

Music Genre Analysis



Connor Casey
Luke Milton
Matthew Matti
Lynn Soors
Dylan Mavencamp

Project 3

- Change in popularity of music genres across 2015 to 2024
- Spotify track streams
- # of total streams a genre has
- Based on the total amount of stream a song has got in its Spotify history
- Divided based on the songs release year

Introduction

- Number of times a song has been played on spotify
- Song must play for at least 30 seconds
- Will not count if <30 seconds
- Complex algorithm to count streams
- Tracks and removes “fraudulent streams”

What counts as a stream on Spotify?

-Process of finding the dataset

Spotify API

Kaggle Top 10000 Spotify Songs from 1950-Now

-Spotify API Problems

Could not figure out how to use Spotify's API

-Limitations with the dataset that we noticed right away

- No stream counts
- Vague popularity score
- Lot of genres
- Harder to clean

Our First Dataset

Used regex to extract the Track URI to be used in a loop to get stream counts

```
import re

df['Track_URI'] = df['Track URI'].str.extract(r'(?:(?:spotify:track:)(\s*\w{22,}))')

df
```

	Track URI	Track Name	Artist URI(s)	Artist Name(s)	Album URI	Album Artist Name(s)	Album Release Date	Track Duration (ms)	Track Preview URL	Explicit	Popularity	Artist Genres	Track_URI
0	spotify:track:0vNPJrUr8nMFdCs8b2MTNG	Fader	spotify:artist:4W48hZAnAHVOC2c8WH8pcq	The Temper Trap	spotify:album:0V59MMtgoruvEqMv18KAOH	The Temper Trap	2009	192373	https://p.scdn.co/mp3-preview/14264bd1501d2723...	False	0	indietronica,modern rock,shimmer pop	0vNPJrUr8nMFdCs8b2MTNG
1	spotify:track:0NpvdCO506uO58D4AbKzki	Sherry	spotify:artist:6mcrZQmgzFGRWf7C0SObou	Frankie Valli & The Four Seasons	spotify:album:0NUEQLaBzavncMEs4buZ	Frankie Valli & The Four Seasons	1/14/2003	152160	https://p.scdn.co/mp3-preview/e3f765262ebc349e...	False	54	adult standards,bubblegum pop,doo-wop,lounge,n...	0NpvdCO506uO58D4AbKzki
2	spotify:track:1MtUq6Wp1eQ8PC6BbPCj8P	I Took A Pill In Ibiza - Seeb Remix	spotify:artist:2KsP6tYlJITBvSUxmwVWw,spotify...	Mike Posner, Seeb	spotify:album:1Tz3Ai1guEF4hV3d9i17K	Mike Posner	5/6/2016	197933	https://p.scdn.co/mp3-preview/7bae6aac6d699135...	True	63	dance pop,pop,pop dance,pop rap,pop dance	1MtUq6Wp1eQ8PC6BbPCj8P
3	spotify:track:59lq75uFlqzUZcgZ4CbqFG	Let Go for Tonight	spotify:artist:7qRIl6DYV06u2VuRPAVqug	Foxes	spotify:album:5AQ7uKRSpAv7SNUI4j24ru	Foxes	5/12/2014	238413	https://p.scdn.co/mp3-preview/84a003d72f9f1468...	False	39	electropop,metropolis,uk pop	59lq75uFlqzUZcgZ4CbqFG
4	spotify:track:7KdcZQ3GJeGdserhK61kf	The Way I Want To Touch You	spotify:artist:78EFMxbaqx6dOpbtIEqScm	Captain & Tennille	spotify:album:3GUXesVyOehlmaxJyCTh6d	Captain & Tennille	1/1/1975	163586	https://p.scdn.co/mp3-preview/9e7a4a7b7dc56dc3...	False	35	mellow gold,soft rock,yacht rock	7KdcZQ3GJeGdserhK61kf

While looking for data to use. Soundcharts was one that we came across.

-Pros and cons of using Soundcharts

Pros: A very very large amount and wide ranging data, and very helpful support.

Cons: It is a paid service

After reaching out to Soundcharts explaining what our group was doing, we were given access for 5 days to the site.

Songs

Sort by

Total

Growth

% Growth

7D

28D

90D

Artist Country All

Song Genres All

Streams 1K - 4.7B

+ Add filters

60,700,696 songs - Updated today - 03:20 am

#	SONG	STREAMS	POPULARITY	VIEWS	VIDEOS	# PLAYS	COUNTS	PLAYS	STREAMS	PLAYS	FAVORITES
1	<div></div> <div>Die With A Smile</div> <div>Pop</div> <div>Bruno Mars Lady Gaga</div>	1.73B	100	720.55M	2.4M	318.35K	9.23M	611.04K	4.98K		
		+0	+0	+472.8M	+700K	+135.8K	+0	+102.95%	+4.5K		
		0%	0%	+138.31%	+141.8%	+174.56%	0%		+996.48%		
2	<div></div> <div>DMF</div> <div>Latin</div> <div>Bad Bunny</div>	284.6M	98			6.45K	1.23M	1.44M			
		+0	+0			+0	+0	+0			
		0%	0%			0%	0%	0%			
3	<div></div> <div>BIRDS OF A FEATHER</div> <div>Alternative</div> <div>Billie Eilish</div>	1.98B	96	351.53M	1.4M	551.06K	7.38M	8.9M	369.78K	657	
		+581.6M	+1	+112.2M	+200K	+245.4K	+2.3M	+137K		+402	
		+11.74%	+1.03%	+66.65%	+16.67%	+80.28%	+36.69%	+35.25%	+158.92%	+157.65%	
4	<div></div> <div>That's So True</div> <div>Pop</div> <div>Gracie Abrams</div>	606.7M	96	10.52M	541.6K	2.4M		571.87K			
		+551.6M	+9	+9.2M	+476.2K	+2.3M		+0			
		11K%	+10.34%	+701.58%	+778.13%			1.89K%	0%		
5	<div></div> <div>APT.</div> <div>Pop</div> <div>Others</div> <div>Bruno Mars ROSÉ</div>	1.06B	95	1.17B	4.54M	991	6.64M	12.81M	5.23K		
		+0	+1	+874.4M	+2.8M	+768	+5.9M	+11.8M	+5K		
		0%	+1.04%	+300.08%	+164.2%	+138.31%	+80.41%	+92.92%	+11.76%		
6	<div></div> <div>BAILE INOLVIDABLE</div> <div>Latin</div> <div>Bad Bunny</div>	189.83M	95		2.27K	423.47K	873.76K				
		+0	+0		+0	+0	+0	+0			
		0%	0%		0%	0%	0%	0%			
7	<div></div> <div>NUEVA VOY</div> <div>Latin</div> <div>Bad Bunny</div>	175.4M	94		4.39K	538.14K	789.12K				
		+0	+0		+0	+0	+0	+0			
		0%	0%		0%	0%	0%	0%			
8	<div></div> <div>Sailor Song</div> <div>Pop</div> <div>Rock</div> <div>Gigi Perez</div>	664.27M	93	7.27M	531.4K	58.49K	1.91M				
		+392.2M	+1	+4.8M	+162.7K	+53.6K	+1.5M				
		+144.24%	+1.09%	+206.27%	+44.13%	+909.09%	+189.01%				
9	<div></div> <div>Timeless</div> <div>Latin</div> <div>Playboi Carti The Weeknd</div>	513.98M	93	15.64M	80.49K	1.22K	2.84M	7.65M	80.81K	607	
		+0	+1	+9.2M	+376K	+436	+2.2M	+4.7M	+34K	+518	
		0%	+1.09%	+141.01%	+875.53%	+155.9%	+355.57%	+157.34%	+172.48%	+1582.02%	
10	<div></div> <div>VOY A LLAVARTE PA PR</div> <div>Latin</div> <div>Bad Bunny</div>	145.19M	93		995	283.58K	679.26K				
		+0	+0		+0	+0	+0	+0			
		0%	0%		0%	0%	0%	0%			

Switching to Soundcharts

```
['R&B' 'Pop, Folk' 'Alternative, Rock' ... 'Rock, Country, Spoken, Pop'  
'Electro, Soundtrack, Classical, Rock' 'Pop, Country, Alternative']
```

```
]: # First, split the 'Song genres' column by commas and remove leading/trailing spaces
```

```
df_cleaned['Song genres'] = df_cleaned['Song genres'].str.split(',')
```

```
# Now explode the 'Song genres' column to create a new row for each genre
```

```
df_cleaned = df_cleaned.explode('Song genres')
```

```
# Remove any extra whitespace from the 'Song genres' column
```

```
df_cleaned['Song genres'] = df_cleaned['Song genres'].str.strip()
```

```
# Check the unique genres after splitting
```

```
print(df_cleaned['Song genres'].unique())
```

```
['R&B' 'Pop' 'Folk' 'Alternative' 'Rock' 'Electro' 'Hip Hop' 'Soundtrack'  
'Metal' 'Country' 'Latin' 'Blues' 'Others' 'Reggae' 'African' 'Jazz'  
'Kids' 'Spirituals' 'Classical' 'Sports' 'Asian' 'Holiday' 'Instrumental'  
'Spoken' 'European']
```

```
]: # This was unique given the data in Release date. Just doing 1 .fillna would not work. We had to first fill NaN values with a string version of 1776 (Just
```

```
# Fill NaN values in 'Release date' with 1776 and convert to string
```

```
df_cleaned['Year Released'] = df_cleaned['Release date'].fillna('1776').astype(str)
```

```
# Extract the year from 'Year Released' using regex
```

```
df_cleaned['Year Released'] = df_cleaned['Year Released'].str.extract(r'(\d{4})')
```

```
# Fill any remaining NaN values in 'Year Released' with 1776
```

```
df_cleaned['Year Released'] = df_cleaned['Year Released'].fillna(1776)
```

```
# Convert the 'Year Released' column to integers so further mathematical processes can be done
```

```
df_cleaned['Year Released'] = df_cleaned['Year Released'].astype(int)
```

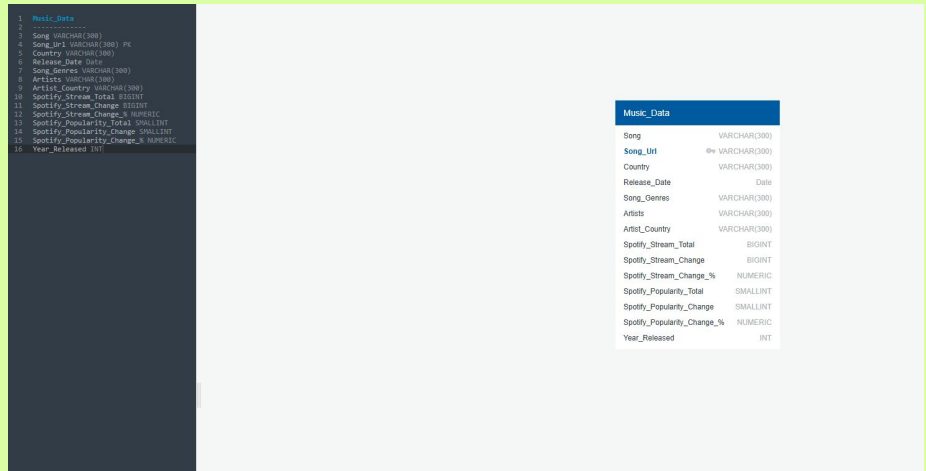

- Taking Matts code and manually cleaning the CSV
- Some dates had errors
- Wanted to include anything that may be of use



Basic CSV Cleaning

- Created a QuickDBD for the Excel file
- Using Postgres created a database
- Tested my database with SELECT queries

Setting up an SQL Database



```
1 Music_Data
2
3 CREATE TABLE
4 Song VARCHAR(300)
5 Song_URI VARCHAR(300) PK
6 Country VARCHAR(300)
7 Release_Date DATE
8 Song_Genres VARCHAR(300)
9 Artists VARCHAR(300)
10 Artist_Country VARCHAR(300)
11 Spotify_Stream_Total BIGINT
12 Spotify_Stream_Change BIGINT
13 Spotify_Popularity_Total SMALLINT
14 Spotify_Popularity_Change SMALLINT
15 Spotify_Popularity_Change_% NUMERIC
16 Year_Released INT
```

Music_Data	
Song	VARCHAR(300)
Song_URI	PK VARCHAR(300)
Country	VARCHAR(300)
Release_Date	DATE
Song_Genres	VARCHAR(300)
Artists	VARCHAR(300)
Artist_Country	VARCHAR(300)
Spotify_Stream_Total	BIGINT
Spotify_Stream_Change	BIGINT
Spotify_Stream_Change_%	NUMERIC
Spotify_Popularity_Total	SMALLINT
Spotify_Popularity_Change	SMALLINT
Spotify_Popularity_Change_%	NUMERIC
Year_Released	INT

Process of creating interactive visualizations

- Read clean csv into Pandas file
- Loop through clean csv in Pandas to create data frame with year and combined streams for genre within that year
- Convert to json file
- Write html and .js files to create interactive visualizations run via live server
- Made graphs using Plotly

Select Year

2024



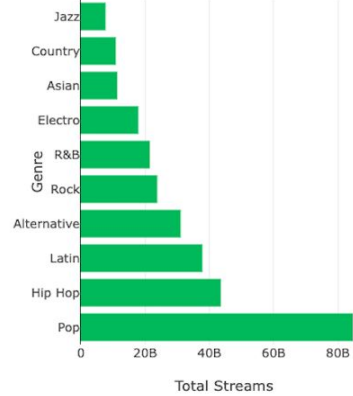
Statistics

Total Songs: 1,026

Total Streams: 238,699,508,210

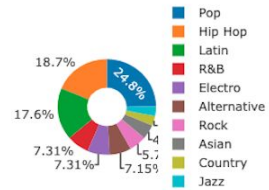
Unique Genres: 101

Top 10 Genres by Streams

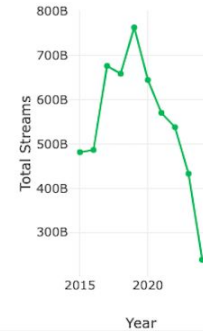


[interactive link](#)

Genre Distribution



Streaming Trends Over Time



Used SQLiteStudio to test queries on the songs.db

```
# Connect to SQLite3 database
conn = sqlite3.connect('songs.db')

start_year = 2015
end_year = 2024
# SQL query for "Rock" genre
query_rock = f"""
SELECT "Song", "Spotify stream Total", "Year Released"
  FROM songs_data
 WHERE "Year Released" BETWEEN {start_year} AND {end_year} AND "Song Genres" LIKE "%Rock%"
 GROUP BY "Year Released"
 ORDER BY "Spotify stream Total" DESC;
"""

# SQL query for "Pop" genre
query_pop = f"""
SELECT "Song", "Spotify stream Total", "Year Released"
  FROM songs_data
 WHERE "Year Released" BETWEEN {start_year} AND {end_year} AND "Song genres" LIKE "%Pop%"
 GROUP BY "Year Released"
 ORDER BY "Spotify stream Total" DESC;
"""

# SQL query for "Hip Hop" genre
query_hip_hop = f"""
SELECT "Song", "Spotify stream Total", "Year Released"
  FROM songs_data
 WHERE "Year Released" BETWEEN {start_year} AND {end_year} AND "Song Genres" LIKE "%Hip Hop%"
 GROUP BY "Year Released"
 ORDER BY "Spotify stream Total" DESC;
"""

# SQL query for "Alternative" genre
query_alt = f"""
SELECT "Song", "Spotify stream Total", "Year Released"
  FROM songs_data
 WHERE "Year Released" BETWEEN {start_year} AND {end_year} AND "Song Genres" LIKE "%Alternative%"
 GROUP BY "Year Released"
 ORDER BY "Spotify stream Total" DESC;
"""

# SQL query for "Country" genre
query_country = f"""
SELECT "Song", "Spotify stream Total", "Year Released"
  FROM songs_data
 WHERE "Year Released" BETWEEN {start_year} AND {end_year} AND "Song Genres" LIKE "%Country%"
 GROUP BY "Year Released"

```

```
# Read the data into pandas DataFrames
```

```
df_query1 = pd.read_sql_query(query_rock, conn) # Rock data
```

```
df_query2 = pd.read_sql_query(query_pop, conn) # Pop data
```

```
df_query3 = pd.read_sql_query(query_hip_hop, conn) # Hip Hop  
data
```

Reading the queries.

Using plotly to create the chart

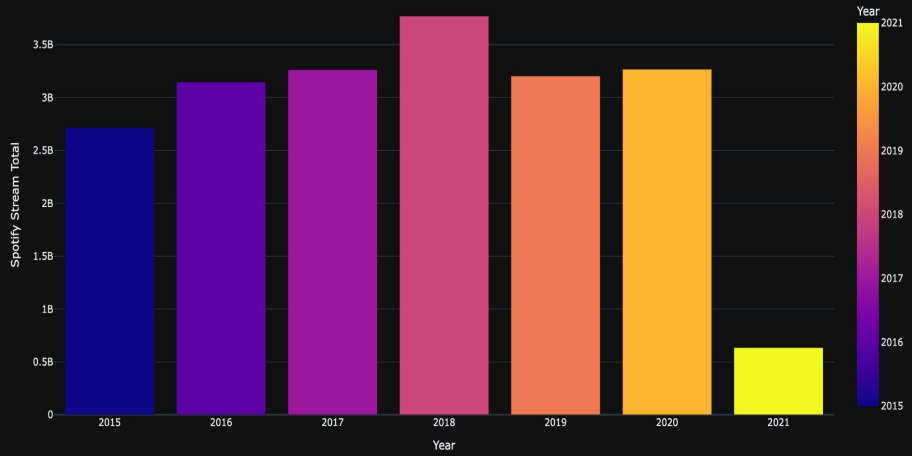
```
import plotly.express as px

# Create the bar chart for Rock genre
fig rock = px.bar(df query1,
                  x='Year Released',
                  y='Spotify stream Total',
                  title="Total Spotify Streams by Year for Rock Genre",
                  labels={'Year Released': 'Year', 'Spotify stream Total': 'Spotify
Stream Total'},
                  color='Year Released', # Optional: color by year
                  template='plotly dark') # Optional: dark theme for better visibility

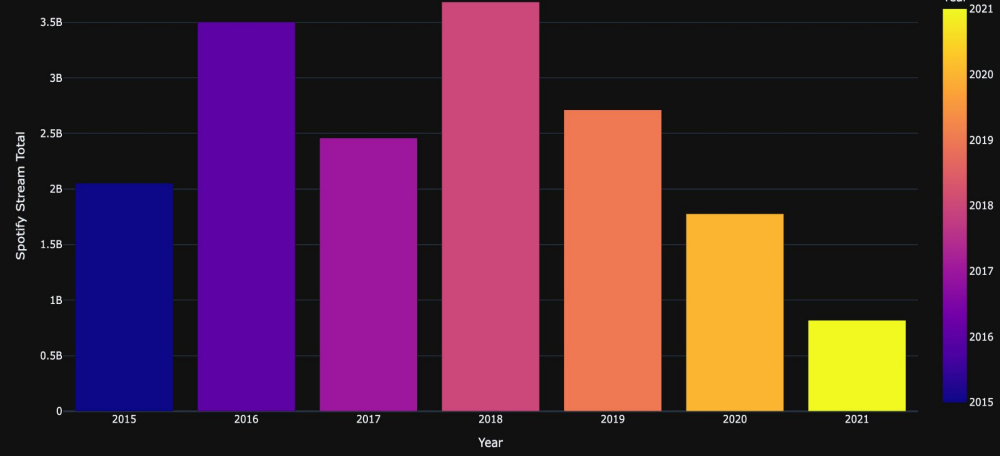
# Save the chart as an interactive HTML file
fig rock.write_html("html/spotify rock streams by year.html")

# Show the plot
fig rock.show()
```

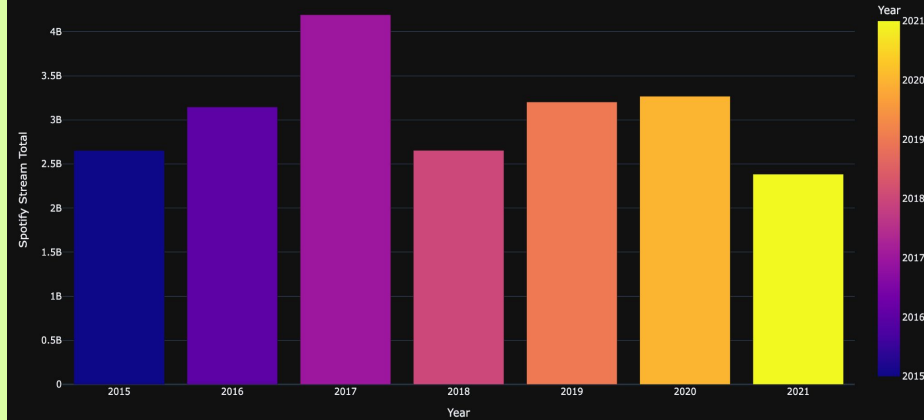
Total Spotify Streams by Year for Rock Genre



Total Spotify Streams by Year for Hip Hop Genre



Total Spotify Streams by Year for Pop Genre




```
# Concatenate DataFrames
```

```
df combined = pd.concat([df query1.assign(Genre="Rock"), df query2.assign(Genre="Pop"),  
df query3.assign(Genre="Hip Hop"), df query4.assign(Genre="Alternative"), df query5.assign(Genre="Country"),  
df query6.assign(Genre="Latin"), df query7.assign(Genre="Electro"), df query8.assign(Genre=("R&B")),  
df query9.assign(Genre="Folk"), df query10.assign(Genre=("Jazz")), df query11.assign(Genre=("Metal")),  
df query12.assign(Genre=("Soundtrack"))])
```

```
# Create a combined bar chart with side-by-side bars
```

```
fig combined = px.bar(df combined,  
x='Year Released',  
y='Spotify stream Total',  
color='Genre', # Color by genre  
title="Total Spotify Streams by Genres",  
labels={'Year Released': 'Year', 'Spotify stream Total': 'Spotify Stream Total'},  
template='plotly dark',  
barmode='group') # Set bars to be side by side
```

```
# Save the combined chart as an interactive HTML file
```

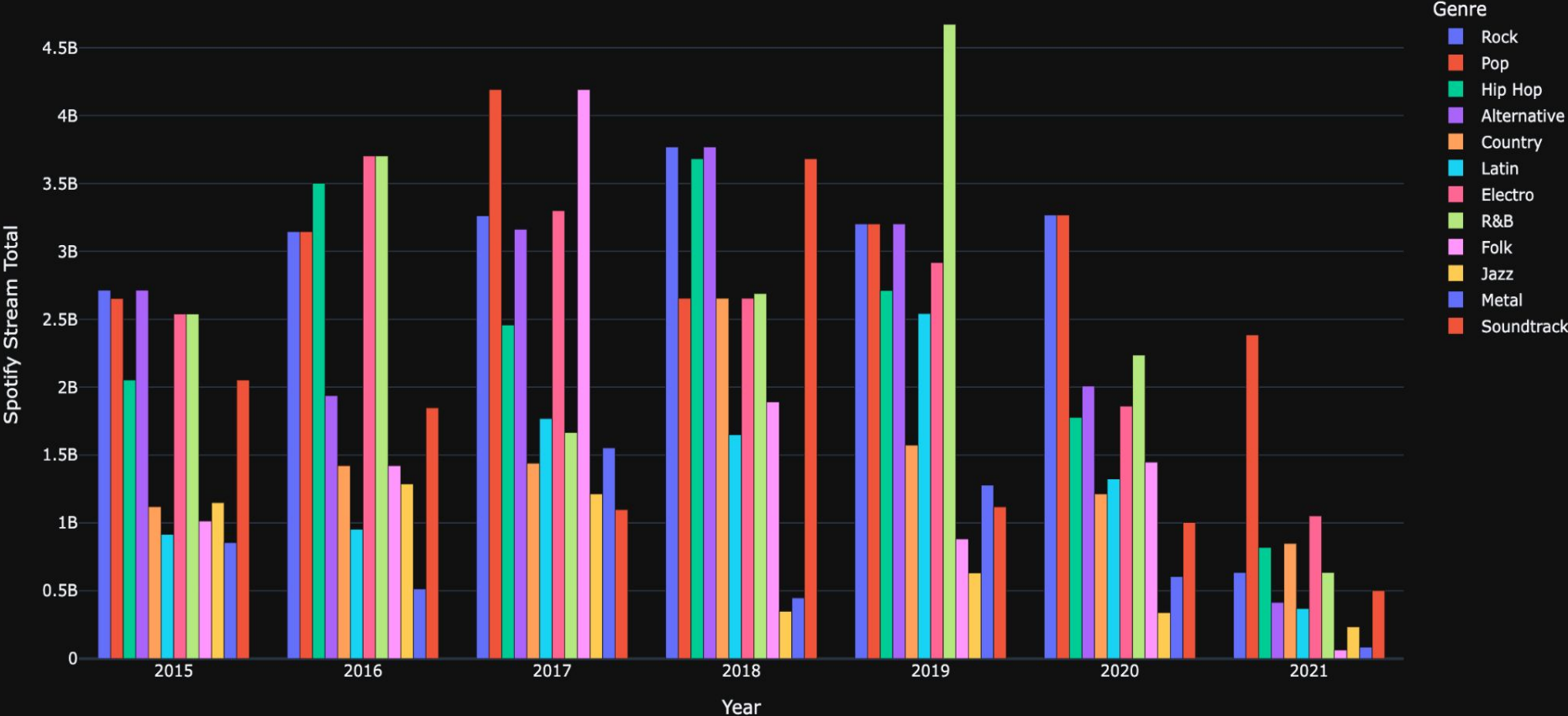
```
fig combined.write_html("html/spotify genre streams by year grouped.html")
```

```
# Show the plot
```

```
fig combined.show()
```

Total Spotify Streams by Genres

[interactive link](#)



Results of Visualizations

- Pop is the most popular genre every year except 2018 where hip-hop surpassed it
- Total streams peaks in 2019
- The decline of rock music is notable being the 2nd most popular in 2015 and 2016 to 6th most popular in 2023
- The rise of the Latin genre from 7th most popular in 2015 to 3rd most in 2023
- 2024 Alternative having more streams than Rock
- Genres in the top 10 every year : Pop, Rock, Electro, R&B, Alternative, Latin, Hip-Hop, and Country

- Initial Data Cleaning.....
- Regex used to sort for just the numerical 4-digit year in `Project_3_test2.ipynb`
 - This was a unique case due to the data in the 'Release date'. A simple `.fillna` approach wouldn't work. We had to first fill NaN values with a string version of '1776' to avoid breaking the regex when encountering an integer instead of a string.
 - Steps:

Fill NaN values in 'Release date' with '1776' and convert to string:

```
df_cleaned['Year Released'] = df_cleaned['Release date'].fillna('1776').astype(str)
```

■

Extract the year from 'Year Released' using regex:

```
df_cleaned['Year Released'] = df_cleaned['Year Released'].str.extract(r'(\d{4})')
```

■

Fill any remaining NaN values in 'Year Released' with '1776':

```
df_cleaned['Year Released'] = df_cleaned['Year Released'].fillna('1776')
```

■

- Convert the 'Year Released' column to integers for further mathematical operations:

```
df_cleaned['Year Released'] = df_cleaned['Year Released'].astype(int)
```

```
from tabulate import tabulate
```

```
# Pretty print the first 20 rows of  
the DataFrame using tabulate
```

```
print(tabulate(df_cleaned.head(20),  
headers='keys', tablefmt='pretty',  
showindex=False))
```

Using tabulate sample

```

tabulate_df = pd.read_csv('Cleaned.csv')
print(tabulate(tabulate_df.head(15), headers='keys', tablefmt='pretty',
showindex=False))

```

✓ 0.1s

Python

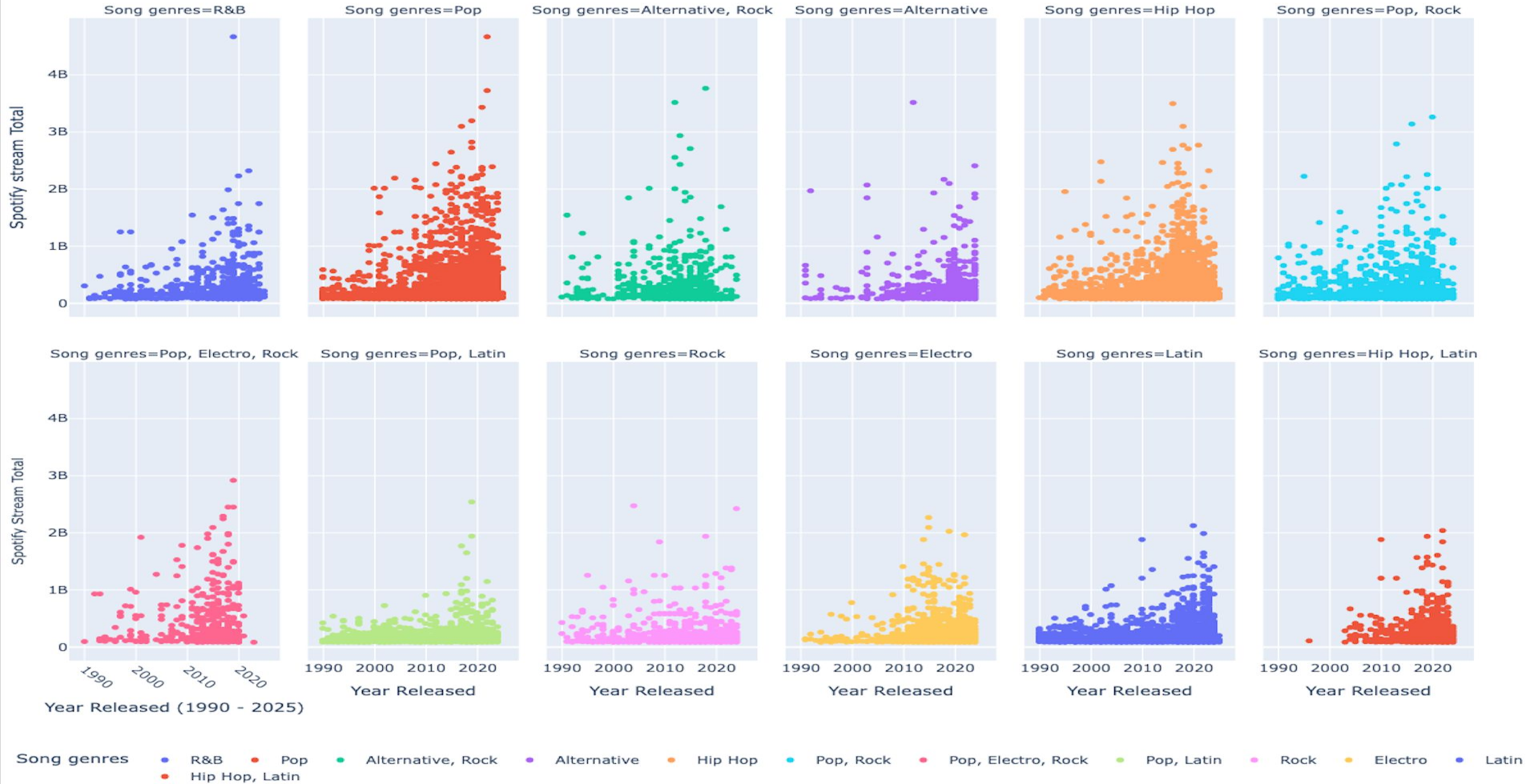
Song genres	Artist(s)	Artist Country	Spotify stream Total
R&B	The Weeknd	Canada	4673836529
Pop	Batusyex, imnotege, Lanceas, rufflws, xSyborg	Turkey, Turkey, Germany, Turkey	4673836529
Pop, Folk	Ed Sheeran	United Kingdom	4193269875
Alternative, Rock	Lewis Capaldi	United Kingdom	3770959284
Pop	Harry Styles	United Kingdom	3731855719
R&B, Electro	Daft Punk, The Weeknd	France, Canada	3704089716
R&B, Electro	Daft Punk, The Weeknd	France, Canada	3704089716
Hip Hop, Soundtrack	Post Malone, Swae Lee	United States, United States	3683693906
Alternative, Rock	The Neighbourhood	United States	3524300277
Alternative	The Neighbourhood	United States	3524300277
Hip Hop	Drake, Kyla, WizKid	Canada, United Kingdom, Nigeria	3503488880
Pop	Justin Bieber, The Kid Laroi	Canada, Australia	3439352937
Pop, Electro, Hip Hop	Justin Bieber, The Kid Laroi	Canada, Australia	3439352937
Electro, Pop	Imagine Dragons	United States	3300915632
Pop, Rock	Glass Animals	United Kingdom	3268317021

2 Quick examples of using facet in the plotly library

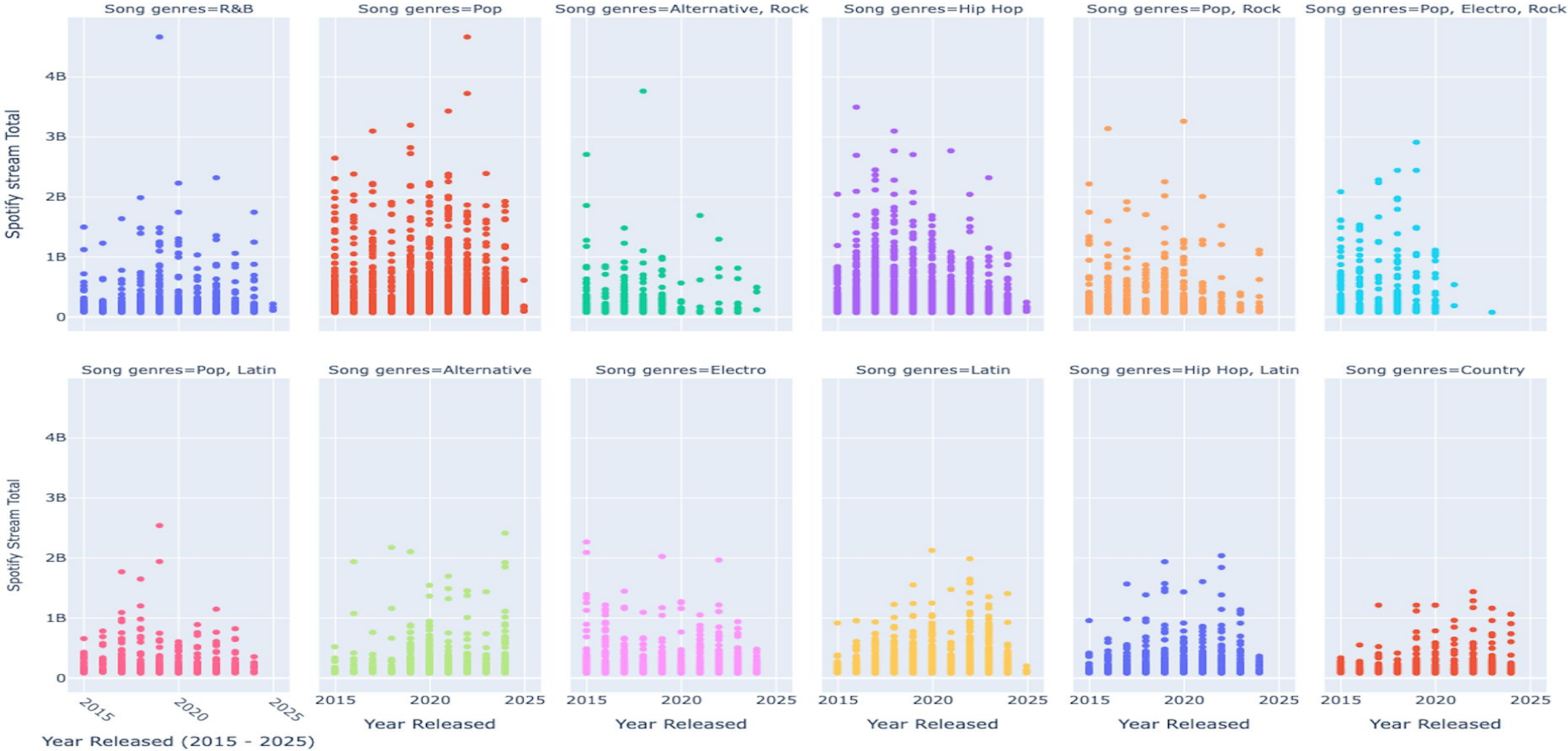
```
import plotly.express as px
```

Using Facets with Plotly examples

Spotify Stream Total by Year Released and Genre



Spotify Stream Total by Year Released and Genre

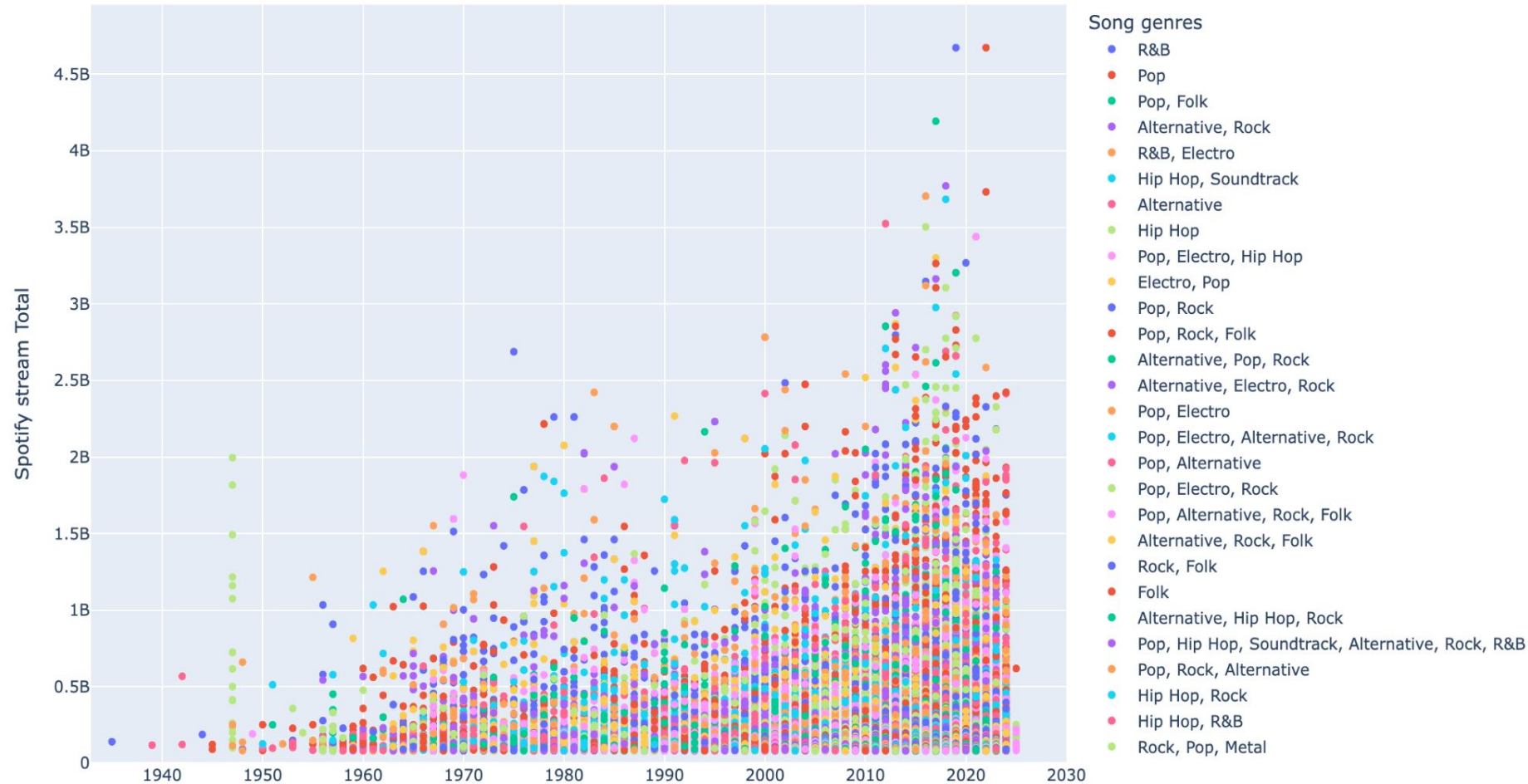


Song genres

- R&B
- Pop
- Alternative, Rock
- Hip Hop
- Pop, Rock
- Pop, Electro, Rock
- Pop, Latin
- Alternative
- Electro
- Latin
- Hip Hop, Latin
- Country

Spotify Stream Total by Year Released and Song Genre

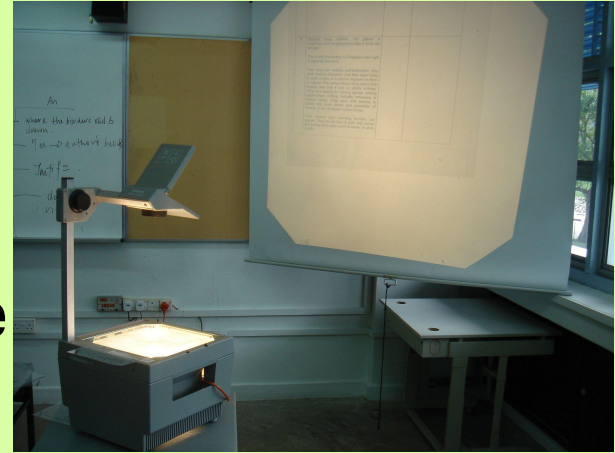
[Spotify Stream Total link](#)



Notes on Visualizations

- Interesting to note what genres often get combined when categorizing a song
- Number of streams peak between 2018 to 2020
- Any guess what those two dots at the peak of the chart are? (hint: they are the same song)

For the record.
Now this is what I think of as a
slide show.
Not to be confused with the picture
slide show.



By mailer_diablo - Self-taken (Unmodified), CC BY-SA 3.0,
<https://commons.wikimedia.org/w/index.php?curid=525127>

Outro