

# CSE 564: Final Project Prelim Report

By Amirhossein Najafzadeh & Iliya Mirzaei

Spring 2025

## Project Title

**Mood in the Music: A Data-Driven Analysis of How Listening Habits Reflect Mental Health**

## Background

Music is a universal medium for emotional expression and psychological reflection. Across cultures, people rely on music not only for entertainment, but also as a tool for emotional regulation and mental well-being. Recent interdisciplinary research, covering psychology, neuroscience, musicology, and data science, has increasingly explored this connection, investigating how music affects mood and supports mental health interventions [1].

Empirical studies suggest that music listening can alleviate symptoms of anxiety, depression, and insomnia, and support emotional processing in clinical settings through music therapy [2, 3]. Furthermore, individuals often gravitate toward particular genres, tempos, or lyrical content based on their psychological state, highlighting the bidirectional relationship between music preference and mental health.

Spotify and Apple Music are two of the most popular music streaming apps that people use to listen to songs every day. Because they are widely used around the world, these platforms offer a huge and diverse collection of songs from many different genres — like pop, rock, jazz, classical, and more. What makes these platforms especially useful for our analysis is that they also provide detailed attributes for each song.

This project aims to uncover patterns and correlations between self-reported mental health conditions and the characteristics that people prefer to listen to music. By analyzing large-scale datasets of music preferences, audio features, and mental health survey responses, this work seeks to identify interpretable trends that may inform mental health assessment tools, therapeutic music selection, and personalized wellness recommendations.

## Original Datasets and Feature Selection

This project integrates multiple datasets related to music, genre classification, listener behavior, and song-level attributes. The datasets serve complementary purposes and vary in size and structure. Table 3 provides an overview of each dataset, including its intended use, number of records (items), total features, and the number of features selected for analysis.

Name	Type	Usage	Entities	Features
apple_music_dataset	CSV	Songs to Genres	10,000	24
mxmh_survey_results	CSV	Genres to Mental Health	1,000	34
spotify_2000_tops	CSV	Songs to Attributes	2,000	15
spotify_music_dataset	CSV	Songs to Attributes	2,000	18
spotify_song_attributes	CSV	Genres to Attributes	10,080	22
universal_top_spotify_songs	CSV	Country to Listening Behavior	22,000+	24

Table 1: Overview of the Original Datasets

### Feature Selection and Filtering

In the first step of the preprocessing pipeline, relevant features were selected from each dataset to reduce noise and ensure alignment with the analytical objectives. The filtering process involved selecting key columns and saving the reduced datasets to a temporary directory (`tmp/`). Table 3.1 provides the details of this filtering.

Dataset	Retained	Columns
apple_music_dataset	2	'trackCensoredName', 'primaryGenreName'
spotify_2000_tops	11	'Title', 'Top Genre', 'Beats Per Minute (BPM)', 'Energy', 'Danceability', 'Loudness (dB)', 'Speechiness', 'Acousticness', 'Liveness', 'Valence', 'Popularity'
spotify_music_dataset	13	'song', 'genre', 'popularity', 'danceability', 'energy', 'loudness', 'mode', 'speechiness', 'acousticness', 'instrumentalness', 'liveness', 'valence', 'tempo'
spotify_song_attributes	12	'trackName', 'genre', 'danceability', 'energy', 'loudness', 'mode', 'speechiness', 'acousticness', 'instrumentalness', 'liveness', 'valence', 'tempo'
mxmh_survey_results	11	'Age', 'Hours per day', 'While working', 'Fav genre', 'Exploratory', 'Foreign languages', 'Anxiety', 'Depression', 'Insomnia', 'OCD', 'Music effects'
universal_top_spotify_songs	5	'country', 'tempo', 'danceability', 'energy', 'valence'

Table 2: Selected Features from Original Datasets

This filtering step significantly reduced the dimensionality of the original data, preparing the datasets for the downstream steps of transformation, integration, and analysis.

# Data Wrangling

## Decoupling

Several datasets in our project contain a **genre** field in which individual entries are associated with multiple genres. These multi-value genre features are often represented as comma-separated strings within a single cell, such as:

Rock, Alternative, Indie

To enhance the accuracy and granularity of our analysis, we performed a decoupling step to normalize these values. Each genre was extracted and assigned to its own row, effectively transforming the dataset into a one-to-many relationship between songs and genres. This approach enables a more detailed understanding of genre-specific patterns and ensures compatibility with analytical models that expect one genre per instance.

The following table summarizes the dataset sizes before and after genre expansion:

Dataset	Before Expansion	After Expansion
Spotify Dataset	2,000	3,704
Apple Music Dataset	10,000	11,629

Table 3: Dataset Size Before and After Genre Field Expansion

This preprocessing step was crucial in allowing genre-based clustering, statistical aggregation, and visualization techniques to function correctly, without being misled by overlapping or conflated genre labels.

## Normalization

Following the decoupling step, we applied normalization to all numerical features across our datasets. This process was essential to ensure that features with varying scales could be meaningfully compared and analyzed. Several numerical attributes, such as **loudness**, had values ranging approximately between  $-40$  and  $40$ , while others, like **danceability**, ranged between  $0$  and  $2$ . Without normalization, features with larger numeric ranges would dominate distance-based metrics and skew the results of clustering or dimensionality reduction techniques. To address this imbalance, we normalized each numerical feature to a common scale, typically using Min-Max normalization:

$$x_{\text{normalized}} = \frac{x - \min(x)}{\max(x) - \min(x)}$$

This transformation maps each value to a range between  $0$  and  $1$ , allowing for consistent comparison across features. By normalizing the data, we established a consistent foundation for further statistical analysis, clustering, and machine learning tasks within the visualization dashboard.

## Fixing Inconsistencies

Before merging our datasets, we noticed that genres were written differently across them. For example, one dataset had **Rock** and another had **Rock and Roll**. To fix this, we used natural language processing and cosine similarity to compare genre names. This helped us group similar genres under one common name.

## Country Inference from Audio Features

The `mxmh_survey_results` dataset does not include a location field, which limits the ability to analyze trends geographically. To resolve this, we leveraged the `universal_top_spotify_songs` dataset, which includes listening trends in 73 countries.

We calculated the average `tempo`, `danceability`, `energy`, and `valence` for each country. Then, for each respondent in the mental health dataset, we computed the Euclidean distance between their BPM (and imputed values for other audio features) and each country's feature vector. The country with the smallest distance was assigned to that respondent. This process effectively mapped user behavior to the closest global listening profile.

This geographic inference process enables the integration of country-level visualizations such as choropleth maps and regional comparisons, enriching the storytelling capacity of the final dashboard.

### `apple_music_dataset`

- Total changes: 7,674
- Average best similarity score: 0.6813
- Top 5 genres changed in this dataset following by the number of changes:
  - Hip-Hop: 1,038
  - Pop: 2,223
  - Alternative: 1,111
  - Soundtrack: 688
  - Dance: 157

### `spotify_2000_tops`

- Total changes: 1,994
- Average best similarity score: 0.6782
- Top 5 genres changed in this dataset following by the number of changes:
  - adult standards: 123
  - album rock: 413
  - alternative hip hop: 2
  - alternative metal: 70
  - classic rock: 51

### `spotify_music_dataset`

- Total changes: 3,704
- Average best similarity score: 0.7611
- Top 5 genres changed in this dataset following by the number of changes:
  - pop: 936
  - rock: 162
  - Pop: 697
  - hip hop: 776
  - R&B: 439

### `spotify_song_attributes`

- Total changes: 8,580
- Average best similarity score: 0.4707

- Top 5 genres changed in this dataset following by the number of changes:
  - alt z: 656
  - pop: 602
  - dance pop: 172
  - alternative metal: 150
  - singer-songwriter pop: 164

### **universal\_top\_spotify\_songs**

- Total records analyzed: 22,000+
- Unique Countries in the result: 38
- Average features used for matching: tempo, danceability, energy, valence

Thanks to our language-based method for fixing inconsistencies, we were able to easily find and remove duplicates. Also, by getting rid of outliers and inconsistent entries, the results of our merging in the next step became more accurate.

### **Data Merging**

In this step, we combined all the datasets into one complete dataset to prepare the data for analysis. We used song names to match songs with their attributes, and then used genres to bring together all songs, attributes, and mental health data. Additionally, to enable geographical visualizations, we used the listening behavior dataset from Spotify to infer countries based on musical feature similarity.

- **Step 1: Merging Songs and Genres**
  - Songs in Apple Music dataset: 11,629
  - Songs in Spotify dataset: 3,704
  - Songs after merging: 3,488
- **Step 2: Merging Songs with Their Attributes**
  - Songs in merged dataset: 3,488
  - Songs in Spotify 2000 dataset: 1,994
  - Songs in Spotify attributes dataset: 10,080
  - Songs after merging: 10,522
- **Step 3: Merging with Mental Health Data (Using Genre)**
  - Songs with attributes: 10,522
  - Mental health records: 1,000
  - Final merged records: 520,778
- **Step 4: Country Inference via Audio Similarity**
  - Used Spotify’s country-level top song data
  - Extracted average tempo, danceability, energy, and valence per country
  - Matched each respondent to closest country using Euclidean distance based on BPM and placeholder audio features
  - Result: Mental health records enriched with `country` column

In the end, our final dataset includes **23 features** and **520,778 records**. This rich dataset allows us to explore the connection between music attributes and mental health issues.

## Filling Missing Values

After creating the final dataset, some values were missing because of the merging process. We cleaned the data by removing duplicate rows and dropping any rows that were completely empty. To fill in the remaining missing values, we used linear regression to make educated guesses based on the patterns in the data.

## Sampling

In the final step of our data processing, we aimed to reduce the dataset size from over half a million entities to 2000 data items. To achieve this, we used a clustering method. First, we applied Principal Component Analysis (PCA) to select the best features for clustering. Then, we used the K-means algorithm along with the knee finding method to group the data into 10 clusters. Finally, we selected around 2000 samples based on the size of each cluster to create our final dataset. In the next figures you can see the results of our clustering based on PCA1 and PCA2.

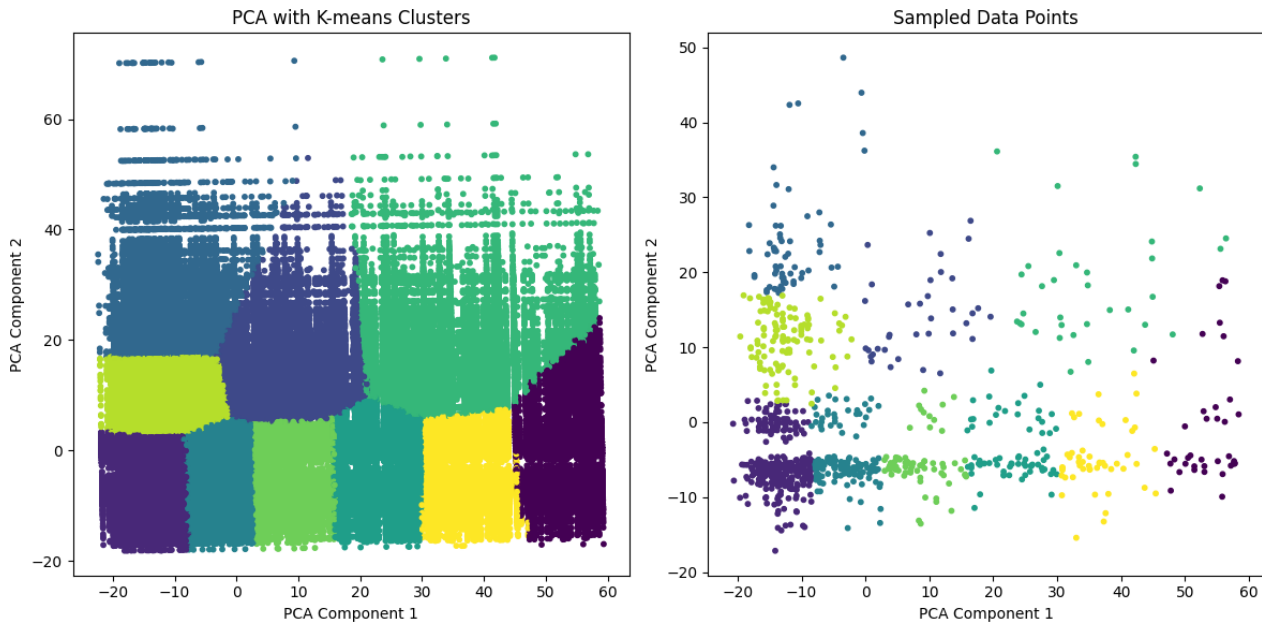


Figure 1: PCA results with 1000 samples

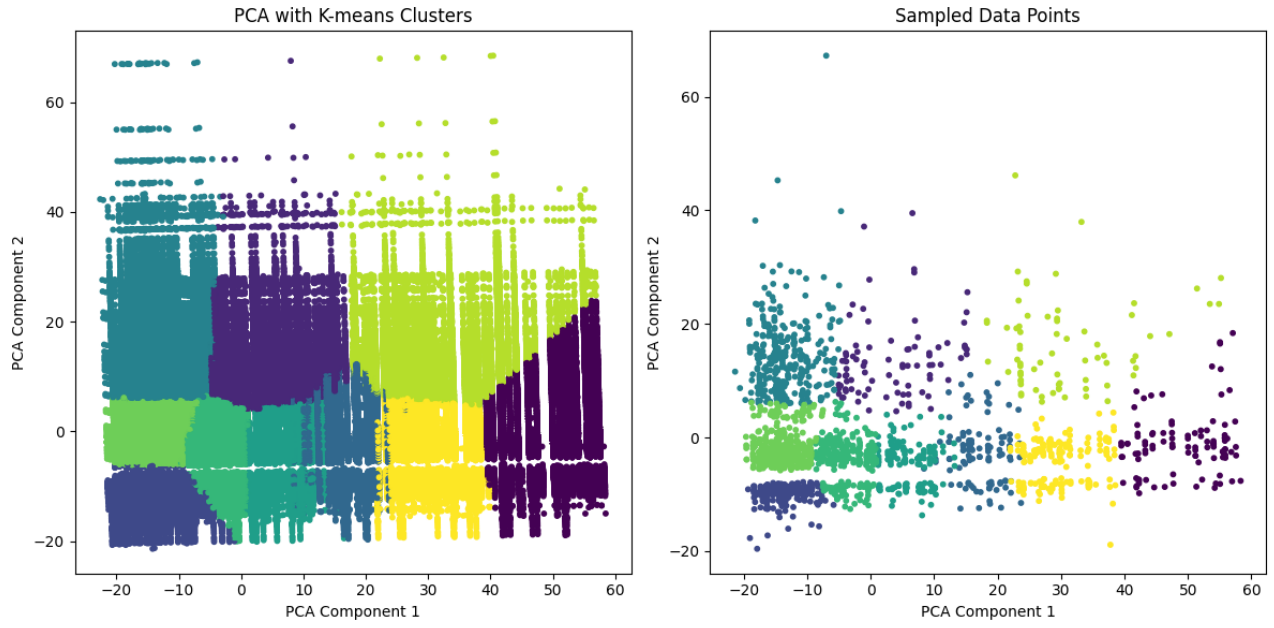


Figure 2: PCA results with normalized data and 2000 samples

## Dataset

The final dataset consists of 23 columns and 1995 rows. The columns contain various features related to music attributes and mental health data. Below are the unique values for each column:

Column Name	Number of Unique Values
acousticness	99
bpm	61
danceability	85
energy	98
genre	11
instrumentalness	63
liveness	90
loudness	651
mode	2
popularity	84
speechiness	63
tempo	68
valence	96
age	70
hours	140
while working	2
exploratory	2
foreign	2
anxiety	73
depression	72
insomnia	69
ocd	60
effects	3
country	38

Table 4: Unique Values in Each Column

The following are the top 5 rows of the dataset, showing a snapshot of the music and mental health attributes:

acousticness	bpm	danceability	energy	genre	instrumentalness	liveness	exploratory	foreign	anxiety	depression	insomnia	ocd	effects	country
0.99	0.65	0.40	0.18	K pop	0.00	0.09	Yes	No	0.45	0.74	0.97	0.17	Improve	PK
0.00	0.65	0.51	0.89	Rock	0.01	0.23	Yes	No	0.19	0.56	0.59	0.07	Improve	IT
0.07	0.65	0.82	0.72	Rock	0.16	0.07	Yes	No	0.49	0.63	0.74	0.30	Improve	PK
0.71	0.67	0.61	0.45	Rock	0.00	0.13	Yes	No	0.90	0.80	0.94	0.37	Improve	CL
0.57	0.65	0.24	0.54	Hip hop	0.37	0.09	No	Yes	0.71	0.42	0.17	0.18	No effect	UY

Table 5: Top 5 Rows of the Dataset



This dataset includes various features such as *acousticness*, *energy*, and *genre*, alongside mental health attributes like *anxiety*, *depression*, and *insomnia*. The data is useful for analyzing the relationship between music characteristics and mental health issues.

## Methods

To build our dataset, we utilized Python and shell scripts. For dataset reading and manipulation, we leveraged the `Pandas` library. For linguistic processing, we used the following imports from the `scikit-learn` library:

- `from sklearn.feature_extraction.text import TfidfVectorizer`
- `from sklearn.metrics.pairwise import cosine_similarity`

For Principal Component Analysis (PCA) and clustering, we used the following modules from `scikit-learn`:

- `from sklearn.decomposition import PCA`
- `from sklearn.cluster import KMeans`

To assign country labels based on proximity in audio feature space, we computed distances using:

- `from scipy.spatial.distance import cdist`

## Next Steps

We started by setting up two main projects: one using Python Flask and the other using React.js. The Python Flask application is responsible for handling the backend tasks. It reads the dataset and processes the data, then provides the necessary information to our frontend application through a web API. On the frontend, we used React.js to build the user interface. In React, we fetch data from the web API that we created with Flask. We then use D3.js, a powerful JavaScript library, to visualize and plot the data on a dashboard, allowing users to interact with and explore the information in a more intuitive way.

## References

- [1] Accelerated construction of stress relief music datasets: National Library of Medicine - <https://pmc.ncbi.nlm.nih.gov/articles/PMC11125514>
- [2] Acoustic Sounds for Wellbeing; A Novel Dataset and Baseline Results: Cornell University - <https://arxiv.org/abs/1908.0167>
- [3] Harnessing Social Music Tags for Characterizing Depression Risk: Cornell University - <https://arxiv.org/abs/2007.13159>