heatmap(correlation_matrix, cmap='coolwarm')
17 plt.show()
18
19
20 absolute_corr_values = correlation_matrix.abs().values.flatten()
21
22 sorted_corr_values = sorted(absolute_corr_values, reverse=True)
23
24
25 # Print the ranked variable pairs
26 print(sorted_corr_values)

```

A.5 PCA and TSNE Visualizations

```
70  #graph popularity
71
72  popularity = data[["popularity"]]
73  print(popularity)
74
75  #plt.bar(range(len(popularity)), popularity)
76  plt.plot(popularity)
77  plt.show()
78
79
80  #want to use more variables, but takes a very long time
81  # Use the scatter_matrix() function to create a scatter matrix
82  pd.plotting.scatter_matrix(tracks_data_nums, figsize=(6, 6), diagonal='kde')
83  plt.show()
84
85
86
87
88
89
90  #k means grouping
91  from sklearn.cluster import KMeans
92  from sklearn.preprocessing import StandardScaler
93
94
95
96  # initializing the standard scaler
97  standard_scaler = StandardScaler()
98
99  # normalize data to a 0-1 value
100 tracks_data_normalized = standard_scaler.fit_transform(tracks_data_nums)
101
102 # Initializing and fitting the data on normalized data using no. of clusters as 10
103 kmeans = KMeans(n_clusters=10, random_state=0, verbose = 1).fit(tracks_data_normalized)
104
105 # store the clustered data
106 clustered_data = kmeans.predict(tracks_data_normalized)
107
108
109 #PCA visualization- less compute intensive
110 from sklearn.decomposition import PCA
111
112
113 #n_components may need adjusting
114 pca = PCA(n_components = 2)
115 reduced_data = pca.fit_transform(tracks_data_normalized)
116 print(reduced_data)
117
118 projection = pd.DataFrame(columns=['x', 'y'], data=reduced_data)
119 projection['title'] = tracks_data['name']
120 projection['cluster'] = clustered_data
121
122 #visualize groupings
123
124
125 import plotly.express as px
126
127 fig = px.scatter(projection, x='x', y='y', color='cluster', hover_data=['x', 'y', 'title'])
128 fig.show()
129
130
131
132
133 #tracks_data_normalized_sliced = tracks_data_normalized[0:50000]
134 #TSNE visualization - too slow for now
135 tracks_data_sliced = tracks_data[0:50000]
136 clustered_data_sliced = clustered_data[0:50000]
137 print(tracks_data_sliced)
```

```
138 from sklearn.manifold import TSNE
140
141 # initialize the tsne component using 2 components so as to visualize it using 2d plots
142 tsne = TSNE(n_components=2, verbose=1, n_jobs=-1)
143 # fit and transform the normalized data
144 tracks_embedding = tsne.fit_transform(tracks_data_normalized_sliced)
145 print("tsne complete")
146 # create dataframe containing the x and y coordinates of tsne data
147 projection = pd.DataFrame(columns=['x', 'y'], data=tracks_embedding)
148 projection['tracks'] = tracks_data_sliced['name']
149 projection['cluster'] = clustered_data_sliced
150
151 # visualizing the clustered data
152 fig = px.scatter(projection, x='x', y='y', color='cluster', hover_data=['x', 'y', 'tracks'])
153 fig.show()
154
155
156
157
```

Appendix B - Correlation Matrix

