

King Saud University  
College of Computer and Information Sciences  
Department of Information Technology



IT 362 Course Project  
Semester-1, 1447H

Student Name	Student ID
Dana Alsalami	443200945
Norah Aljayan	444200832
Aljohara Alsultan	444203635
Layan Alhowaimel	444200969
Alya Almuqren	444200610

Course Instructor: **Dr. Abeer Aldayel**

## Table of Contents

Introduction .....	3
Data sources.....	3
Objectives.....	5
Method.....	5
Challenges .....	9
Raw Data File(s).....	10
A Jupyter Notebook .....	10
Logbook.....	<b>Error! Bookmark not defined.</b>

## Table of Figures

Figure 1:Boxplot .....	7
Figure 2:Countplot .....	7
Figure 3:Scatter plot.....	8

## Introduction

The music industry has witnessed increasing fluidity in genre boundaries, with many artists exploring multiple styles rather than remaining confined to a single category. This raises an important question about the relationship between artistic versatility and success: **Do artists who span multiple genres have higher popularity compared to those who focus on a single genre?**

To address this research question, we employ data retrieved from Spotify Charts and enriched with the Spotify API. First, we collected weekly regional chart CSVs from multiple countries, each containing the top 200 tracks. From these charts, we extracted unique artist names. Using a custom script, we then queried the Spotify API to gather additional metadata for each artist, including their Spotify ID, associated genres, and popularity scores (0–100 scale). This integrated dataset forms the basis for analyzing whether artists spanning multiple genres tend to achieve higher popularity among listeners.

The study will evaluate popularity trends across multi-genre versus single-genre artists, and highlight implications for the modern music industry. The overall aim is not only to assess whether versatility contributes to greater popularity, but also to offer insights into how genre diversity and digital streaming platforms shape contemporary musical success.

## Data sources

For our project on artist popularity and genre diversity, we relied primarily on two data sources: **Spotify Charts** and **Spotify Web API**.

### Spotify Charts (CSV files)

Spotify Charts are weekly and daily charts published by Spotify, listing the top ten tracks in various countries and globally. We downloaded the weekly top tracks of the following countries: Global, US, Saudi Arabia (SA), Belarus (BY), Italy (IT), Brazil (BR), India (IN), United Kingdom (UK), Taiwan (TW), Switzerland (CH), South Korea (KR), Japan (JP), Venezuela (VE), Uruguay (UY), Turkey (TR), Thailand (TH), Australia (AU), United Arab Emirates (AE), Hong Kong (HK), Egypt (EG).

Data collected: Artist names, track names, chart rank, streams, weeks on chart, peak rank.

Number of observations: After consolidating all downloaded CSV files, there are approximately 1785 unique artists.

Features & types:

Track\_name (string): Name of track

artist\_names (string): Name of the artist

rank(integer): Current rank of the chart (integer)

streams(integer): Number of streams of the chart

weeks\_on\_chart(integer): Number of weeks the track has been on the chart

peak\_rank(integer): Highest rank achieved on the chart

previous\_rank(integer): the tracks position in the previous week

source (string): chart source (country/region) The only relevant feature to our hypothesis is the artist\_names feature

## **Spotify API**

We used the Spotipy Python library with a developer Spotify account to retrieve detailed artist information.

Data collected: Artist ID, genres, and popularity metric

Features & types: id (String): Spotify unique identifier for the artist

name (string): artist name

genres (list): list of genres associated with the artist

popularity(integer): Spotify popularity score (0-100) calculated based on recent streams and engagement

Representation Bias: Each chart contains 200 entries, but the underlying listener populations differ across regions. Large markets like the US, UK, and Global charts reflect millions of streams, while smaller markets like Uruguay or Thailand reflect fewer streams. The dataset includes only **unique artists**, so each artist appears once even if present in multiple charts. This reduces the overrepresentation of globally popular artists. However, artists from smaller markets are still less likely to appear, and genres or languages favored in large markets may dominate.

Measurement Bias: Spotify popularity is influenced by recent streams and engagement. New releases may temporarily inflate an artist's popularity, while older or niche artists may be underrepresented.

Data Limitations: Artists with multiple pseudonyms or collaborations may appear under different names. Artists without genres listed in the Spotify API are excluded, biasing the dataset toward mainstream or well-categorized artists.

## Objectives

1. **To investigate the relationship between genre diversity and popularity** by comparing artists who span multiple genres with those limited to a single genre.
2. **To collect and analyze data from the Spotify API**, including artist names, IDs, genres, and popularity scores, in order to build a reliable dataset for evaluation.
3. **To identify potential trends or correlations** between cross-genre versatility and higher popularity metrics.
4. **To examine possible biases or limitations** in the dataset, such as unequal representation of artists or genre classifications.
5. **To provide insights into the role of genre diversity** in shaping artist success and its implications for the modern music industry

## Method

### Data Collection

Weekly Spotify chart data from 20 regions, including Global, US, UK, Saudi Arabia, Japan, Brazil, and others, were collected by downloading CSV files containing track names, artist names, chart ranks, streams, weeks on chart, and peak positions. The CSV files were programmatically imported and merged into a single dataset using Python. Duplicate artist entries were removed to ensure that each artist appeared only once across all charts. From this consolidated dataset, all unique artist names were extracted. Using the Spotipy Python library with Spotify Web API credentials, detailed artist metadata was retrieved for each unique artist,

including Spotify ID, genres, and popularity scores. Artists without genre information or duplicate IDs were excluded to maintain data consistency. Error handling was implemented to manage connection errors and API exceptions during data retrieval. The resulting dataset was stored in a Pandas Data Frame containing only artists with non-empty genres and valid popularity metrics.

### **Data Preprocessing**

Preprocessing consisted of merging multiple CSV files, filtering out duplicate artists, and removing any artists with missing genre information. This ensured that the dataset was consistent and ready for analysis. All files were stored locally to allow reproducibility without repeated API requests.

## Data Analysis and Visualization by Objective

**Objective 1:** *Investigate the relationship between genre diversity and popularity.*

- We created a new feature `genre_count` to indicate how many genres each artist spans. Artists were categorized as **single-genre** or **multi-genre**. Popularity distributions were compared using **boxplots**, and the number of artists per category was visualized using **countplots**.

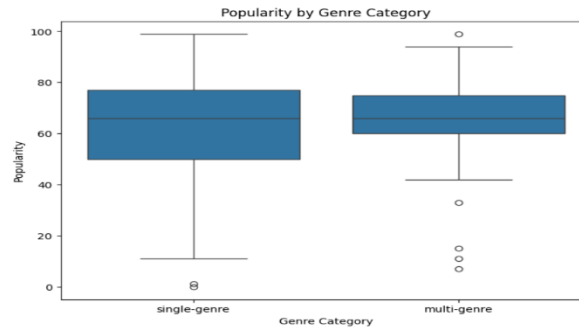


Figure 1:Boxplot

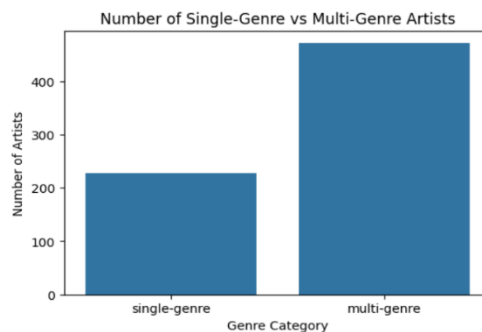


Figure 2:Countplot

**Objective 2:** *Collect and analyze data from the Spotify API, including artist names, IDs, genres, and popularity scores.*

- CSV files from 20 regions were merged, and unique artist names were queried from the Spotify API using Spotipy to retrieve Spotify IDs, genres, and popularity scores. Artists

without genres or duplicates were excluded. The final dataset was stored in a Pandas DataFrame for analysis.

**Objective 3:** *Identify potential trends or correlations between cross-genre versatility and popularity metrics.*

- Correlation between genre\_count and popularity was examined. Trends were visualized using **scatter plots** and **regression plots** to assess whether spanning more genres is associated with higher popularity.

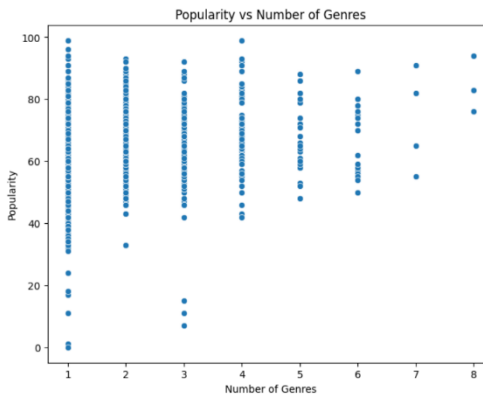


Figure 3: Scatter plot

**Objective 4:** *Examine potential biases or limitations in the dataset.*

- The dataset includes only unique artists across all charts, which reduces the overrepresentation of globally popular artists but also means that artists from smaller regions are less likely to appear overall. Since each chart contains 200 entries but represents very different listener populations, genres and languages from larger markets (such as the US or UK) may dominate the dataset compared to smaller markets like Uruguay or Thailand. Spotify's popularity score is also shaped by recent streams, so new releases may temporarily inflate popularity while older or niche artists are underrepresented. Finally, artists without listed genres in the Spotify API were excluded, which introduces a bias toward mainstream or well-categorized artists, and artists using



pseudonyms or frequent collaborations may appear multiple times under different names. Visualizations such as **countplots** of genre categories help illustrate potential representation biases.

**Objective 5:** *Provide insights into the role of genre diversity in shaping artist success.*

- Descriptive statistics, histograms of popularity, boxplots, scatter plots, and regression plots were used to explore and compare popularity across single-genre and multi-genre artists. These analyses reveal patterns linking genre diversity to popularity and provide insights into its impact on modern music industry trends.

## Tools and Libraries

Python libraries used in the project include:

- **spotipy**: Retrieve artist metadata from Spotify API
- **pandas**: Data manipulation and DataFrame creation
- **os** and **glob**: File handling and importing multiple CSVs
- **requests**: Downloading CSV files
- **time**: Managing pauses between API requests
- **spotipy.exceptions** and **requests.exceptions**: Error handling during API calls
- **matplotlib** and **seaborn**: Data visualization for descriptive and comparative analysis

## Challenges

### 1. API Limitations and Errors

Collecting data for hundreds of artists required many requests, which often hit rate limits or caused timeouts.

### 2. Incomplete Genre Information

Some artists did not have genre data listed in the Spotify API, leading to their exclusion and reducing the dataset's diversity.

### 3. Representation Bias

Larger markets (e.g., US, UK, Global) dominated the charts, while smaller markets were underrepresented, potentially biasing results toward globally popular genres.

### 4. Limited Data for Emerging Artists

Less popular or newer artists often have incomplete data or fewer genre tags, which can bias comparisons between artists with many genres and those with few.

## **5. Duplicate/alias issues**

Artists with multiple pseudonyms, collaborations, or spelling variations may appear under different IDs, fragmenting the dataset.

### Raw Data File(s)

- <https://github.com/aljoharas/datascience/tree/main/charts>

### A Jupyter Notebook

- <https://github.com/aljoharas/datascience/blob/main/spotifyproject.ipynb>