

MUSIC RECOMMENDATION DATABASE

PROF. Nick Brown

Contributors:

Abhishek Hegde [NU ID: 002744522]

Lokesh Jeswani [NU ID: 002795957]

Team: "Music Junkies"

Abstract:

A Music Database is created with carefully chosen entities such as songs, albums, genres, artists, and record labels. The project has a command-line interface that offers database capabilities like adding and deleting data, altering values, sorting by an attribute, building views, and performing various queries. Choose any song based on your preference or mood. The database will display a playlist, or a list of recommended songs based on the genre, song or artist chosen. Here, we've taken datasets from the "Spotify" music service, kept the pertinent entities that the database schema required, and imported them.

Introduction:

Many music listeners have turned to listen to online music. The database technology has made it possible that music listeners could get access to music as they want. Online service of music subscription has gained immense popularity in the era of cloud computing. The advancement of cloud techniques eases users to get access to an unlimited number of songs. Here we will be creating "Music Recommendation Database", which will provide recommendations based on the genre and various other parameters of the song, which is currently being played.

Music Recommendation System is used to recommend songs based on genre. The system allows users to create playlists based on genre and popularity whenever they are logged in. Recommendations are also made based on genre and various other relational entities. The system can provide instantaneous recommendations for each song which is being played.

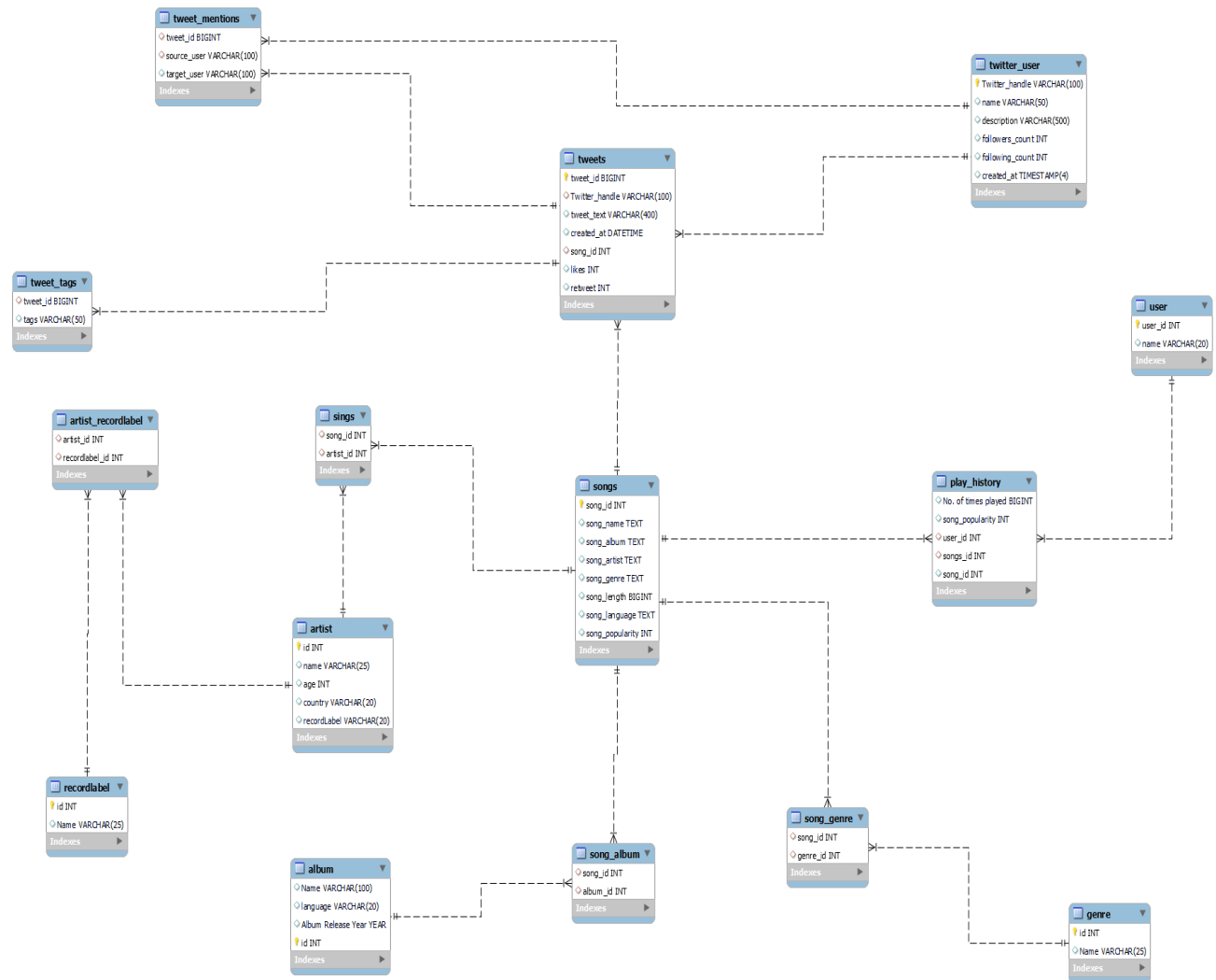
Use-case for un-authorized users:

1. Search songs and albums through database using input keywords.
2. View recently added albums and its details.
3. View custom music playlist and get recommendations based on it.

Use-case for authorized users(admin):

1. Can add and remove albums from the playlists and database.
2. Add new music/albums to the database.
3. The currently authorized user or administrator can see the updates to the uploaded data.

ER Diagram:



Snippets from Database created:

Songs Table:

	song_id	song_name	song_album	song_artist	song_genre	song_length	song_language	song_popularity
▶	1	Beat It	Thriller 25 Super Deluxe Edition	Michael Jackson	Pop	3	English	81
	2	Smooth Criminal - 2012 Remaster	Bad 25th Anniversary	Michael Jackson	Pop	3	English	79
	3	Gimme More	Blackout	Britney Spears	Pop	3	English	78
	4	Eye of the Tiger	Rocky IV	Survivor	Rock	3	English	76
	5	Everybody Wants To Rule The World	Songs From The Big Chair (Super Deluxe Edition)	Tears For Fears	R&B	3	English	86
	6	Everybody Talks	Picture Show	Neon Trees	R&B	2	English	79
	7	Earned It (Fifty Shades Of Grey)	Beauty Behind The Madness	The Weeknd	Pop	3	English	74
	8	In Da Club	Get Rich Or Die Tryin'	50 Cent	Hip-Hop	2	English	82
	9	Black and Yellow	Rolling Papers	Wiz Khalifa	Rap	3	English	75
	10	Superman	The Eminem Show	Eminem, Dina Rae	Hip-Hop	4	English	0
	11	In Da Club	Get Rich Or Die Tryin'	50 Cent	Hip-Hop	2	English	75

Album Table:

Name	language	Release_year	id
Thriller 25 Super Deluxe Edition	English	1982	1
Bad 25th Anniversary	English	1987	2
Blackout	English	2007	3
Rocky IV	English	1905	4
Songs From The Big Chair (Super Deluxe Edition)	English	1985	5
Picture Show	English	2012	6
Beauty Behind The Madness	English	2015	7
Get Rich Or Die Tryin'	English	2003	8
Rolling Papers	English	2011	9
The Eminem Show	English	2002	10
Get Rich Or Die Tryin'	English	2003	11

Artist Table:

id	name	age	country	recordLabel
1	Michael Jackson	50	USA	Sony
2	Britney Spears	30	USA	Warner
4	Tears For Fears	34	Ireland	CBS
5	Neon Trees	23	Germany	Universal
6	The Weeknd	35	USA	BMG
7	50 Cent	36	Mexico	ARMIND
8	Wiz Khalifa	32	USA	Sony
9	Eminem	33	USA	Warner
10	C. SHIROCK	28	Ireland	CBS
11	Post Malone	25	USA	Universal
12	DJ Snake	26	USA	BMG

Genre Table:

id	Name
1	Pop
2	Rock
3	R&B
4	Hip-Hop
5	Rap
6	Electronic
7	Funk
8	Bollywood
9	Anthem
10	House
11	Club

Played history Table:

Played_count	song_popularity	song_id
4680	81	1
2055	79	2
2397	78	3
7647	76	4
9455	86	5
2635	78	6
2554	74	7
7498	82	8
7111	75	9
2652	75	11
1283	80	12

RecordLabel Table:

id	Name
1	Sony
2	Warner
3	CBS
4	Universal
5	BMG
6	ARMIND
7	T-series
8	RCA
9	Epic
10	Atlantic
11	Republic

Below are the tables created using foreign keys of two tables.

Sings Table – shows which artist is the contributor to which song.

song_id	artist_id
1	1
2	1
3	2
5	4
6	5
7	6
8	7
9	8
11	7
13	10
17	14

- One artist can have many songs, for ex, artist with id 1 have sung songs with id 1 and 2.

Song-album Table – [shows which song belongs to which album].

song_id	album_id
1	1
39	2
38	2
2	2
3	3
4	4
5	5
6	6
7	7
76	8
11	8

- *One album can have multiple songs in it, for ex, album_id 2 has song_id's 2,38 and 39.*

Artist_recordlabel Table - [connects artist and recordlabel tables]

artist_id	recordlabel_id
1	1
2	2
4	3
5	4
6	5
7	6
8	1
9	2
10	3
11	4
12	5

- *Each artist can have songs released under multiple record labels.*

SNAPSHOT OF VIEWS AND USE CASES:

- 1) List few songs and its artist: Eminem with album name 'Recovery' and release year '2010'.

```
CREATE VIEW Eminem_2010 AS
select distinct a.song_name, a.song_album, b.name, c.Release_year from songs a
JOIN artist b on b.name = a.song_artist
JOIN album c on c.name = a.song_album
where b.name = "Eminem" and c.Release_year="2010";
```

```
select * from Eminem_2010;
```

	song_name	song_album	name	Release_year
▶	Not Afraid	Recovery	Eminem	2010

- 2) Display few artists and their age whose name starts with 'L' and age is above 40?

```
CREATE VIEW L40 AS
select distinct a.name, a.age from artist a Join songs b
on a.name = b.song_artist
where a.name like 'L%' and age>40;
```

```
select * from L40;
```

	name	age
▶	Laura Branigan	67
	Lionel Richie	48

- 3) Show some popular bollywood songs sung by 'Arijit Singh' with least popularity ratings.

```
CREATE VIEW less_popular_arijit_singh_songs AS select distinct
a.song_name,a.song_artist,a.song_popularity from songs a join artist b on b.name =
a.song_artist join album c on c.name = a.song_album where b.name = 'Arijit Singh' and
a.song_genre = 'Bollywood' order by song_popularity ASC;
```

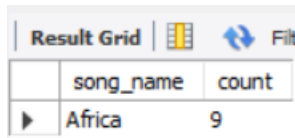
```
Select * from less_popular_arijit_singh_songs;
```

	song_name	song_artist	song_popularity
▶	Palat - Tera Hero Idhar Hai	Arijit Singh	50
	Kabhi Jo Baadal Barse	Arijit Singh	59

- 4) List the songs with the least number of tweets between 2022-11-1 between 2022-11-12.

```
CREATE VIEW least_tweet_songs AS
Select m.song_name, count(t.song_id) as count from songs m inner join tweets t on
t.song_id=m.song_id group by m.song_name
order by count limit 1;
```

```
select * from least_tweet_songs;
```

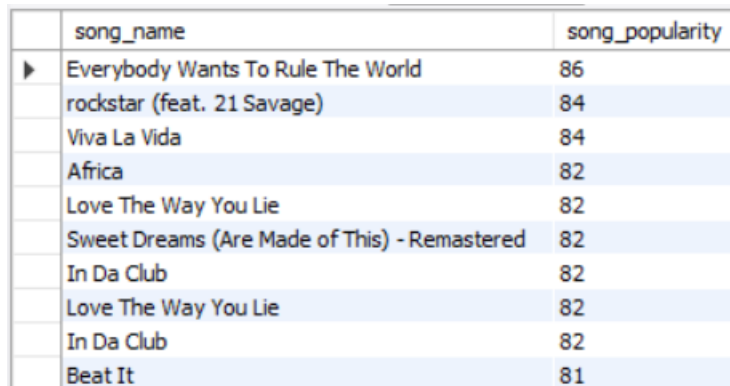


	song_name	count
▶	Africa	9

- 5) What are the top 10 songs with highest popular ratings and in English language?

```
CREATE VIEW top10_popular_songs AS
select song_name,song_popularity from music.songs
where song_language = 'english' order by song_popularity desc limit 10;
```

```
select * from top10_popular_songs;
```



	song_name	song_popularity
▶	Everybody Wants To Rule The World	86
	rockstar (feat. 21 Savage)	84
	Viva La Vida	84
	Africa	82
	Love The Way You Lie	82
	Sweet Dreams (Are Made of This) - Remastered	82
	In Da Club	82
	Love The Way You Lie	82
	In Da Club	82
	Beat It	81

- 6) How many unique artists are there in the database where genre is 'Pop'?

```
CREATE VIEW pop_genre_artists AS
select distinct name from artist a JOIN songs b
on a.name = b.song_artist where song_genre = 'Pop';
```

```
select * from pop_genre_artists;
```

	name
▶	Michael Jackson
	Britney Spears
	The Weeknd
	Dua Lipa
	Laura Branigan
	TOTO
	George Michael
	Whitney Houston
	Madonna
	Ricky Martin

- 7) Name few artists whose name starts with 'M' and age is above 30?

```
CREATE VIEW artistM30 AS
```

```
select distinct a.name from artist a join songs b on a.name = b.song_artist  
where a.name like 'M%' and age>30;
```

```
select * from artistM30;
```

	name
▶	Michael Jackson
	Madonna

- 8) List songs with rating above 80 and release year between 2000 and 2015.

```
CREATE VIEW high_rating80 AS
```

```
select a.song_name from songs a JOIN album b  
on a.song_album = b.name where song_popularity > 80  
and b.Release_year between '2000' and '2015';
```

```
select * from high_rating80;
```

	song_name
▶	In Da Club
	Love The Way You Lie
	Yeah! (feat. Lil Jon & Ludacris)
	Toxic
	Promiscuous
	Viva La Vida

- 9) How many unique albums does artist 'Michael Jackson' have in the database?

```
CREATE VIEW mj_unique_albums AS  
select count (distinct song_album) as unique_albums from songs  
where song_artist = 'Michael Jackson';
```

```
select * from mj_unique_albums;
```

	unique_albums
▶	3

- 10) Return songs sung by Arijit Singh with record label 'T-series' and release year between 2013 and 2015.

```
CREATE VIEW arijit_singh AS  
select a.song_name, b.recordLabel, c.Release_year from songs a JOIN artist b  
on b.name = a.song_artist JOIN album c on c.name = a.song_album  
where b.name = 'Arijit Singh' and b.recordLabel = 'T-series'  
and c.Release_year between '2013' and '2015';
```

```
select * from Arijit_singh;
```

	song_name	recordLabel	Release_year
▶	Kabhi Jo Baadal Barse	T-series	2013
	Palat - Tera Hero Idhar Hai	T-series	2014

- 11) Show all songs in the database where Eminem has been a contributor.

```
CREATE VIEW eminem_songs AS  
select distinct song_name as unique_contributions from songs  
where song_artist like '%Eminem%';
```

```
select * from Eminem_songs;
```

	unique_contributions
▶	Superman
	Not Afraid
	Love The Way You Lie
	The Monster
	Venom - Music From The Motion Picture
	Rap God

12) List unique songs with rating less than 80 and release year is 1987.

```
CREATE VIEW 1987_80less AS
select distinct a.song_name, b.release_year from songs a JOIN album b
on a.song_album = b.name
where song_popularity < 80 and b.Release_year = '1987';

select * from 1987_80less;
```

	song_name	release_year
▶	Smooth Criminal - 2012 Remaster	1987
	Man in the Mirror - 2012 Remaster	1987
	Dirty Diana - 2012 Remaster	1987

13) What are the top 10 songs with highest popular ratings and artist name starts with 'M' and in English language?

```
CREATE VIEW 10m_popular_songs AS select song_name, song_popularity, song_artist from
music.songs where song_language = 'english' and song_artist like 'M%' order by song_popularity
desc limit 10;
```

```
Select * from 10m_popular_songs;
```

song_name	song_popularity	song_artist
Beat It	81	Michael Jackson
Smooth Criminal - 2012 Remaster	79	Michael Jackson
La Isla Bonita	73	Madonna
4 Minutes (feat. Justin Timberlake & Timbaland)	70	Madonna, Justin Timberlake, Timbaland
Man in the Mirror - 2012 Remaster	69	Michael Jackson
Billie Jean	63	Michael Jackson
Beat It	61	Michael Jackson

14) List few songs and its artist: Michael Jackson with album name 'Thriller' and release year '1982'.

```
CREATE VIEW MJ1982Thriller AS
select distinct a.song_name, a.song_artist, a.song_album, c.Release_year from songs a
JOIN artist b on b.name = a.song_artist
JOIN album c on c.name=a.song_album
where b.name = 'Michael Jackson' and c.name = 'Thriller' and c.Release_year = '1982';

select * from MJ1982Thriller;
```

	song_name	song_artist	song_album	Release_year
►	Billie Jean	Michael Jackson	Thriller	1982
	Beat It	Michael Jackson	Thriller	1982

- 15) What are the top 5 artists of all time, whose song_language is english and genre = 'Pop' based on highest popular ratings?

```
CREATE VIEW top5_artists AS
select distinct song_artist from music.songs
where song_language = 'English' and song_genre = 'Pop'
order by song_popularity desc limit 5;
```

Select * from top5_artists;

	song_artist
►	Coldplay
	Eurythmics, Annie Lennox, Dave Stewart
	TOTO
	Michael Jackson
	Nelly Furtado, Timbaland

- 16) Display song name and its popularity where the number of times song played is greater than 9000.

```
CREATE VIEW most_played AS
Select distinct a.song_name, a.song_popularity from songs a
JOIN play_history b on a.song_id = b.song_id
Where b.played_count>9000;
```

Select * from most_played;

	song_name	song_popularity
►	Everybody Wants To Rule The World	86
	Summer Of '69	81
	...Baby One More Time	78
	4 Minutes (feat. Justin Timberlake & Timbaland)	70
	Samjhawan	64
	Ilahi	67
	Love The Way You Lie	82
	Breakaway	64

- 17) Find a tweet which mentions a song whose artist name is 'eminem'.

```
CREATE VIEW eminem_songs AS
SELECT a.tweet_text,b.song_name from music.tweets a
join music.songs b on a.song_id = b.song_id
where a.tweet_text like '%Eminem%';

select * from eminem_songs;
```

	tweet_text	song_name
▶	o tanto que eu ouço essa música até hoje e até...	Love The Way You Lie
	At a winery with live music...this group is doing ...	Love The Way You Lie
	Eminem - Love The Way You Lie ft. Rihanna htt...	Love The Way You Lie

- 18) Search for a song from our database and find it in twitter tweets and show the number of re-tweets for that tweet.

```
CREATE VIEW song_retweets AS
select a.tweet_text, b.song_name, b.song_artist, a.retweet
from music.tweets a join music.songs b on a.song_id = b.song_id
where a.tweet_text like '%Sweet%';

select * from song_retweets;
```

	tweet_text	song_name	song_artist	retweet
▶	#NowPlaying Eurythmics - Sweet Dreams (Are ...	Sweet Dreams (Are Made of This) - Remastered	Eurythmics, Annie Lennox, Dave Stewart	0
	Eurythmics, Annie Lennox, Dave Stewart - Swe...	Sweet Dreams (Are Made of This) - Remastered	Eurythmics, Annie Lennox, Dave Stewart	0
	Eurythmics-Sweet Dreams (Are Made of This) (...	Sweet Dreams (Are Made of This) - Remastered	Eurythmics, Annie Lennox, Dave Stewart	1
	Eurythmics, Annie Lennox, Dave Stewart - Swe...	Sweet Dreams (Are Made of This) - Remastered	Eurythmics, Annie Lennox, Dave Stewart	0
	# Sweet Dreams (Are Made of This) (Remaster...	Sweet Dreams (Are Made of This) - Remastered	Eurythmics, Annie Lennox, Dave Stewart	0

- 19) Fetch songs and retweets tweeted by spotify twitter handle and find similar results in your music database.

```
CREATE VIEW spotify_twitter_songs AS
select a.tweet_text,b.song_name,b.song_artist,a.retweet
from music.tweets a join music.songs b on a.song_id = b.song_id
where a.twitter_handle like '%AM25spotify%';

select * from spotify_twitter_songs;
```

	tweet_text	song_name	song_artist	retweet
▶	Everybody Wants To Rule The World - Tears Fo...	Everybody Wants To Rule The World	Tears For Fears	0

- 20) List the songs with genre where played count is more than 5000

```
CREATE VIEW 5000played AS
select b.song_name, b.song_genre, c.played_count
from music.songs b Join music.play_history c on b.song_id = c.song_id
where c.played_count > 5000;
```

```
select * from 5000played;
```

	song_name	song_genre	played_count
►	Eye of the Tiger	Rock	7647
	Everybody Wants To Rule The World	R&B	9455
	In Da Club	Hip-Hop	7498
	Black and Yellow	Rap	7111
	Let Me Love You	Pop	8560
	New Rules	Pop	5967
	Around the World (La La La La La) - Radio Version	Pop	5263
	Not Afraid	Hip-Hop	7407

Source of Data:

Finding the right source of data will help to get appropriate and good quality of data for database. For this project, our source of data is spotify and we have used exportify to export the data of spotify.

CODE:

The below code is the code we have written for scraping the twitter and getting the data and inserting into the tables we had created for twitter data in our database. As mentioned in the comments, we are making database connections, authenticating the user, getting all the tweets for a particular time period, and then insert the tweets or data we retrieved into respective tables in the database. And some queries regarding the data we got from twitter scraping.

```
1  import tweepy # for twitter api
2  import time
3  import mysql.connector
4  from mysql.connector import Error
5  import re
6  import requests
7  from requests_oauthlib import OAuth1
8  import json
9  import datetime
10 from datetime import datetime, timedelta
11
12 #database connection
13 def create_connection(host_name, user_name, user_password, db):
14     connection = None
15     try:
16         connection = mysql.connector.connect(
17             host=host_name,
18             user=user_name,
19             passwd=user_password,
20             db=db
21         )
22         print("Connection to MySQL DB successful")
23     except Error as e:
24         print("The error '{e}' occurred")
25     return connection
26 connection = create_connection("localhost", "root", "admin", "music")
27 cursor = connection.cursor()
28
29 # Authenticate
30 APP_KEY = "xZ40SNBnmkE9zN67GUBhPCq5B"
31 APP_SECRET = "tb57HaWSEzgGUU5pDC43JkRVxZ9f3YbEU0F7AEWFjuIUIQ7GSU"
32 USER_OAUTH_TOKEN = "1587541442548838400-yzs5qXr7g5vmAhlQ5N3GPQapKhd1TH"
33 USER_OAUTH_TOKEN_SECRET = "g5XvNK7Lg8wmyg6UCXE2QJ3FgpGrdXjkZcn51KbR9aUe9"
34 BEARER_TOKEN = "AAAAAAAAAAAAAAAAABFziwEAAAAA3mCapctlxrQEXK9%2FiGk6qNpTnXI%3D8ctkG70yB1Tm7Mj4LSf07I6q5I4qPh1fmt0en0dUF06Z5Kwha1"
35 authApi = OAuth1(APP_KEY, APP_SECRET,
36                 USER_OAUTH_TOKEN, USER_OAUTH_TOKEN_SECRET)
37
38 auth = tweepy.OAuthHandler(APP_KEY, APP_SECRET)
39 api = tweepy.API(auth)
40
41 # Retrieve Tweets
42 cursor.execute('SELECT * FROM songs s where s.song_popularity > 84')
43 songs = cursor.fetchall()
```

```

44
45 # Map song name and song id
46 song_connect= {}
47 for song in songs:
48     song_connect[song[1]]=song[0]
49
50 keywords=[]
51 for i in songs:
52     keywords.append(i[1])
53     print(i[1])
54 print(keywords)
55
56 for names in keywords:
57
58     public_tweets = api.search_tweets(names)
59     song_id = song_connect[names]
60     print("-----", public_tweets)
61
62     for tweet in public_tweets:
63         # Using tweets only from 2022-11-01 to 2022-11-12
64         t1=datetime.strptime("2022-11-01", "%Y-%m-%d")
65         t2=datetime.strptime("2022-11-12", "%Y-%m-%d")
66         t1date=t1.date()
67         t2date=t2.date()
68         checkdate= tweet.created_at.date()
69
70         if(checkdate<t1date and checkdate>t2date):
71             break
72         else:
73             tweet_id = tweet.id
74             created_at = tweet.created_at
75             tweet_text = tweet.text
76             username = tweet.user.screen_name
77             name = tweet.user.name
78             userId = tweet.user.id
79             follower_count = tweet.user.followers_count
80             following_count = tweet.user.friends_count
81             twitter_handle = tweet.user.screen_name
82             profile_image_url = tweet.user.profile_image_url_https
83             description = tweet.user.description
84             userCreated_at = tweet.user.created_at
85

```

```

86
87     cursor.execute('select * from tweets WHERE tweet_id = %s' , (tweet_id, ))
88     findTweet = cursor.fetchone()
89     print("findTweet-----", findTweet)
90     if findTweet:
91         print("Tweet already exists")
92     else:
93         cursor.execute('select * from twitter_user where twitter_handle = %s' , (twitter_handle, ))
94         findUser = cursor.fetchone()
95         print("findUser-----", findUser)
96         if findUser:
97             print("User already exists")
98         else:
99             #insert into twitter_user table
100             cursor.execute('insert into twitter_user (twitter_handle, name, description, followers_count, following_count, created_at) values (%s, %s, %s, %s, %s, %s)' , (twitter_handle, name, description, follower_count, following_count, userCreated_at))
101
102
103         tweets = api.get_status(tweet_id)
104         favorite_count = tweets.favorite_count
105         retweet_count = tweets.retweet_count
106
107         #insert into tweets table
108         cursor.execute('insert into tweets (tweet_id, twitter_handle, tweet_text, created_at, song_id, likes, retweet) values (%s, %s, %s, %s, %s, %s, %s)' , (tweet_id, twitter_handle, tweet_text, created_at, song_id, favorite_count, retweet_count))
109
110
111         if(len(tweet.entities['user_mentions']) > 0):
112             for mention in tweet.entities['user_mentions']:
113                 target_user = mention['screen_name']
114                 #insert into tweet_mentions table
115                 cursor.execute('insert into tweet_mentions (tweet_id, source_user, target_user) values (%s, %s, %s)' , (tweet_id, twitter_handle, target_user))
116
117         if(len(tweet.entities['hashtags']) > 0):
118             for tag in tweet.entities['hashtags']:
119                 tag = tag['text']
120                 #insert into tweet_tags table
121                 cursor.execute('insert into tweet_tags (tweet_id, tags) values (%s, %s)' , (tweet_id, tag))
122
123         connection.commit()
124         print(cursor.rowcount, "was inserted.")
125

```

```

127 # Q1 What user posted this tweet?
128 anyTweet = input("Enter tweet id: ")
129 first = cursor.execute("SELECT u.name, t.tweet_text FROM twitter_user as u INNER JOIN tweets as t ON u.twitter_handle = t.twitter_handle WHERE t.tweet_id = %s" % anyTweet)
130 first = cursor.fetchone()
131 print("Q1: ", first)
132
133 # Q2 When did the user post this tweet?
134 second = cursor.execute("SELECT u.name, t.tweet_text, t.created_at FROM twitter_user as u INNER JOIN tweets as t ON u.twitter_handle = t.twitter_handle WHERE t.tweet_id = %s" % anyTweet)
135 second = cursor.fetchone()
136 print("Q2: ", second)
137
138 # Q3 What tweets have this user posted in the past 24 hours?
139 anyUser = input("Enter user handle: ")
140 now = datetime.now()
141 prev = now + timedelta(days=-1)
142 now = datetime.strftime(now, "%Y-%m-%d %H:%M:%S")
143 prev = datetime.strftime(prev, "%Y-%m-%d %H:%M:%S")
144 print("time-----", type(prev), now)
145 fetchThree = api.user_timeline(screen_name=anyUser, count=100)
146 for tweet in fetchThree:
147     created_at_date = datetime.strftime(tweet.created_at, "%Y-%m-%d %H:%M:%S")
148     print(created_at_date)
149     if(created_at_date > prev and created_at_date < now):
150         print("tweet in 24 hours")
151         cursor.execute("insert into tweets (tweet_id, twitter handle, tweet text, created at, likes, retweet) values (%s, %s, %s, %s, %s, %s);",
152             (tweet.id, anyUser, tweet.text, tweet.created_at, tweet.favorite_count, tweet.retweet_count))
153         connection.commit()
154         print(cursor.rowcount, "was inserted.")
155     else:
156         print("tweet not in 24 hours")
157 qthree = cursor.execute("SELECT u.name, t.tweet_text, t.created_at FROM twitter_user as u INNER JOIN tweets as t ON u.twitter_handle = t.twitter_handle WHERE t.created_at > %s and t.created_at < %s" % (prev, now))
158 qthree = cursor.fetchall()
159 print("Q3: ", qthree)
160
161 # Q4 How many tweets have this user posted in the past 24 hours?
162 fourth = cursor.execute("SELECT u.name, count(t.tweet_text), t.created_at FROM twitter_user as u INNER JOIN tweets as t ON u.twitter_handle = t.twitter_handle WHERE t.created_at > %s and t.created_at < %s" % (prev, now))
163 fourth = cursor.fetchall()
164 print("Q4: ", fourth)
165
166 # Q5 When did this user join Twitter?
167 qfive = cursor.execute("SELECT u.twitter_handle, u.created_at FROM twitter_user as u where u.twitter_handle = %s" % (anyUser, ))
168 qfive = cursor.fetchall()
169 print("Q5: ", qfive)
170 # Q6 What keywords/ hashtags are popular?
171 qsix = cursor.execute("SELECT t.tags, count(t.tags) FROM tweet_tags as t GROUP BY t.tags ORDER BY count(t.tags) DESC LIMIT 10")
172 qsix = cursor.fetchall()
173 print("Q6: ", qsix)
174 # Q7 What tweets are popular?
175 qseven = cursor.execute("SELECT t.tweet_text, t.likes, t.retweet FROM tweets as t ORDER BY t.retweet DESC LIMIT 10")
176 qseven = cursor.fetchall()
177 print("Q7: ", qseven)

```

The scraped data was raw data that we audited, cleaned, and validated with completeness. This involved downloading and reformatting the raw data.

DATABASE PROJECT QUALITY

Audit Validity/ Accuracy, Audit Completeness, Audit Consistency/Uniformity:

A data audit is a step-by-step process that examines every step of the data science process. Problems can be introduced at any step of this process, so a full audit requires close examination at each step. Generally, auditing refers to inspecting an item or a process.

Once data is collected it is important to validate/audit that data to ensure that it can be used for effective decision-making. The intent of Completeness testing is to evaluate whether you have received a full set of data - data anomalies can be explored further after the preliminary testing. This also involved sanity checks. Some of the things considered during this process are below.


```
#import the libraries required
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import missingno as msno
%matplotlib inline

#read the csv file
import pandas as pd
data = pd.read_csv("D:\Music-recommendation-database-main\Music-recommendation-database-main\Data\sourceData.csv")
data
```

	Song_Id	Track Name	Artist URI(s)	Artist Name(s)	Album URI	Album Name	Album Artist URI(s)	Track URI	Album Artist Name(s)	Album Release Date	
0	1	Beat It	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	Michael Jackson	spotify:album:1C2h7mLnPSeVY6iMRTF4s	Thriller 25 Super Deluxe Edition	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:1OQtq8RnDM8kG2gqUPJaj	Michael Jackson	11/30/1982	https://i.scdn.co/image/ab6761
1	2	Smooth Criminal - 2012 Remaster	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	Michael Jackson	spotify:album:24TAupSNVWSAHL0R7n7vm	Bad 25th Anniversary	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:2kCQHf9g3G5BNdvUEiEnNk	Michael Jackson	8/31/1987	https://i.scdn.co/image/ab6761
2	3	Gimme More	spotify:artist:26d5oYcwsYlMAKD3tpOx4	Britney Spears	spotify:album:1ePKYGH5ZQCb1b4QeEDj	Blackout	spotify:artist:26d5oYcwsYlMAKD3tpOx4	spotify:track:6ic8OILUNEATToEFU3xmaH	Britney Spears	10/25/2007	https://i.scdn.co/image/ab6761
3	4	Eye of the Tiger	spotify:artist:26bcq2my5G87uRr558Qg	Survivor	spotify:album:3138bpFjGx0X4j5CRLLA	Rocky IV	spotify:artist:0YlQWjT6nXaFlPZqre9Of	spotify:track:2KH16WveTQWT6KQG9Rge2	Various Artists	6/7/1905	https://i.scdn.co/image/ab6761
4	5	Everybody Wants To Rule The World	spotify:artist:4bthk9UfsYUldcfyqumSUU	Tears For Fears	spotify:album:3myPwaMYjdwhbz9nf-gcG6W	Songs From The Big Chair (Super Deluxe Edition)	spotify:artist:4bthk9UfsYUldcfyqumSUU	spotify:track:4RlWVPyQ5RL0ao9LPZeSouE	Tears For Fears	2/25/1985	https://i.scdn.co/image/ab6761
--	--	--	--	--	--	--	--	--	--	--	--
315	316	Beat It / State of Shock - Immortal Version	spotify:artist:3fMbDgg4jU18AjlCKBhRSm, spotify_	Michael Jackson, The Jacksons, Mick Jagger, Ke-	spotify:album:5ReKddpdZXgpVlecQhLnEH	Immortal	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:3ejivP4oU6R0yIqSz5lwU	Michael Jackson	11/21/2011	https://i.scdn.co/image/ab6761
316	317	Thriller - Immortal Version	spotify:artist:3fMbDgg4jU18AjlCKBhRSm, spotify_	Michael Jackson, Kevin Antunes	spotify:album:5ReKddpdZXgpVlecQhLnEH	Immortal	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:4V79u5lp8Se1heEA2vWHO	Michael Jackson	11/21/2011	https://i.scdn.co/image/ab6761
317	318	Smooth Criminal - Immortal Version	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	Michael Jackson	spotify:album:5ReKddpdZXgpVlecQhLnEH	Immortal	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:4fF904O3cxc2wNMmMXAwlw	Michael Jackson	11/21/2011	https://i.scdn.co/image/ab6761
318	319	Thriller	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	Michael Jackson	spotify:album:7pMVCMwGyue9rzTHxLcm	Michael Jackson's This Is It	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:1D9KE0XlrmPUkMtdfzqgX4	Michael Jackson	10/26/2009	https://i.scdn.co/image/ab6761
319	320	Slave to the Rhythm - Audien Remix Radio Edit	spotify:artist:3fMbDgg4jU18AjlCKBhRSm, spotify_	Michael Jackson, Audien	spotify:album:64pX2Zupl8msRz1cnx0x5s	Slave to the Rhythm (Audien Remix Radio Edit)	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:2SL1bha5spc1q7U9gB40UC	Michael Jackson	8/12/2014	https://i.scdn.co/image/ab6761
320 rows x 12 columns											

#sort the data according to no of times played
data = data.sort_values('No of times Songs Played')
data

Python

...

	Song_id	Track Name	Artist URI(s)	Artist Name(s)	Album URI	Album Name	Album Artist URI(s)	Track URI	Album Artist Name(s)	Album Release Date
261	262	Gravitas	spotify:artist:7gz2f6L2XRnJsL0Ytve	Advakit	spotify:album:7y1Yf5KxQ5TLEZvuYVY	Gravitas	spotify:artist:7gz2f6L2XRnJsL0Ytve	spotify:track:2H4vzm58twRmDF5lessU	Advakit	5/1/2020 https://i.scdn.co
273	274	Aesthetic Arrest - Live	spotify:artist:5XMQ9F7g4nCx2CtmPm7d	Evanoff	spotify:album:3alq6KE5bVJef52hw3wFN	Evanoff Live	spotify:artist:5XMQ9F7g4nCx2CtmPm7d	spotify:track:5aC7n8JxU5nX2j6p6Y	Evanoff	10/11/2018 https://i.scdn.co
282	283	On My Way	spotify:artist:38aAEKHAmZwECZgFuUdC8, spotify:...	Teddy Beats, Mon RoVla	spotify:album:55cq5ZxKGw0RInUvAaIPI	On My Way	spotify:artist:38aAEKHAmZwECZgFuUdC8, spotify:...	spotify:track:2kqAlghIqg0OVWetDmw90ct	Teddy Beats, Mon RoVla	4/10/2020 https://i.scdn.co
37	38	Man in the Mirror - 2012 Remaster	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	Michael Jackson	spotify:album:24TAupSNVWSAHL0R7n7vm	Bad 25th Anniversary	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:1k0NativDusOZiR29W0Uj	Michael Jackson	8/31/1987 https://i.scdn.co
292	293	Raveolution	spotify:artist:53UXMZwcwQyV4j7zZaVF58, spotify:...	Sandro Silva, Graham Bell	spotify:album:6RgnP8RqKAwA8M8kaK4KA	Raveolution	spotify:artist:53UXMZwcwQyV4j7zZaVF58, spotify:...	spotify:track:3zs77DXDUubQXUdLuTHNb	Sandro Silva, Graham Bell	3/16/2020 https://i.scdn.cc
--	--	--	--	--	--	--	--	--	--	--
97	98	(Your Love Keeps Lifting Me) Higher & Higher	spotify:artist:4VnomLkTm9Ahe1zTmZju	Jackie Wilson	spotify:album:4OqYlmCwz3VK7A9GTWHN	Higher And Higher	spotify:artist:4VnomLkTm9Ahe1zTmZju	spotify:track:5gqg1H5OPMfuVZQ1wQNo7	Jackie Wilson	1967 https://i.scdn.co
128	129	The Other Way	spotify:artist:4Yep00c4kxZTyaRgzvhhTb, spotify:...	Bastique, Encure	spotify:album:5cvelZ7YfYmwQnWc3Zeymz	The Other Way	spotify:artist:4Yep00c4kxZTyaRgzvhhTb, spotify:...	spotify:track:6AAfNPXVE6SCFMjvotf7db	Bastique, Encure	6/4/2021 https://i.scdn.co
318	319	Thriller	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	Michael Jackson	spotify:album:7pMVCmwGyue9rzTHxLcm	Michael Jackson's This Is It	spotify:artist:3fMbDgg4jU18AjlCKBhRSm	spotify:track:1D9KE0XlrmPUkMtdfzqgX4	Michael Jackson	10/26/2009 https://i.scdn.co
232	233	Don't Let Me Down	spotify:artist:69GGBaA162ItqCwzG5jJp, spotify:...	The Chainsmokers, Daya	spotify:album:25BypSK8eZ2pasalwchczf	Don't Let Me Down	spotify:artist:69GGBaA162ItqCwzG5jJp, spotify:...	spotify:track:1t1f6WsaMmKE84TzqbK	The Chainsmokers	2/5/2016 https://i.scdn.co
223	224	Inception	spotify:artist:586NPHdKGokvNjUeIeyoS	Aleton	spotify:album:2HY0u9L9ib0dRuZkXbsg	Inception	spotify:artist:586NPHdKGokvNjUeIeyoS	spotify:track:79Hg6pgXlHeHfOTdId3d	Aleton	7/3/2020 https://i.scdn.cc

320 rows × 18 columns

```
for column in data:
    print(data[column].values)
```

[5]

Python

```
corelation = data.corr()
```

[6]

Python

```
... C:\Users\Abhishek\AppData\Local\Temp\ipykernel_25168\1518246052.py:1: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric_only to silence this warning.
    corelation = data.corr()
```

```
print(corelation.columns)
```

[7]

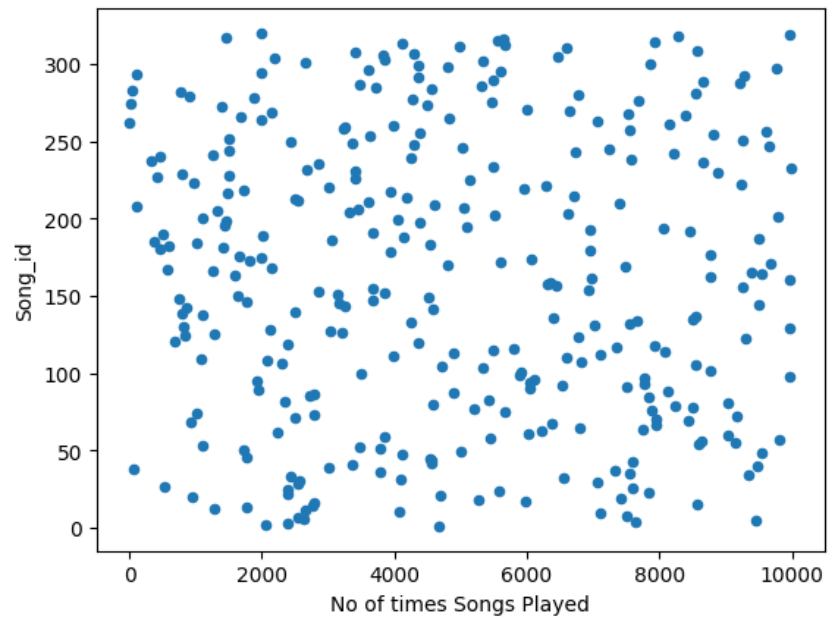
Python

```
... Index(['Song_id', 'Disc Number', 'Track Number', 'Track Duration (ms)',
        'Popularity', 'No of times Songs Played'],
        dtype='object')
```

```
data.plot(kind='scatter', y = 'Song_id', x='No of times Songs Played');
```

[8]

Python

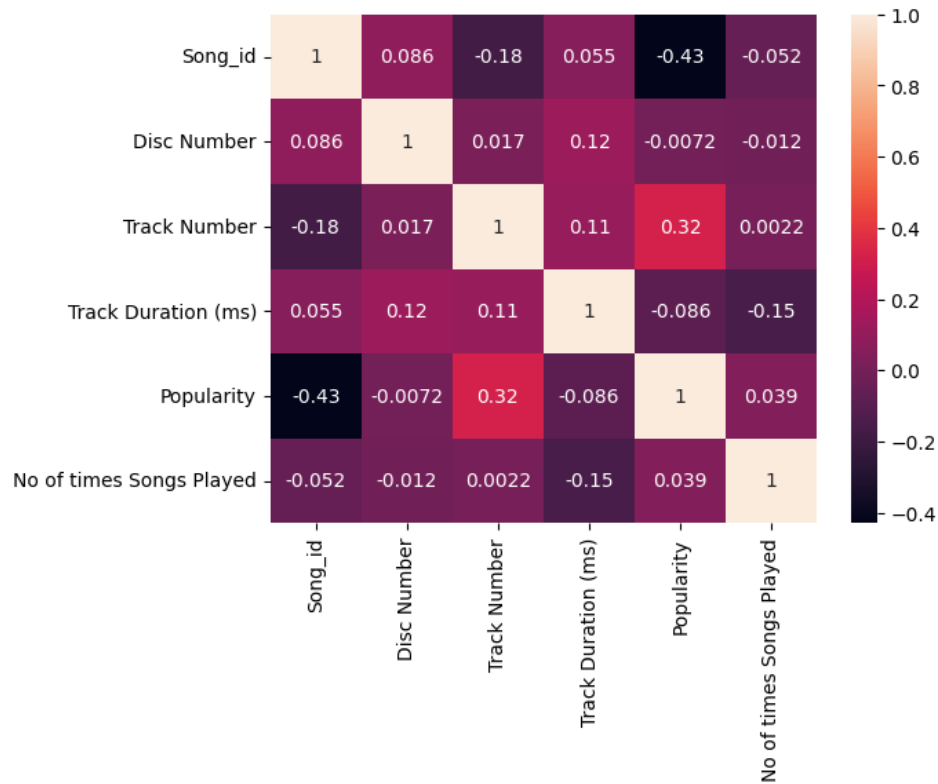


```
sns.heatmap(corelation, xticklabels=corelation.columns, yticklabels=corelation.columns, annot=True)
```

[9]

Python

```
... <AxesSubplot: >
```



```

# Audit completeness
data2 = pd.read_csv("D:\Music-recommendation-database-main\Music-recommendation-database-main\Data\sourceData.csv")
data2['Track Name'].describe()

```

[22] Python

```

count          316
unique          293
top    Around the World (La La La La La) - Radio Version
freq           4
Name: Track Name, dtype: object

```

data2.isnull().sum()

[11] Python

```

Song_id          0
Track Name       4
Artist URI(s)    5
Artist Name(s)   3
Album URI        7
Album Name       4
Album Artist URI(s) 6
Track URI        5
Album Artist Name(s) 4
Album Release Date 4
Album Image URL  6
Disc Number      5
Track Number     1
Track Duration (ms) 2
Track Preview URL 78
Explicit         3
Popularity       25
No of times Songs Played 1
dtype: int64

```

data3 = pd.read_csv("D:\Music-recommendation-database-main\Music-recommendation-database-main\Data\cleaned_data.csv")

data3

[23]

Python

...

	Song_Id	Song_Name	Artist_Name	Album_Name	Album_Artist_Name	Release_Year	Duration	Explicit	Popularity	Count_of_Plays
0	1	Beat It	Michael Jackson	Thriller 25 Super Deluxe Edition	Michael Jackson	1982	4:18	False	81	4680
1	2	Smooth Criminal - 2012 Remaster	Michael Jackson	Bad 25th Anniversary	Michael Jackson	1987	4:18	False	79	2055
2	3	Gimme More	Britney Spears	Blackout	Britney Spears	2007	4:11	False	78	2397
3	4	Eye of the Tiger	Survivor	Rocky IV	Various Artists	1905	4:06	False	76	7647
4	5	Everybody Wants To Rule The World	Tears For Fears	Songs From The Big Chair (Super Deluxe Edition)	Tears For Fears	1985	4:11	False	86	9455
...
290	316	Beat It / State of Shock - Immortal Version	Michael Jackson, The Jacksons, Mick Jagger, Ke...	Immortal	Michael Jackson	2011	3:09	False	37	5641
291	317	Thriller - Immortal Version	Michael Jackson, Kevin Antunes	Immortal	Michael Jackson	2011	3:38	False	42	1461
292	318	Smooth Criminal - Immortal Version	Michael Jackson	Immortal	Michael Jackson	2011	2:00	False	56	8275
293	319	Thriller	Michael Jackson	Michael Jackson's This Is It	Michael Jackson	2009	5:57	False	60	9973
294	320	Slave to the Rhythm - Audien Remix Radio Edit	Michael Jackson, Audien	Slave to the Rhythm (Audien Remix Radio Edit)	Michael Jackson	2014	3:14	False	32	1983

295 rows x 10 columns

#check for null data

data3.isnull().sum()

[13]

Python

...

Song_Id0

Song_Name0

Artist_Name0

Album_Name0

Album_Artist_Name0

Release_Year0

Duration0

Explicit0

Popularity0

Count_of_Plays0

dtype: int64

#plot a bar graph to represent null data from source data

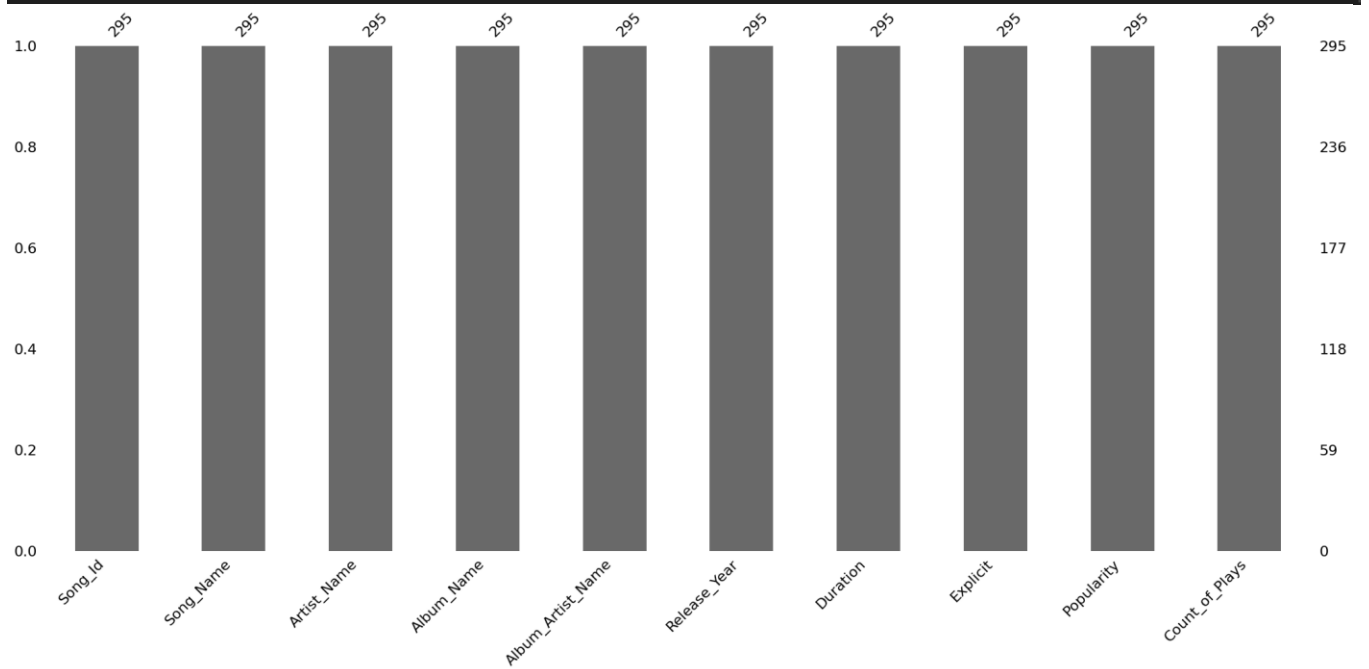
msno.bar(data2)

[14]

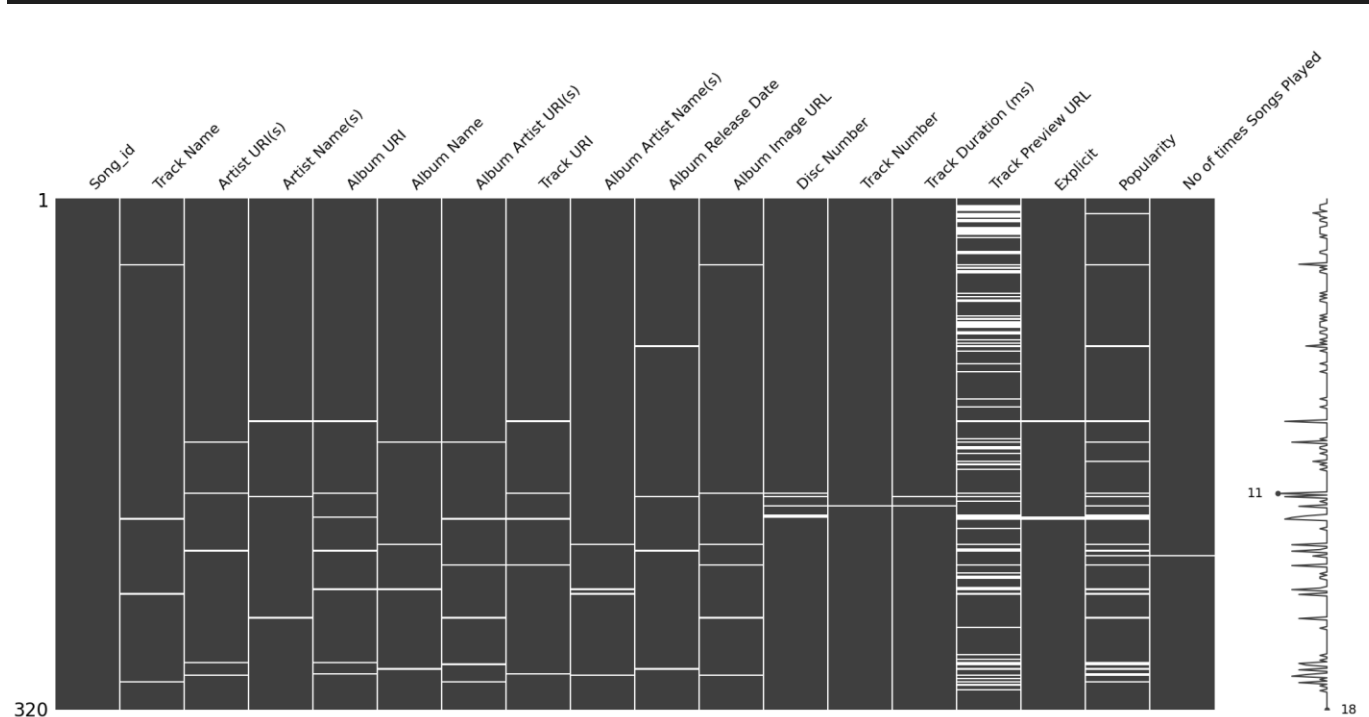
Python

Feature	Count of Null Values
Song_id	0
Track Name	0
Artist URI(s)	0
Artist Name(s)	0
Album URI	0
Album Name	0
Album Artist URI(s)	0
Track URI	0
Album Artist Name(s)	0
Album Release Date	0
Album Image URL	0
Disc Number	0
Track Number	0
Track Duration (ms)	0
Track Preview URL	242
Explicit	0
Popularity	295
No of times Songs Played	0

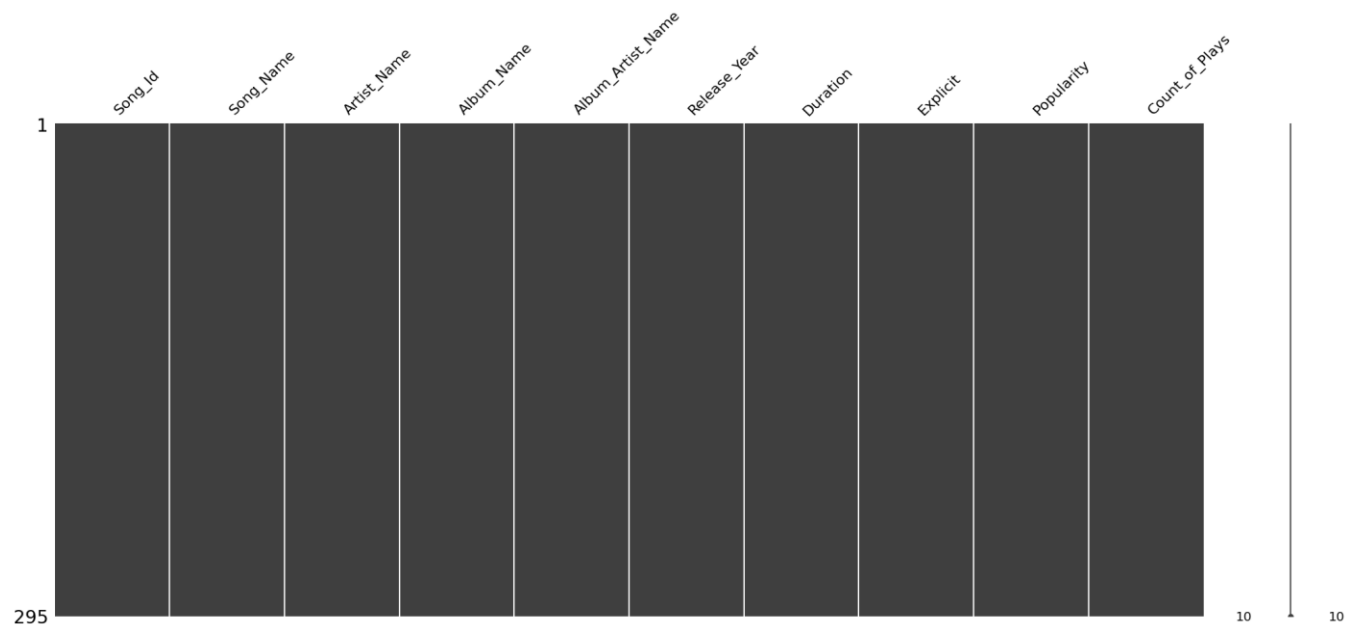
```
[15] #plot a bar graph to show null data after cleaning
msno.bar(data3) Python
```



```
[16] #plot a matrix graph to show null data from source data
msno.matrix(data2) Python
```



```
Python
#plot a matrix graph to show null data after cleaning of data
msno.matrix(data3)
```

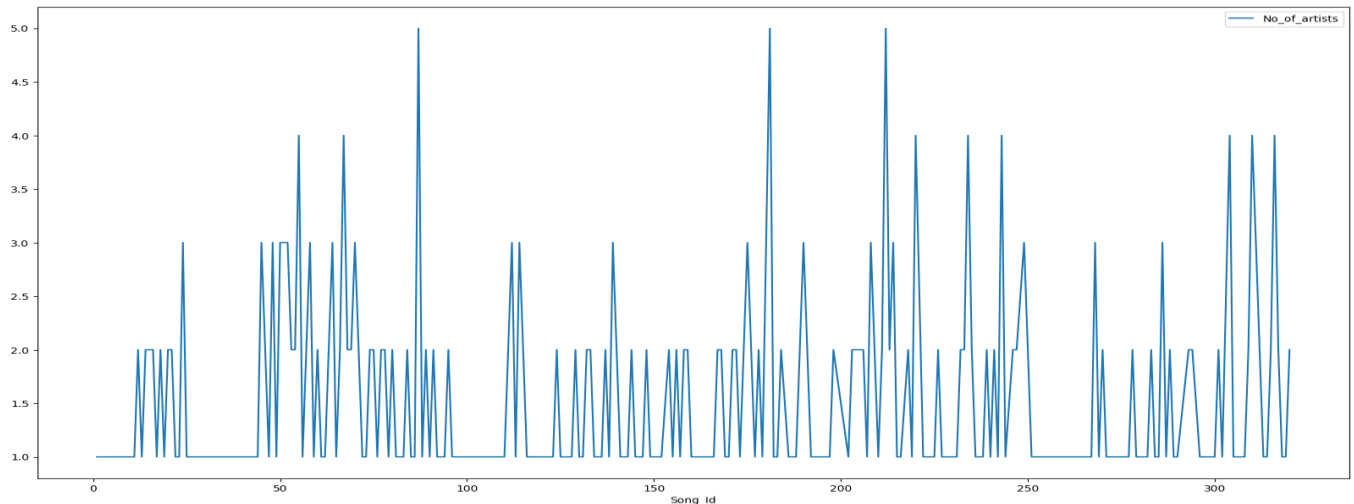


```
Python
#Audit consistency
data4 = pd.read_csv("D:\Music-recommendation-database-main\Music-recommendation-database-main\Data\Consistency.csv")
data4
```

	Song_Id	Song_Name	No_of_artists
0	1	Beat It	1
1	2	Smooth Criminal - 2012 Remaster	1
2	3	Gimme More	1
3	4	Eye of the Tiger	1
4	5	Everybody Wants To Rule The World	1
...
290	316	Beat It / State of Shock - Immortal Version	4
291	317	Thriller - Immortal Version	2
292	318	Smooth Criminal - Immortal Version	1
293	319	Thriller	1
294	320	Slave to the Rhythm - Audien Remix Radio Edit	2

295 rows x 3 columns

```
data4.plot(x="Song_Id", y="No_of_artists", kind="line", figsize=(20, 10))
plt.show()
#Each record of song has at least one artist
```



DATA CLENSING PROCESS BEFORE THE FORMATION OF FINAL DATABASE:

- After analyzing the data in the dataset exported from the Spotify platform, relevant columns and data were retained as per the table structure in the database.
- After the structural cleansing, the data was fixed at an individual level by adding constant values, to avoid complete data loss.
- Hence, Maximum Data was retained from the imported dataset.

```

1  import pandas as pd
2  import matplotlib.pyplot as plt
3  import seaborn as sns
4  import numpy as np
5  import mysql.connector
6  from mysql.connector import Error
7  import random
8  import os
9  import os.path
10
11  d = os.getcwd()
12  # os.chdir("..")
13  o = [os.path.join(d,o) for o in os.listdir(d) if os.path.isdir(os.path.join(d,o))] # Gets all directories in the folder as a tuple
14
15  #database connection
16  def create_connection(host_name, user_name, user_password, db):
17      connection = None
18      try:
19          connection = mysql.connector.connect(
20              host=host_name,
21              user=user_name,
22              passwd=user_password,
23              db=db
24          )
25          print("Connection to MySQL DB successful")
26      except Error as e:
27          print("The error '{e}' occurred")
28      return connection
29  connection = create_connection("localhost", "root", "admin", "music")
30  cursor = connection.cursor()

```

```

31
32 df = pd.read_csv("C:\\Users\\Abhishek\\Downloads\\new Source data.csv",encoding='latin-1')
33
34 #Check for null values and remove them
35 # df.describe()
36 print(df.isnull().sum())
37 print(df.describe)
38 #remove records with null values
39 modifiedData = df.dropna()
40 print(modifiedData.isnull().sum())
41 print(modifiedData.describe)
42 modifiedData.to_csv("C:\\Users\\Abhishek\\Downloads\\dataAfterCleaning.csv",index=False)
43

```

Pre-Clean Dataset Snapshot:

	A	B	C	D	E
1	Song_id	Track Name	Artist URI(s)	Artist Name(s)	Album URI
2	1	Beat It	spotify:artist:3fMbdgg4jU18AjLCKBhRSm	Michael Jackson	spotify:album:1C2h7mLntPSeVYciMRTF4a
3	2	Smooth Criminal - 2012 Remaster	spotify:artist:3fMbdgg4jU18AjLCKBhRSm	Michael Jackson	spotify:album:24TAupSNVWSAHLOR7n71v
4	3	Gimme More	spotify:artist:26dSoYclwsYLMaKD3tpOr4	Britney Spears	spotify:album:1ePkYch5ZQCb1b4tQeiEDj
5	4	Eye of the Tiger	spotify:artist:26bcq2nyj5GB7uRr558iQg	Survivor	spotify:album:3t3BbpFJiGcXl4jI5CRLLA
6	5	Everybody Wants To Rule The World	spotify:artist:4bthk9UfsYUYdcFyqxmSUU	Tears For Fears	spotify:album:3myPwaMYjdwhqtQnFgeG6
7	6	Everybody Talks	spotify:artist:0RpddSzUhfncUWNJXKOsJy	Neon Trees	spotify:album:0uRFz92JmjwDbZbB7hEBlr
8	7	Earned It (Fifty Shades Of Grey)	spotify:artist:1Xyo4u8uXC1ZmMpatF05PJ	The Weeknd	spotify:album:0P3oVJBFOv3TDXIRhGL7s
9	8	In Da Club	spotify:artist:3q7HBObVc0L8jNeTe5Gofh	50 Cent	spotify:album:5G5rgQHdzQnw32SI0WJlo5
10	9	Black and Yellow	spotify:artist:137W8MRPWKq5mrBGDBFSop	Wiz Khalifa	spotify:album:6ZOxIVL8rmk2ATHJiFJhiD
11	10	Superman	spotify:artist:7dGJo4pcD2V6oG8kP0tJRR, s	Eminem, Dina Rae	spotify:album:1ftvBBcu7jYlvXyt3JWB8S
12	11	In Da Club	spotify:artist:3q7HBObVc0L8jNeTe5Gofh	50 Cent	spotify:album:4ycNE7y1rp5215g1kkq1P
13	12	Candy Shop	spotify:artist:3q7HBObVc0L8jNeTe5Gofh, sp	50 Cent, Olivia	spotify:album:2pidzXTaHV4WaIJYRxKDCH
14	13	All We Have Is This Moment	spotify:artist:6rcldmstz3sdi2HTLzXJwg	C. SHIROCK	spotify:album:62oOarigBpkjNsQktNUdTY
15	14	rockstar (feat. 21 Savage)	spotify:artist:246dkjv51zLTtiykXe5h60, spoti	Post Malone, 21	spotify:album:6trNtQUgC8cgbWcqpMYkC
16	15	Let Me Love You	spotify:artist:540vlaP2JwJQb9dm3aArA4, sp	DJ Snake, Justin B	spotify:album:19zH1vHqzkxEsUvN7GaaH
17	16	Blame (feat. John Newman)	spotify:artist:7CajNmpbOovFoOoasH2HaY, s	Calvin Harris, Jo	spotify:album:48zisMeiXniWLzOQghbPqS
18	17	New Rules	spotify:artist:6M2wZ9GZgrQXHCffjv46we	Dua Lipa	spotify:album:01sfgrNbnnPUEYz6GZyt9
19	18	Around the World (La La La La La) - Re	spotify:artist:5wTdspmxb8V4ZjvDodpBo, sp	A Touch Of Class	spotify:album:2kBFEC9a71fNRXbRW5xO
20	19	Not Afraid	spotify:artist:7dGJo4pcD2V6oG8kP0tJRR	Eminem	spotify:album:47BiFcV59TQi2s9SkBo2pb
21	20	Love The Way You Lie	spotify:artist:7dGJo4pcD2V6oG8kP0tJRR, sp	Eminem, Rihanna	spotify:album:47BiFcV59TQi2s9SkBo2pb

Post-Clean Dataset Snapshot:

	A	B	C	D	E	F	G	H
1	song_id	song_name	song_album	song_artist	song_genre	song_length	song_language	song_popularity
2	1	Beat It	Thriller 25 Super Deluxe Edition	Michael Jackson	Pop	3	English	81
3	2	Smooth Criminal - 2012 Remaster	Bad 25th Anniversary	Michael Jackson	Pop	3	English	79
4	3	Gimme More	Blackout	Britney Spears	Pop	3	English	78
5	4	Eye of the Tiger	Rocky IV	Survivor	Rock	3	English	76
6	5	Everybody Wants To Rule The World	Songs From The Big Chair (Super D	Tears For Fears	R&B	3	English	86
7	6	Everybody Talks	Picture Show	Neon Trees	R&B	2	English	79
8	7	Earned It (Fifty Shades Of Grey)	Beauty Behind The Madness	The Weeknd	Pop	3	English	74
9	8	In Da Club	Get Rich Or Die Tryin'	50 Cent	Hip-Hop	2	English	82
10	9	Black and Yellow	Rolling Papers	Wiz Khalifa	Rap	3	English	75
11	10	Superman	The Eminem Show	Eminem, Dina Rae	Hip-Hop	4	English	0
12	11	In Da Club	Get Rich Or Die Tryin'	50 Cent	Hip-Hop	2	English	75
13	12	Candy Shop	The Massacre	50 Cent, Olivia	Hip-Hop	2	English	80
14	13	All We Have Is This Moment	All We Have Is This Moment	C. SHIROCK	R&B	3	English	9
15	14	rockstar (feat. 21 Savage)	beerbongs & bentleys	Post Malone, 21 Savage	R&B	3	English	84
16	15	Let Me Love You	Encore	DJ Snake, Justin Bieber	Pop	2	English	49
17	16	Blame (feat. John Newman)	Motion	Calvin Harris, John Newman	Electronic	2	English	76
18	17	New Rules	Dua Lipa (Deluxe)	Dua Lipa	Pop	2	English	80
19	18	Around the World (La La La La La) - Re	Planet Pop	A Touch Of Class, Pete Koneman	Pop	2	English	67
20	19	Not Afraid	Recovery	Eminem	Hip-Hop	3	English	80
21	20	Love The Way You Lie	Recovery	Eminem, Rihanna	Hip-Hop	3	English	82

NORMALIZATION

Normalization is the process of organizing data in a database. This includes creating tables and establishing relationships between those tables according to rules designed both to protect the data and to make the database more flexible by eliminating redundancy and inconsistent dependency.

The main goal of database normalization is to restructure the logical model of a database to:

- > Eliminate redundancy.
- > Organize data efficiently.
- > Reduce the potential for data anomalies (data anomalies are inconsistency in data stored as a result of database operation)

1st Normalization form:

- > Each table has a primary key.
- > The values in column of a table are atomic (no multi-value attributes).
- > There are no repeating groups.

How we converted our database into 1st NF:

Step 1: Recognizing and eliminating the repeating groups in tables.

Step 2: Identifying the primary key for each table.

Step 3: Identifying all the dependencies in each table.

Below is songs table from our music database.

song_id	song_name	song_album	song_artist	song_genre	song_length	song_language	song_popularity
1	Beat It	Thriller 25 Super Deluxe Edition	Michael Jackson	Pop	3	English	81
2	Smooth Criminal	Bad 25th Anniversary	Michael Jackson	Pop	3	English	79
3	Gimme More	Blackout	Britney Spears	Pop	3	English	78
4	Eye of the Tiger	Rocky IV	Survivor	Rock	3	English	76
5	Everybody Wants To Rule The World	Songs From The Big Chair (Super Deluxe Edition)	Tears For Fears	R&B	3	English	86

6	Everybody Talks	Picture Show	Neon Trees	R&B	2	English	79
7	Earned It (Fifty Shades Of Grey)	Beauty Behind The Madness	The Weeknd	Pop	3	English	74
8	In Da Club	Get Rich Or Die Tryin'	50 Cent	Hip-Hop	2	English	82
9	Black and Yellow	Rolling Papers	Wiz Khalifa	Rap	3	English	75
10	Superman	The Eminem Show	Eminem, Dina Rae	Hip-Hop	4	English	0

All the records in the table are atomic and there are no repeating groups.

2nd Normalization Form:

- > All requirements from 1st NF must be met.
- > Redundant data across multiple rows of table must be moved to a separate table.
- > The resulting table must be related to each other by use of foreign key.

Table is in 2nd Normal form when:

- It is in 1 Normal Form, and
- It includes no partial dependencies, i.e., no attribute is dependent on only a portion of primary key. In other words, every non-primary-key attribute is fully functionally dependent on the primary key.

How we converted our database into 2 NF:

After 1 NF,

Step 1: Writing each key component on a separate line.

Step 2: Assigning corresponding dependent attributes.

Below is the artist table from our music database

id	name	age	country	recordLabel
1	Michael Jackson	50	USA	Sony
2	Britney Spears	30	USA	Warner
4	Tears For Fears	34	Ireland	CBS
5	Neon Trees	23	Germany	Universal
6	The Weeknd	35	USA	BMG
7	50 Cent	36	Mexico	ARMIND

8	Wiz Khalifa	32	USA	Sony
9	Eminem	33	USA	Warner

For example, recordLabel is dependent is fully functionally dependent on the primary key which is artist_id. And all the records in artist table are atomic and no repeating groups.

3rd Normalization Form:

- > All requirements of 2 NF must be met.
- > Eliminate the keys that do not depend on primary key. i.e., any key that is dependent not only on primary key but also on another field must be moved to another table.

Table is in 3rd Normal Form when following are true:

- It is in 2nd Normal Form.
- It contains no transitive dependencies.
- There are no fields or attributes that does not depend on the key, i.e., non-key attributes must be dependent on key(s) but and only on the key(s).

How we converted our database into 3rd Normal Form:

Step 1: Identified each new determinant (determinant is any attribute whose value determines other values within a row.

Step 2: Identifying attributes dependent on each determinant identified in step 1 and identify dependency.

Step 3: Removing dependent attributes from transitive dependency.

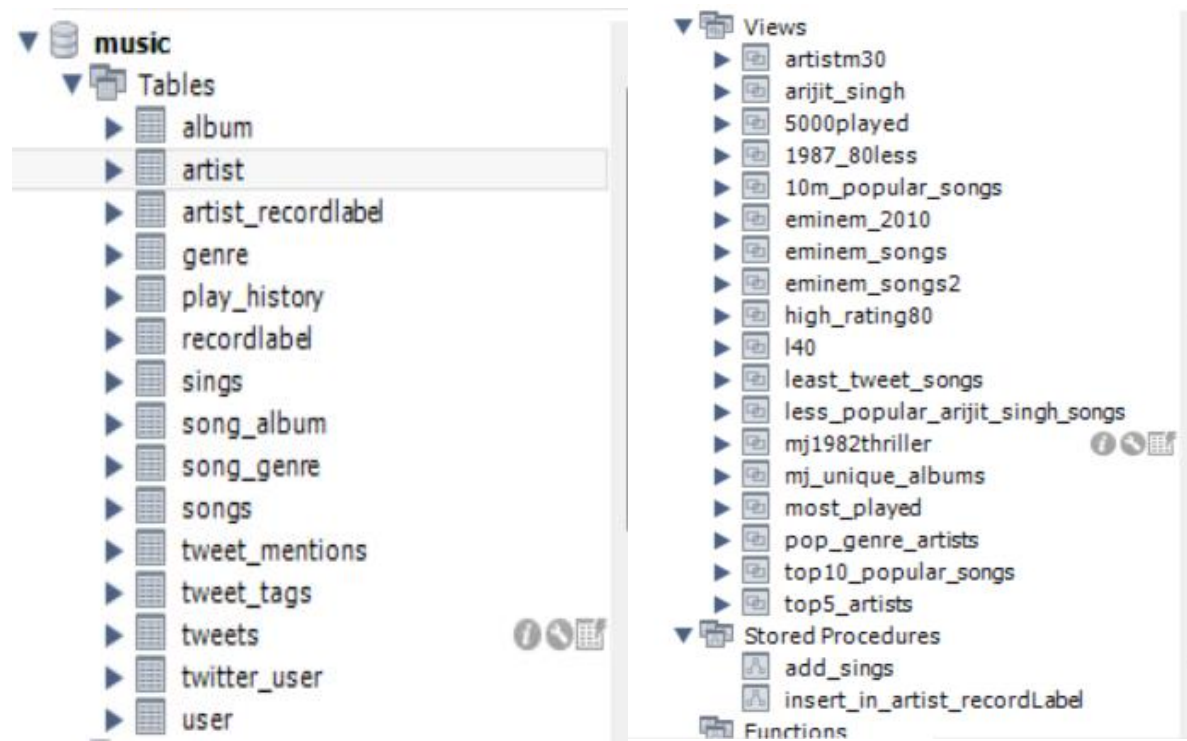
Below is the view from sings table.

song_name	artist_name
Beat It	Michael Jackson
Smooth Criminal - 2012 Remaster	Michael Jackson
Gimme More	Britney Spears
Everybody Wants To Rule The World	Tears For Fears
Everybody Talks	Neon Trees
Earned It (Fifty Shades Of Grey)	The Weeknd
In Da Club	50 Cent
Black and Yellow	Wiz Khalifa
In Da Club	50 Cent

'sings' table contains only songs_id and artist_id as attributes which are foreign keys from artist and songs table respectively, using those id's we have fetched the names of songs and artists accordingly.

This shows that there is no partial as well as transitive dependencies present in either songs or artist table.

DATABASE SCHEMA



Below are some examples of SQL Create and insert statements:

Songs Table:

```
CREATE TABLE `songs` (  
  `song_id` int NOT NULL AUTO_INCREMENT,  
  `song_name` text,  
  `song_album` text,  
  `song_artist` text,  
  `song_genre` text,  
  `song_length` bigint DEFAULT NULL,  
  `song_language` text,  
  `song_popularity` int DEFAULT NULL,  
  PRIMARY KEY (`song_id`)  
);
```

```
INSERT INTO `music`.`songs`
(`song_id`,`song_name`,`song_album`,`song_artist`,`song_genre`,`song_length`,`song_language`,`song_
popularity`)
VALUES(<{song_id: }>,<{song_name: }>,<{song_album: }>,<{song_artist: }>,<{song_genre:
}>,<{song_length: }>,<{song_language: }>,<{song_popularity: }>);
```

Eg:

```
Insert into Songs Values('In da club','Get rich or die trying','50 cent','Hip-hop',3:13,'English','85');
```

Artist Table:

```
CREATE TABLE `artist` (
  `id` int NOT NULL AUTO_INCREMENT,
  `name` varchar(25) DEFAULT NULL,
  `age` int DEFAULT NULL,
  `country` varchar(20) DEFAULT NULL,
  `recordLabel` varchar(20) DEFAULT NULL,
  PRIMARY KEY (`id`)
);
```

```
INSERT INTO `music`.`artist`
(`id`,`name`,`age`,`country`,`recordLabel`)
VALUES
(<{id: }>,<{name: }>,<{age: }>,<{country: }>,<{recordLabel: }>);
```

```
-- insert into sings values
insert into sings(
select a.song_id,b.id from songs a join artist b on a.song_artist = b.name);
```

```
-- insert into song_album values
insert into song_album(
select a.song_id,b.id from songs a join album b on a.song_album = b.name);
```

```
-- insert into song_genre values
insert into song_genre(
select a.song_id,b.id from songs a join genre b on a.song_genre = b.name);
```

```
-- insert into artist)recordlabel values
insert into artist_recordlabel(
select a.id,b.id from artist a join recordlabel b on a.recordLabel = b.name);
```