



A Case Study of Extract, Transform and Load: Spotify

Steven Green

Greyson Moore

Saatvi Rajgopal

Lucas da Silva

Background and Inspiration:

Spotify is a cloud-based music streaming service that allows users to listen to millions of songs directly from their phone, computer or internet-connected device. This type of music streaming has become the new normal for music fans and Spotify is the most popular music streaming platform in the world, servicing over 248 million users. Some of the advantages that this type of music platform offers are the ability to present the users with curated music charts based on genre, similar artists and popular songs near the location of the listener. It also has charts to showcase new and potentially popular songs, such as the “Viral 50” list, that is determined by compiling streaming data along from Spotify’s platform with analytics of trending music and artists from social media sites and blogs.

Proposal Topic:

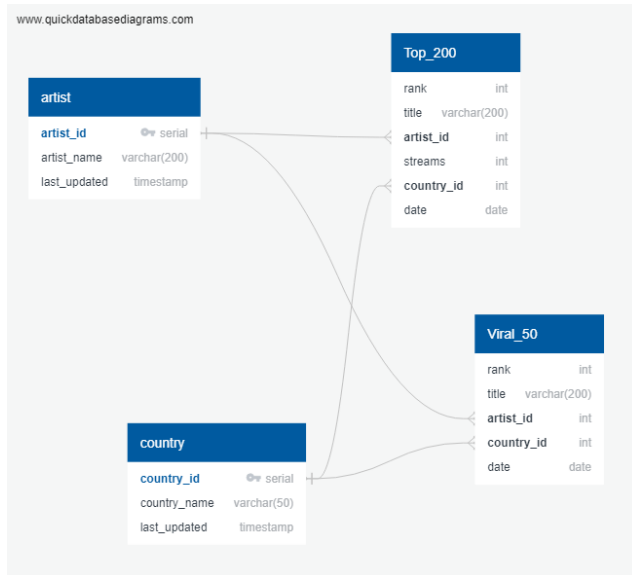
We will be creating a database on Google Cloud to store song information from Spotify’s most popular charts: the “Top 200” and the “Viral 50”. We will be analyzing the data from these charts for the most populous countries in the world where Spotify offers their services (United States, Thailand, India, Brazil, Mexico, Japan, Philippines, Egypt, Vietnam, and Turkey) and we will be comparing it with the overall global data that Spotify provides for these two charts. This data is compiled weekly and the timeframe we’ll be studying is from July 29, 2021, to September 1, 2021.

Process:

In order to create the database, the group will be web scraping song information from the Spotify charts website. The website contains lists of the “Top 200” streamed songs and the “Viral 50” songs in a week and it provides the ranking, song title, artist and streaming numbers for those songs. Using the Python library “BeautifulSoup”, the group will gather the information for each song over the selected time frame and will organize it into data frames using the Python library “pandas.” This data will be cleaned, transformed and uploaded to a PostgreSQL database. The group will then upload this database to the Google Cloud.

ERD:

We will be creating 4 tables: Artist, Country, Top 200 and Viral 50. The artist and country tables will have the names of the countries/artist and their primary keys (artist_id and country_id). These primary keys will feed into the Top 200 and the Viral 50 tables that will already have the ranking, title of the song, stream count, date and last updated information. Looking like this:



QUICKDBD

www.quickdatabasediagrams.com

Diagram Documentation

artist

Field	Description	Type	Default	Other
artist_id		serial		PK
artist_name		varchar(200)		
last_updated		timestamp		

country

Field	Description	Type	Default	Other
country_id		serial		PK
country_name		varchar(50)		
last_updated		timestamp		

Top_200

Field	Description	Type	Default	Other
rank		int		
title		varchar(200)		
artist_id		int		FK
streams		int		
country_id		int		FK
date		date		

Viral_50

Field	Description	Type	Default	Other
rank		int		
title		varchar(200)		
artist_id		int		FK
country_id		int		FK
date		date		

Web Scraping and Data Cleaning:

The group chose to use Python's Beautiful Soup to extract the data, and Pandas to help transform it from a disorganized array into an easy-to-read data-frame. The base URL chosen was [https://spotifycharts.com/regional/{countriesList\[i\]}/weekly/{start_date}--{end_date}](https://spotifycharts.com/regional/{countriesList[i]}/weekly/{start_date}--{end_date}), the country list and dates are present above in the Proposal Topic section of this document. After the data is extracted, Pandas organizes it in a table with the following column headers: Rank, Title, Artist, Streams, Country and Date:

	rank	title	artist	streams	country	date
0	1	STAY (with Justin Bieber)	The Kid LAROI	68764542	global	2021-09-03
1	2	INDUSTRY BABY (feat. Jack Harlow)	Lil Nas X	43663527	global	2021-09-03
2	3	Bad Habits	Ed Sheeran	34182580	global	2021-09-03
3	4	Beggin'	Måneskin	32620597	global	2021-09-03
4	5	Hurricane	Kanye West	30609671	global	2021-09-03
...

Then, the charts extracted from Spotify are merged and transformed with the country and artist tables, so they generate a table with chart data, primary keys, foreign keys and time stamp:

	rank	title	streams	date	country_id	last_updated	artist_id
0	1	STAY (with Justin Bieber)	68764542	2021-09-03	1	2021-09-14 01:38:45.297945	1
1	1	STAY (with Justin Bieber)	68764542	2021-09-03	1	2021-09-14 01:38:45.297945	1763
2	1	STAY (with Justin Bieber)	68764542	2021-09-03	1	2021-09-14 01:38:45.297945	3525
3	1	STAY (with Justin Bieber)	68764542	2021-09-03	12	2021-09-14 01:56:28.124757	1
4	1	STAY (with Justin Bieber)	68764542	2021-09-03	12	2021-09-14 01:56:28.124757	1763

Data Engineering and Queries

After the tables were loaded, the group ran a few queries to better understand some basic aspects of the data collected. The first query was to identify the artists that appeared the most in the Top 200 charts during the period analyzed.

	artist	count
0	Olivia Rodrigo	291
1	BTS	286
2	Doja Cat	205
3	The Weeknd	162
4	Billie Eilish	150
...

```

: #Most songs on the top 200 query
pd.read_sql_query("""
select
    a.artist,
    count(tt.artist_id)
from
    "Top_200" tt
join
    "artist" a on
    tt.artist_id = a.artist_id
group by
    a.artist
order by
    count(tt.artist_id) desc
limit 100;

""", con=engine)

```

The purpose of this query was to identify the most popular artists in Spotify's Top 200 charts during the month of August by checking the overall count of times they appeared in the charts. Olivia Rodrigo won this category.

The second query the group ran was to determine the artists that had been streamed the most globally during the month of August. So, the query counted the number of streams by artist globally:

```
#Top 10 Streams Globally
pd.read_sql_query("""
select distinct
    a.artist,
    sum(tt.streams) as "total streams globally"
from
    "artist" a
join
    "Top_200" tt on
    tt.artist_id = a.artist_id
join
    "country" c on
    tt.country_id = c.country_id
where
    c.country = 'global'
group by
    a.artist
order by
    sum(tt.streams) desc

limit 10;

""",con=engine)
```

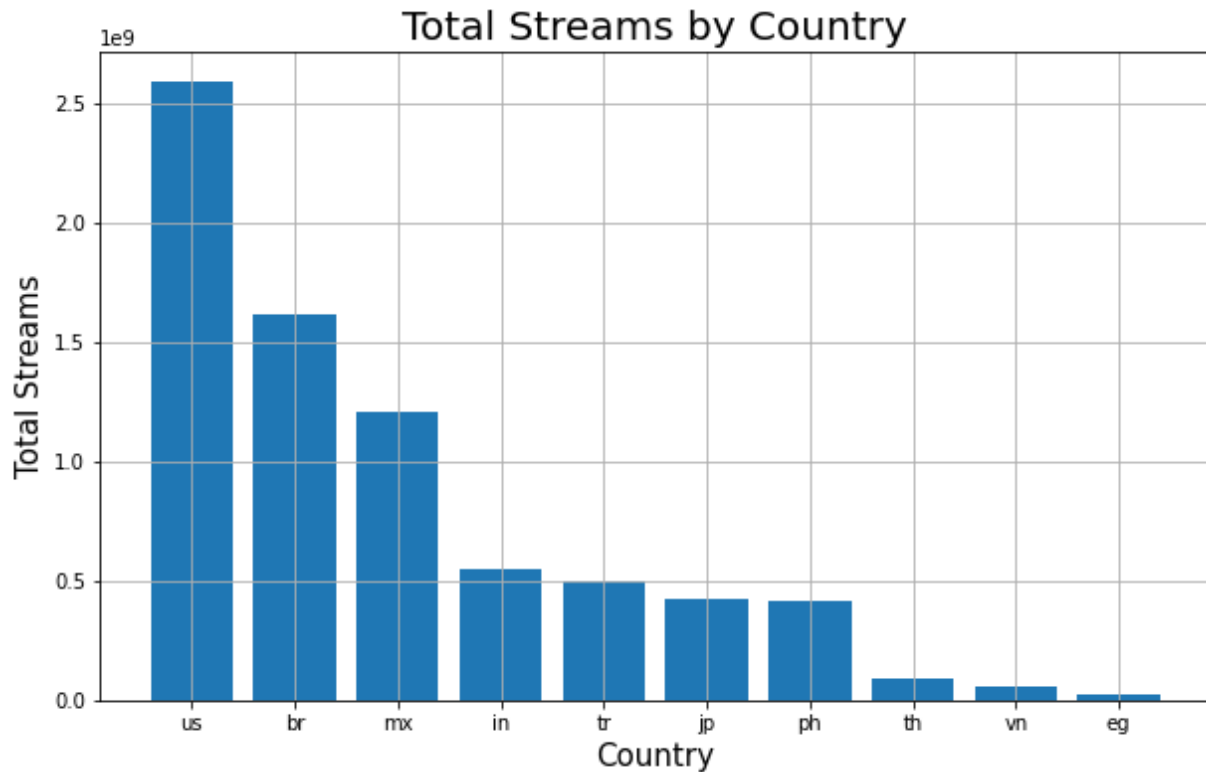
	artist	total streams globally
0	Olivia Rodrigo	613848271
1	Doja Cat	486483965
2	Billie Eilish	420394002
3	The Kid LAROI	387723854
4	Kanye West	365255505
5	Lil Nas X	335014898
6	Måneskin	318737165
7	The Weeknd	291351290
8	Ed Sheeran	267400555
9	BTS	215541918

Olivia Rodrigo was the most streamed artist globally during the period studied again, which begins to solidify Olivia as the most popular artist in Spotify during the last month.

The third query performed, and graphed, was meant to count the total number of streams by country during the last month, which could help us understand what countries have the most influence in the two types of charts the group analyzed, and it also identifies the countries where Spotify sees the most activity (or has the most active users):

```
plt.figure(figsize=(10,6))
plt.bar(countryStreamsDf.country,countryStreamsDf.Total_Streams)
plt.grid()
plt.title("Total Streams by Country",fontsize=20)
plt.xlabel("Country",fontsize=15)
plt.ylabel("Total Streams",fontsize=15)

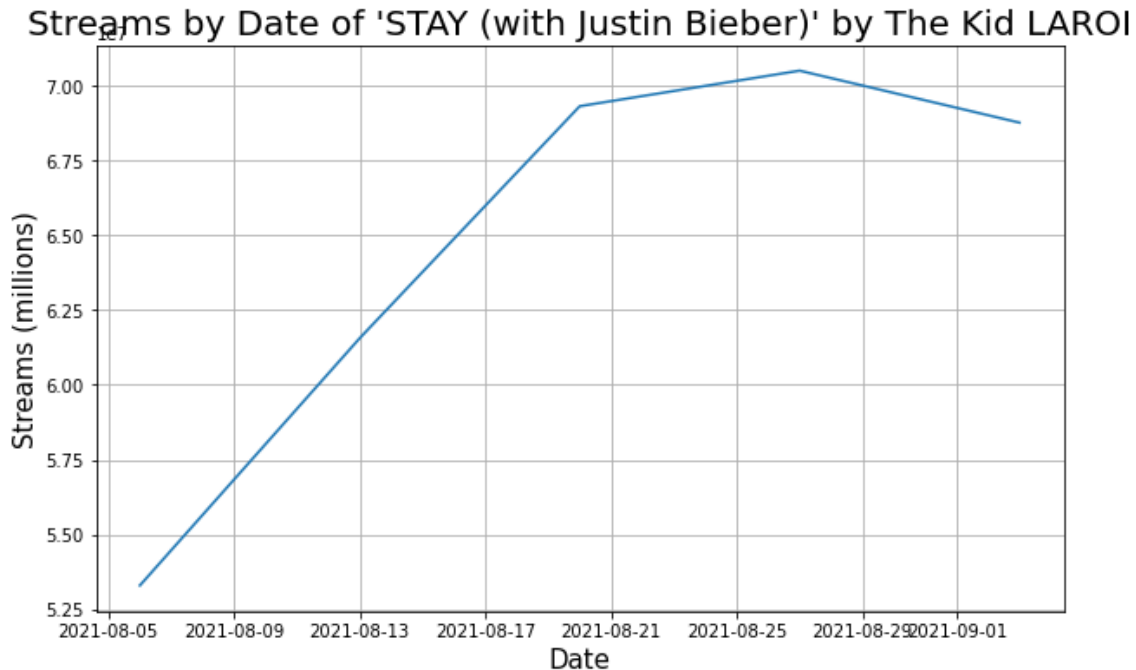
plt.show()
```



So, according to the graph, the United States is the country with the most number of streams among the songs that appeared in the Top 200 and Viral 50 charts, followed by Brazil. This would indicate that the taste for music of the Spotify users in these two countries are probably the most influential in the charts studied.

A fourth and final query was an attempt at understanding how popular a song was during the period studied, so we looked at the number of streams for a specific song, and how it varied, during the month. The group decided to analyze the most popular song in the Top 200 charts globally, which happened to be STAY, by Justin Bieber:

```
#Top 200 table query
pd.read_sql_query("""select tt.title,a.artist, c.country,tt.streams,tt.rank,tt.date
from
    "Top_200" tt
Join
    "country" c on
        tt.country_id = c.country_id
join artist a on
        tt.artist_id = a.artist_id;
""", con=engine)
```



The graph shows that this song started the month with a little more than 5 million streams, but in the following 25 days it shot up to over 7 million streams. So this song rapidly grew in popularity among Spotify users.

Conclusion:

The answers to these questions are heavily influenced by the limitations of the dataset, limitations of time and constraints of this project. If time permitted, the group would have liked to include more information from more datasets, analyzing longer periods of time, to the tables we created to have a better understanding of long-term trends and relationships between top artists and songs in the Top 200 chart and Viral 50 chart.

Some other interesting questions the group would have liked to explore, if it weren't for the limitations listed above, are the following: based on stream counts and country data, can we predict what artist/song is going to have the most streams for the month in the global charts? Can we project which artists are going to rise and fall during the year? Is there seasonality to music? Are some artists/genres more likely to be at the top during the summer? Or during the winter?

With a greater dataset and more time available, we could conceivably answer all of the questions above with a much stronger statistical support for any hypothesis that the group was interested in exploring.

Sources:

<https://spotifycharts.com/regional/global/weekly/latest>

<https://jennifer-franklin.medium.com/how-to-scrape-the-most-popular-songs-on-spotify-using-python-8a8979fa6b06>

<https://www.worldometers.info/world-population/population-by-country/>