# Exploratory Data Analysis on Spotify Tracks

Subtitle: Data Science Project

Name: Arka Koley

Tools Used: Python (Pandas, Seaborn, Matplotlib)

Date: October 2025

# Introduction

- **This project explores Spotify track data using Exploratory Data Analysis (EDA).**

- **EDA is crucial for understanding patterns, trends, and relationships within a dataset before building any formal models.**

- **Objective: To identify the features that contribute to a song's popularity and to analyze how musical trends have evolved over time.**

# About the Dataset

The dataset contains a comprehensive collection of Spotify tracks, each detailed with over 20 distinct feature

## Total Records: ~50,000 tracks

## Key Columns:
Identification: track_name, artist_name, year
Performance Metric: popularity

- 
- 
- 
- Audio Features: danceability, energy, acousticness, loudness, speechiness, valence, tempo, liveness, instrumentalness

# EDA Objectives

**Our analysis is structured around the following objectives:**

1. **Univariate Analysis:** Explore the distribution and characteristics of individual features.
2. **Bivariate Analysis:** Study the relationships and interactions between pairs of features.
3. **Multivariate & Correlation Analysis:** Examine the correlations between multiple features simultaneously to uncover complex patterns.
4. **Time Series Analysis:** Track and analyze how musical features and popularity have changed across different years.
5. **Generate Insights & Recommendations:** Synthesize findings to provide actionable insights.
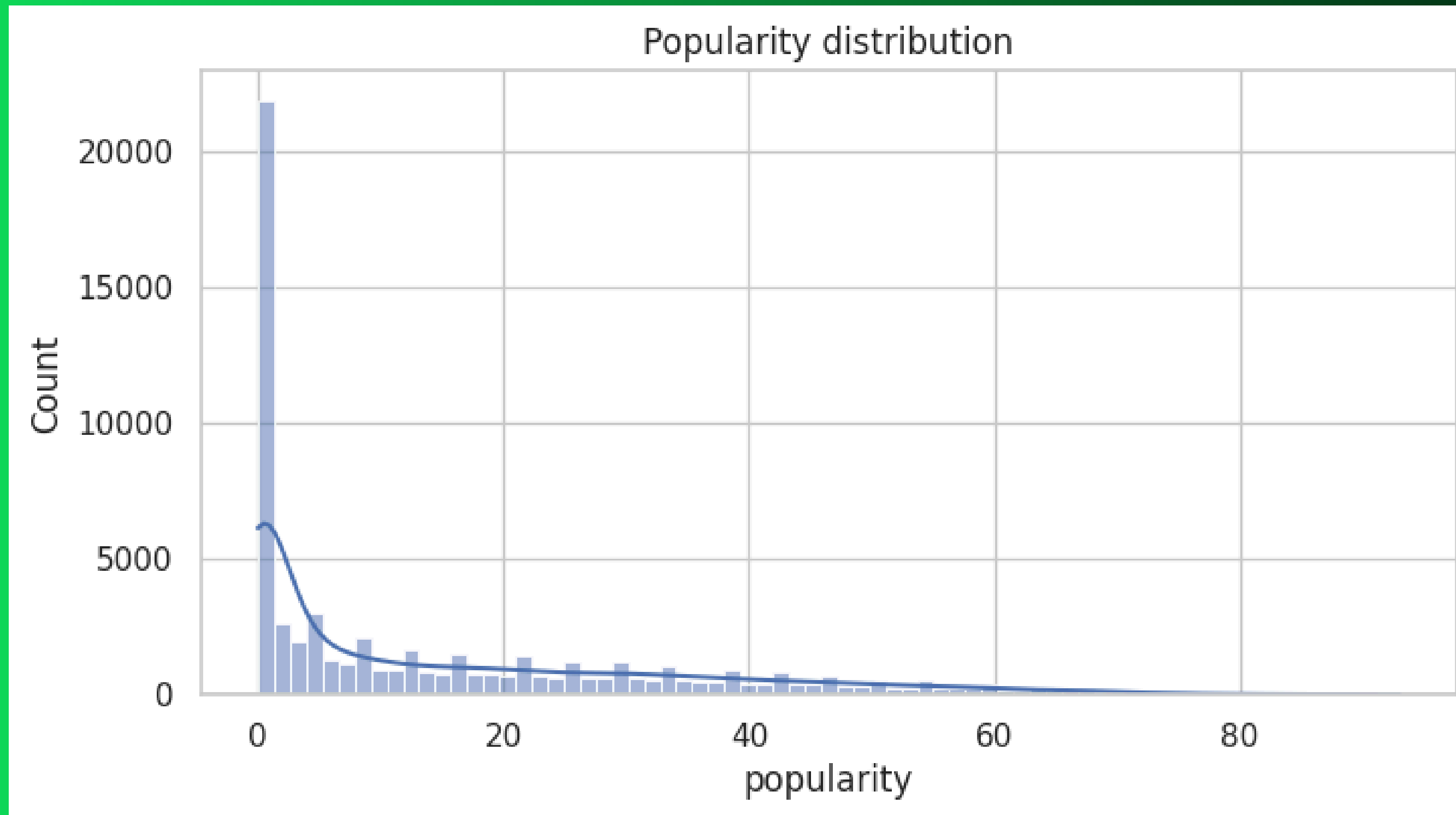
# Data Cleaning & Preprocessing

To ensure the quality and reliability of our analysis, the following data cleaning steps were performed:

- Duplicate Removal: All duplicate rows were identified and removed.
- Missing Values: The dataset was checked for any missing values, which were handled accordingly.
- Data Type Conversion: Numeric columns were converted to their appropriate data types for accurate calculations.
- Data Verification: Column values were checked to ensure they fell within expected ranges and formats.

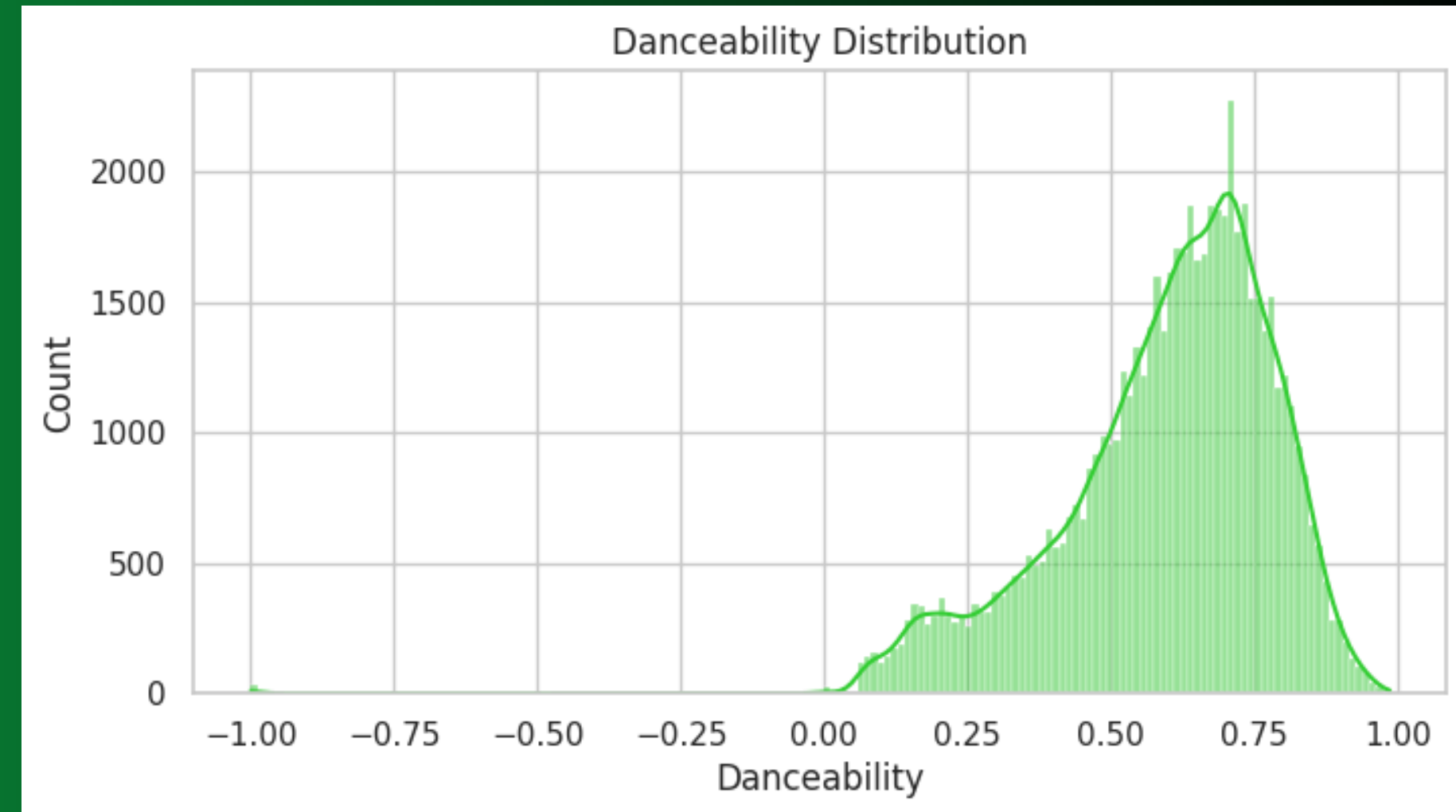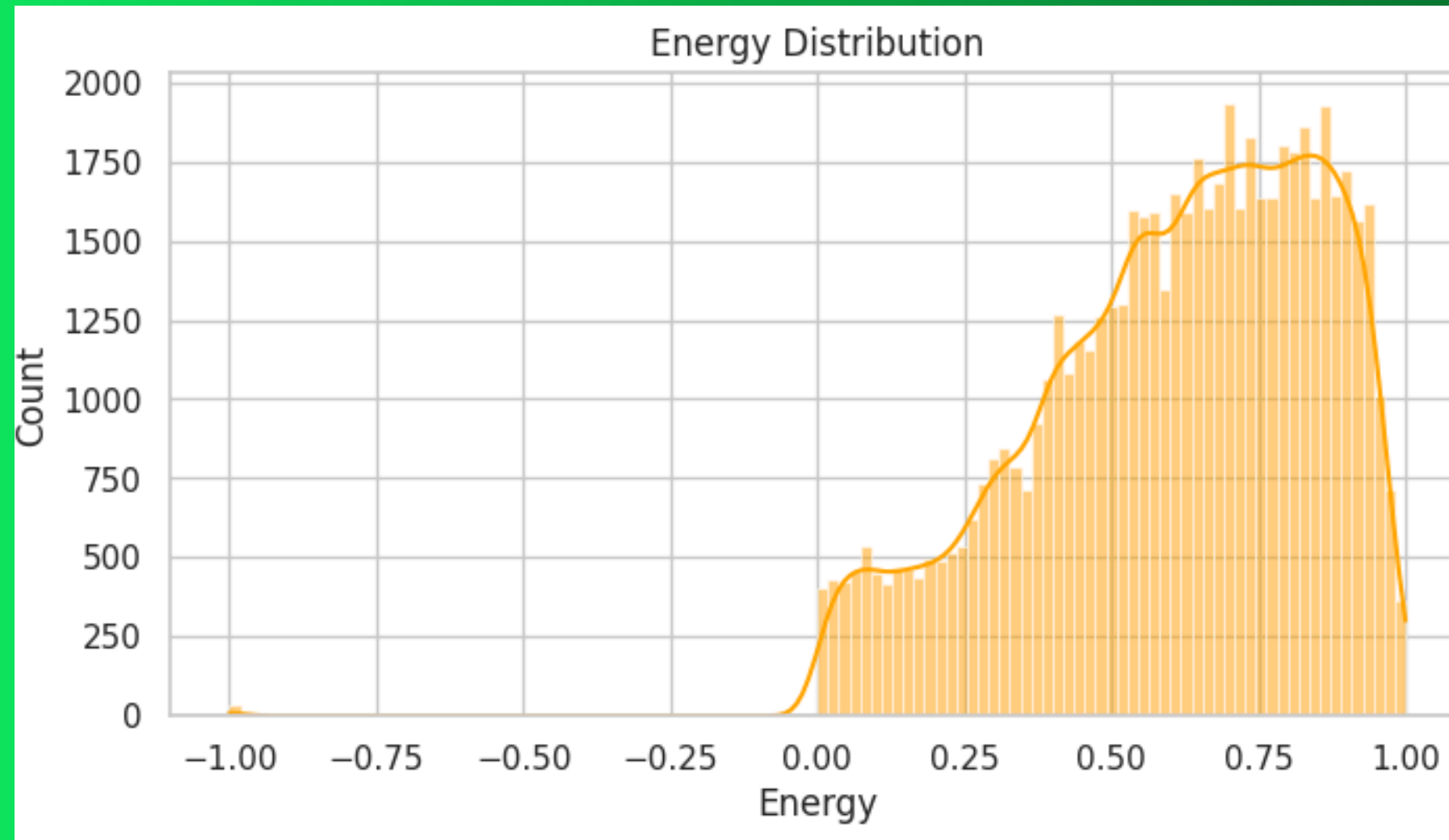# Univariate Analysis



Popularity distribution

Most songs have moderate popularity, clustering in the 40–70 range. This means that scoring an extremely high popularity (>90) is rare, making the hits stand out significantly.
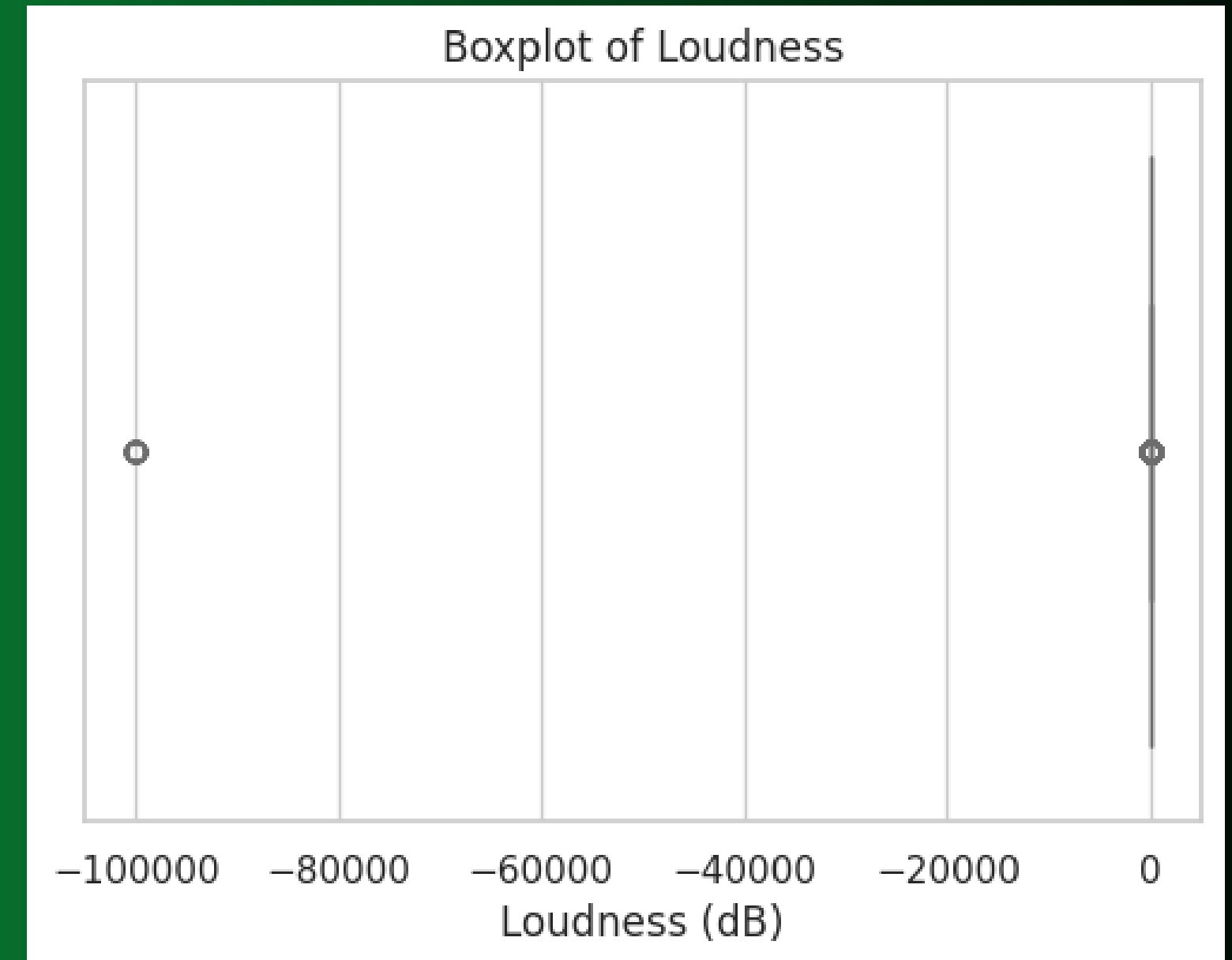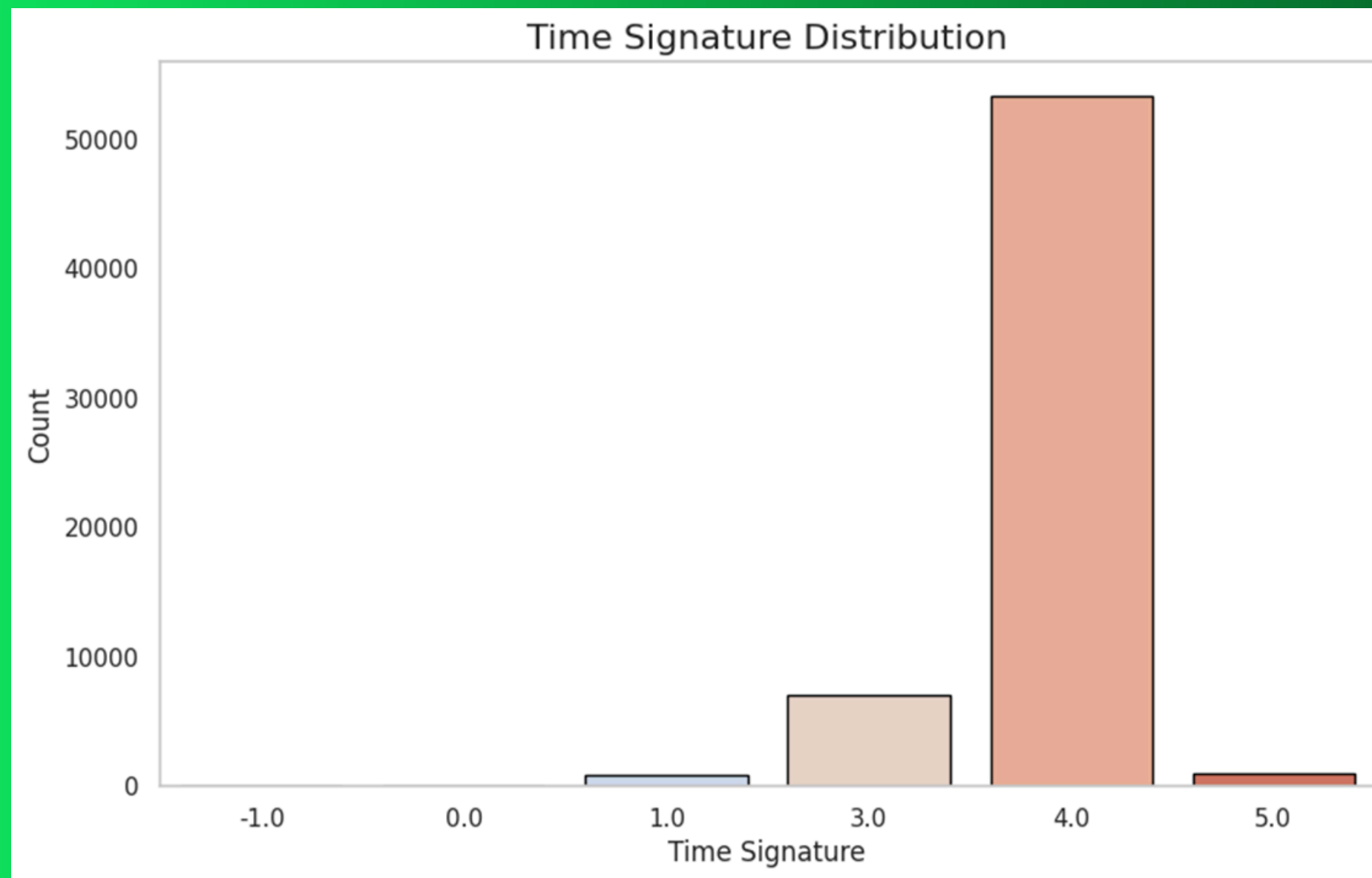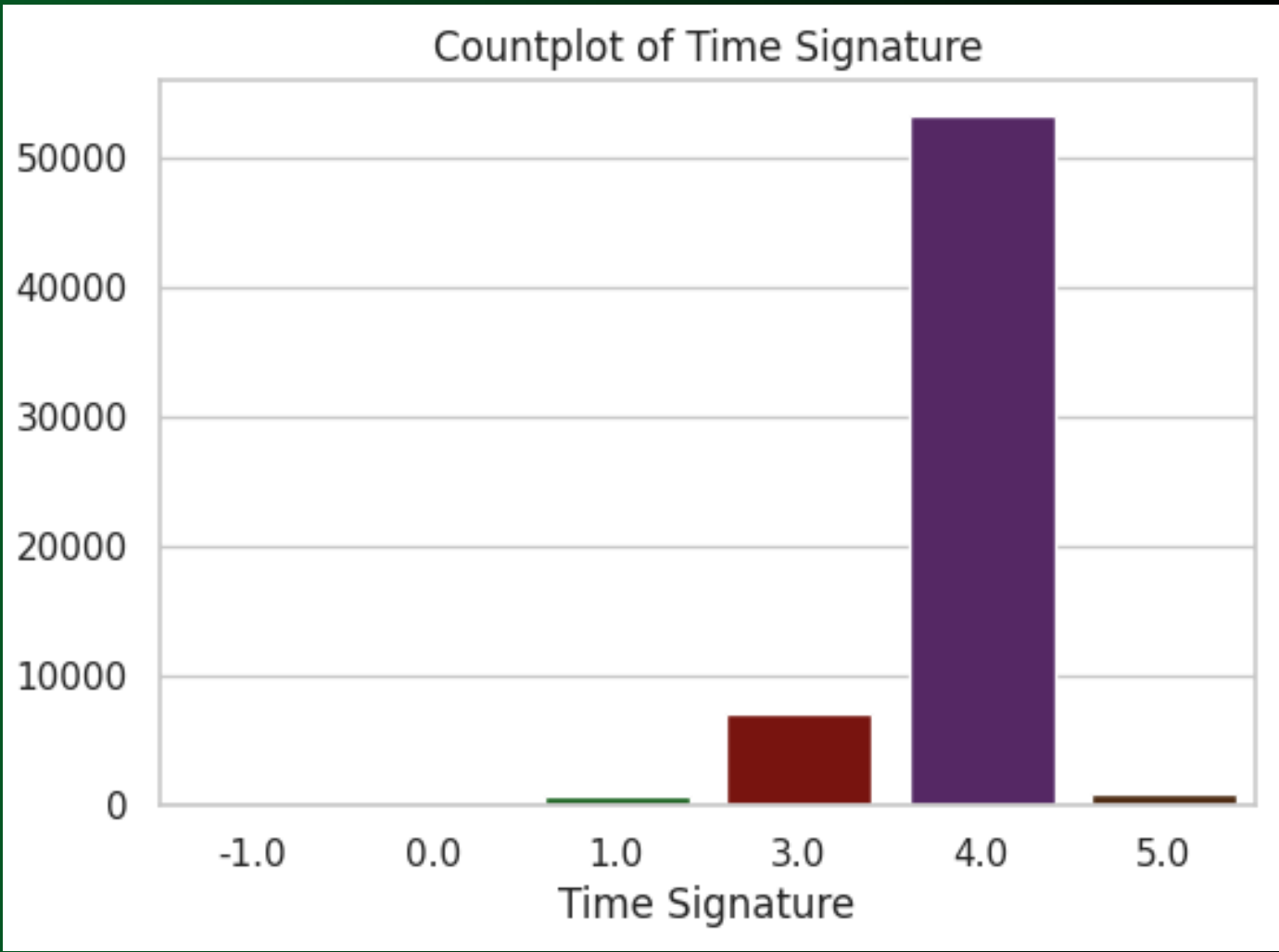
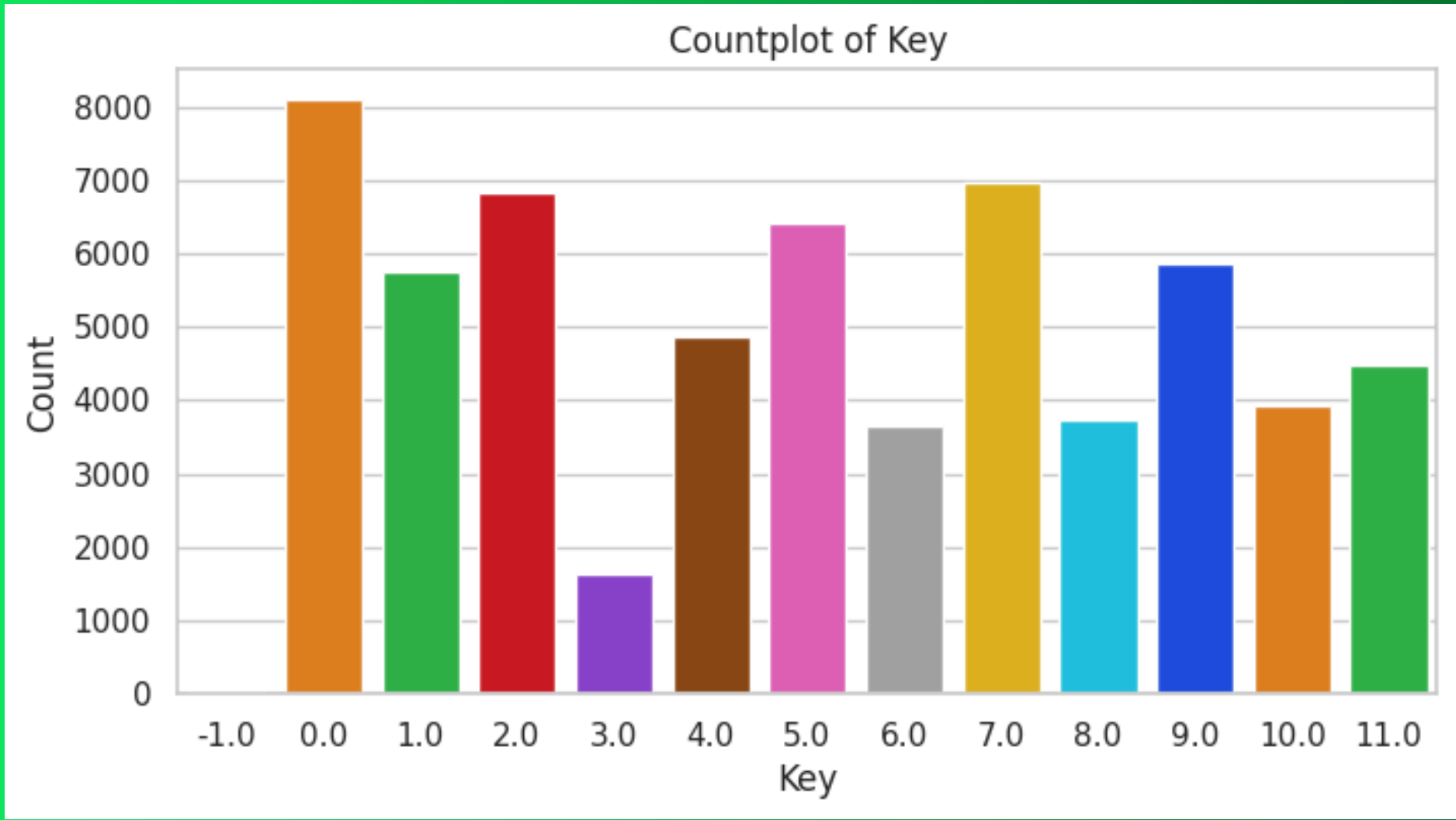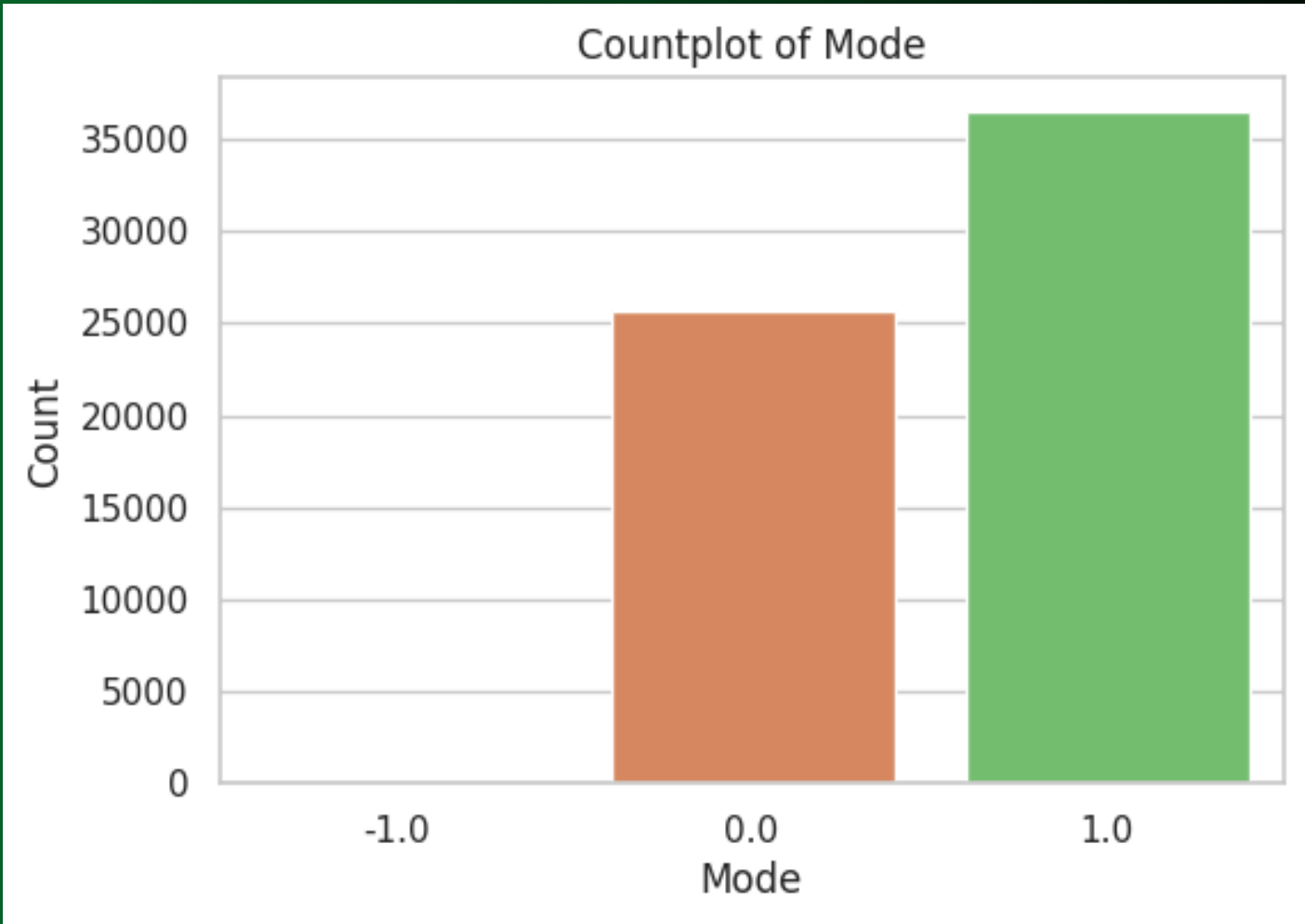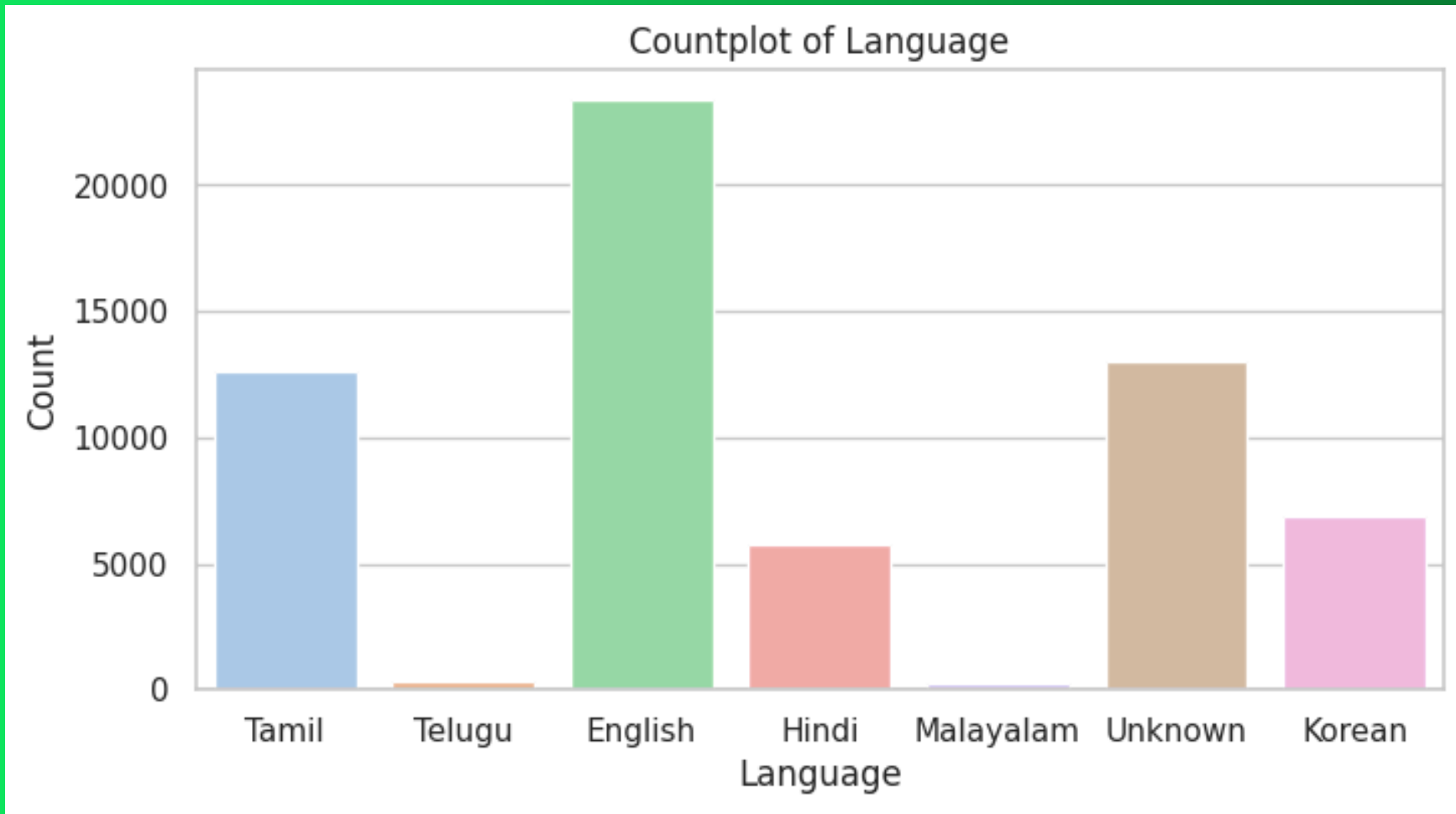# The Vibe Features: Energy, Danceability, and Mood



The majority of tracks are characterized by high danceability (scores >0.6).Similar to danceability, songs tend to have high energy (scores >0.7).
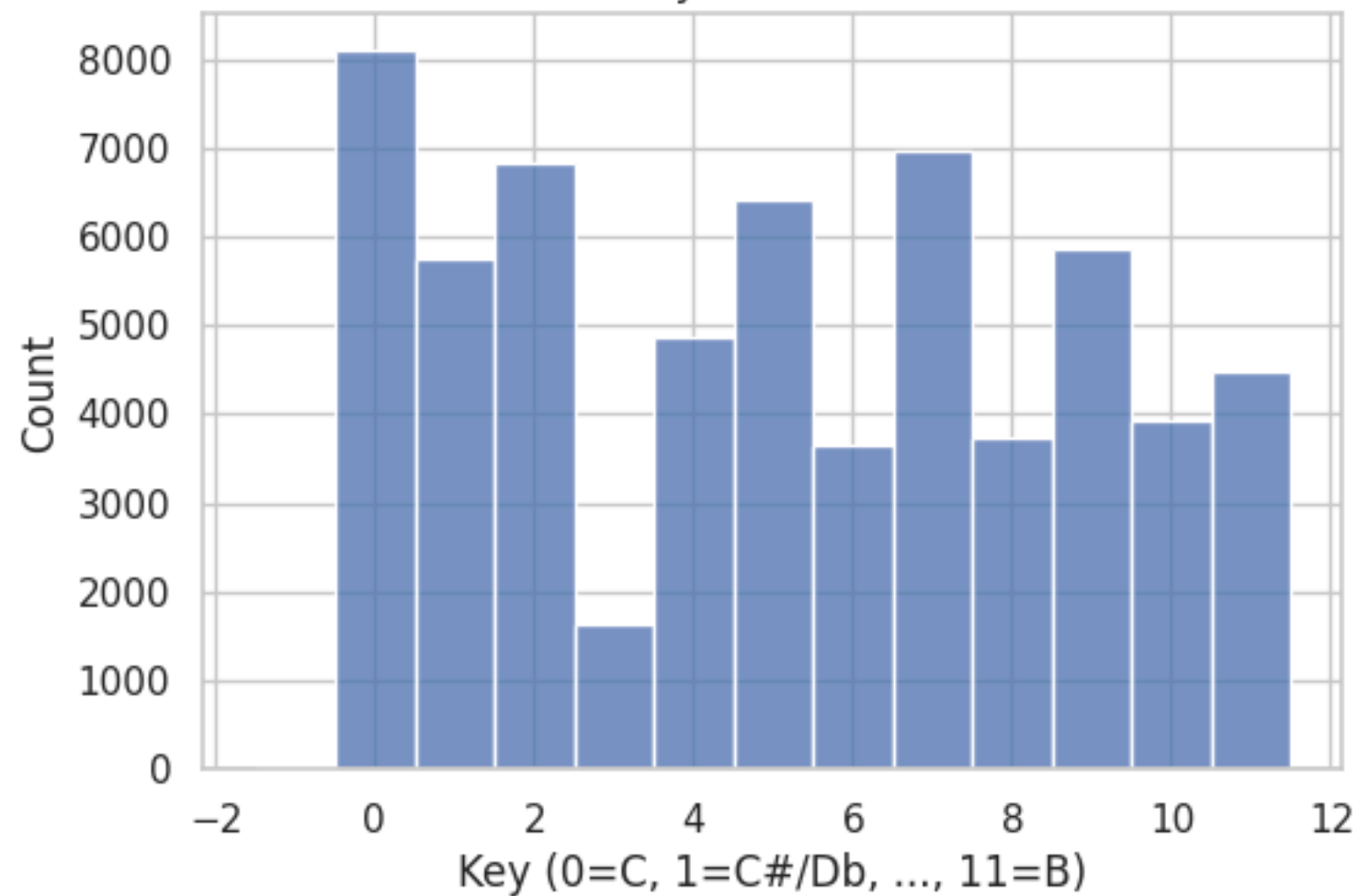
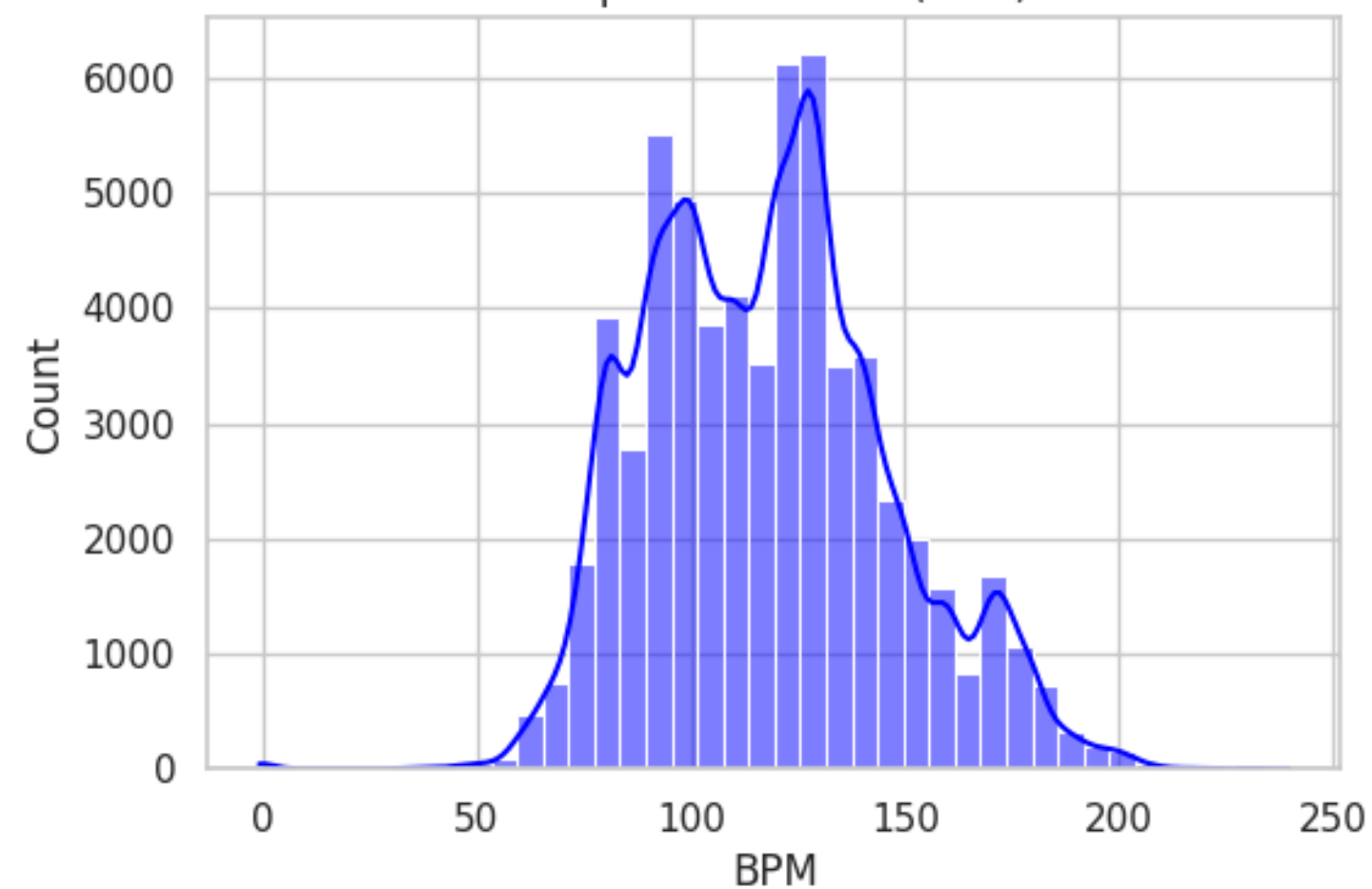# Structural Features: Tempo, Duration, and Loudness



The average song length is around 3.5 minutes. Loudness distribution will likely be tightly packed in the high-decibel range (e.g., -5 dB to -10 dB).Tempo histogram confirms that music often follows popular beats per minute.
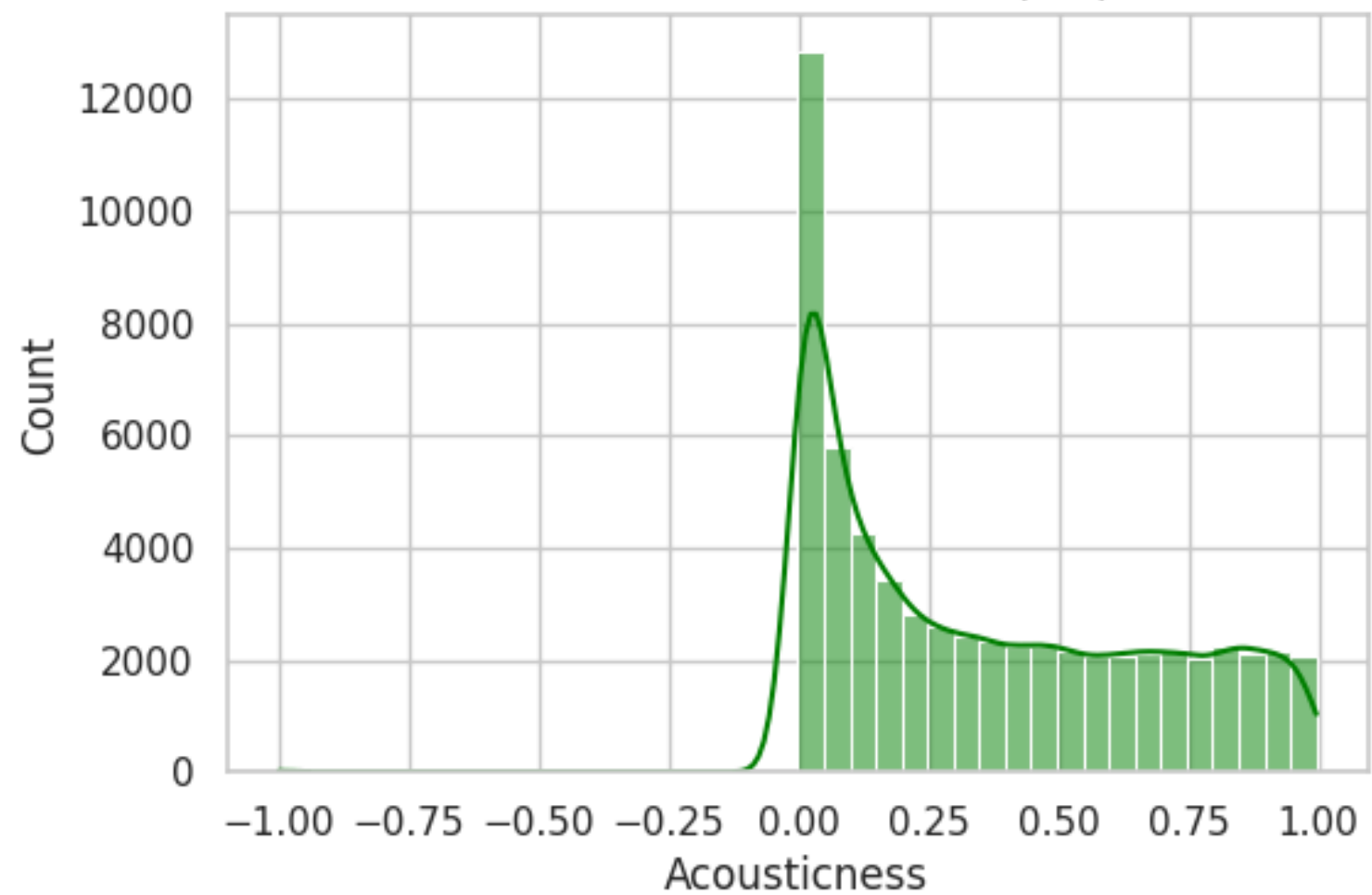
Countplot of Language

Countplot of Mode

Countplot of Key

Countplot of Time Signature
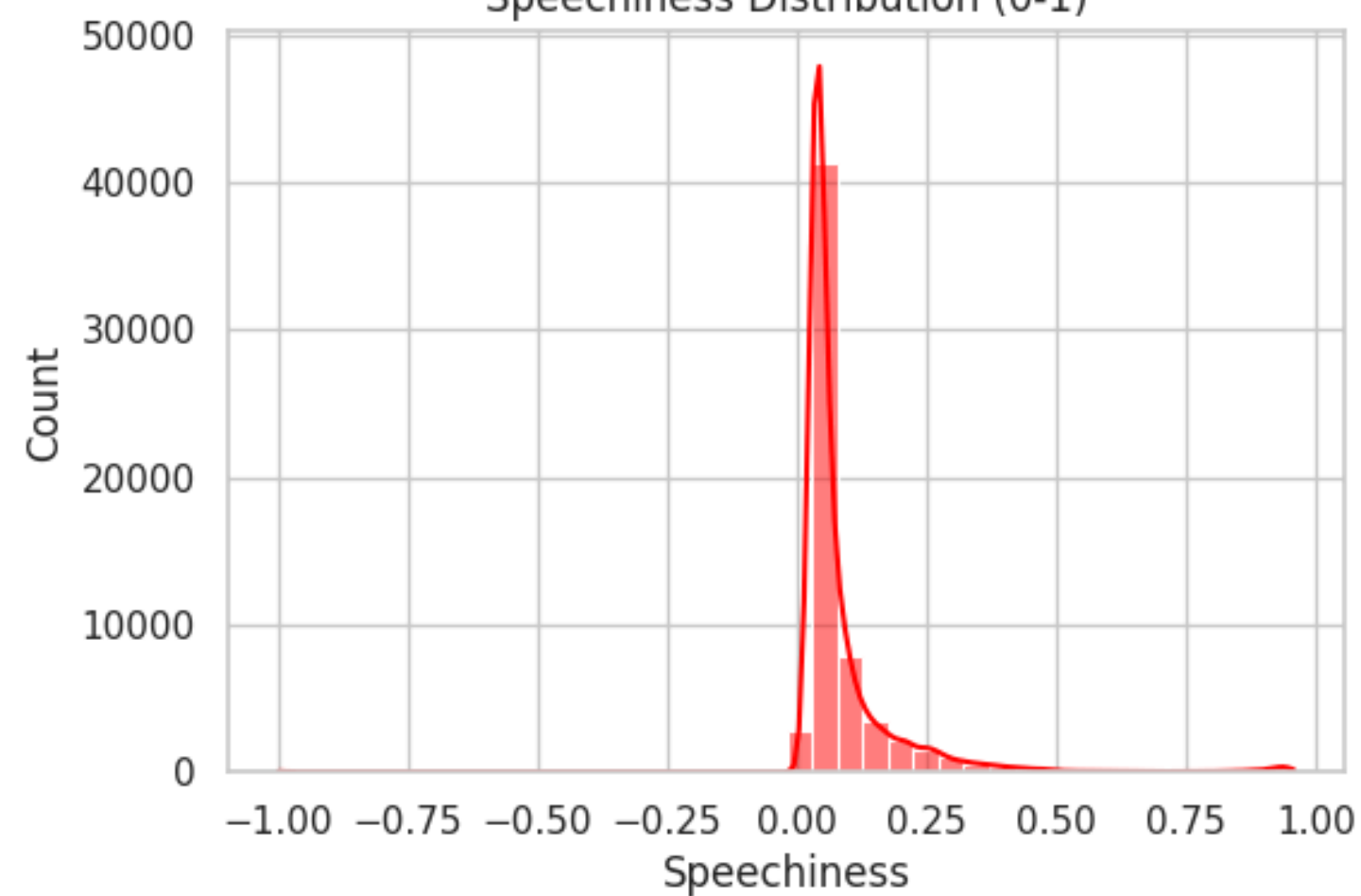
Musical Key Distribution (0-11)

Tempo Distribution (BPM)

Acousticness Distribution (0-1)

Speechiness Distribution (0-1)

**Univariate Analysis (Categorical Variable)**



Share of Languages in Spotify Tracks

Korean 11.1%

Hindi 9.2%

Other 1.0%

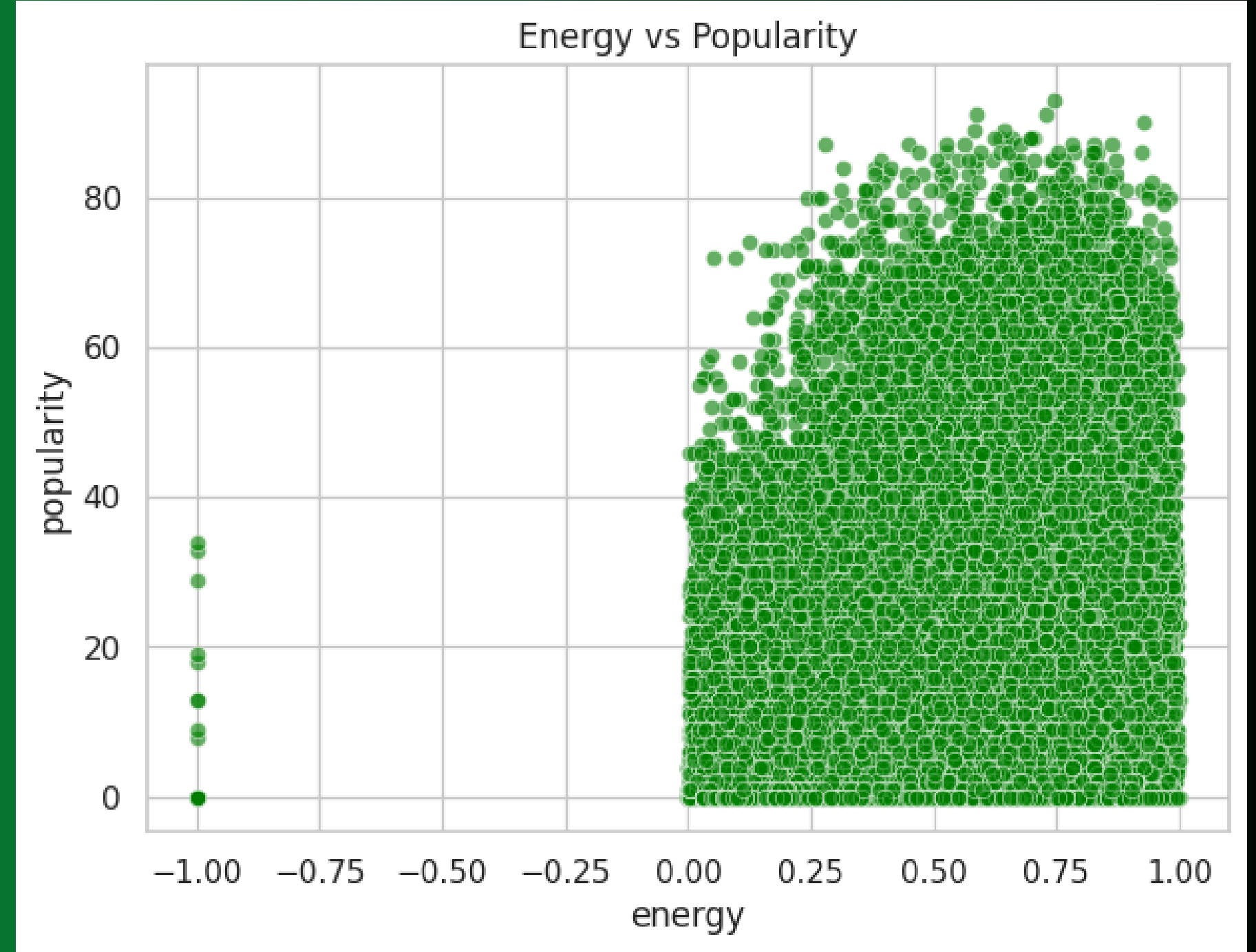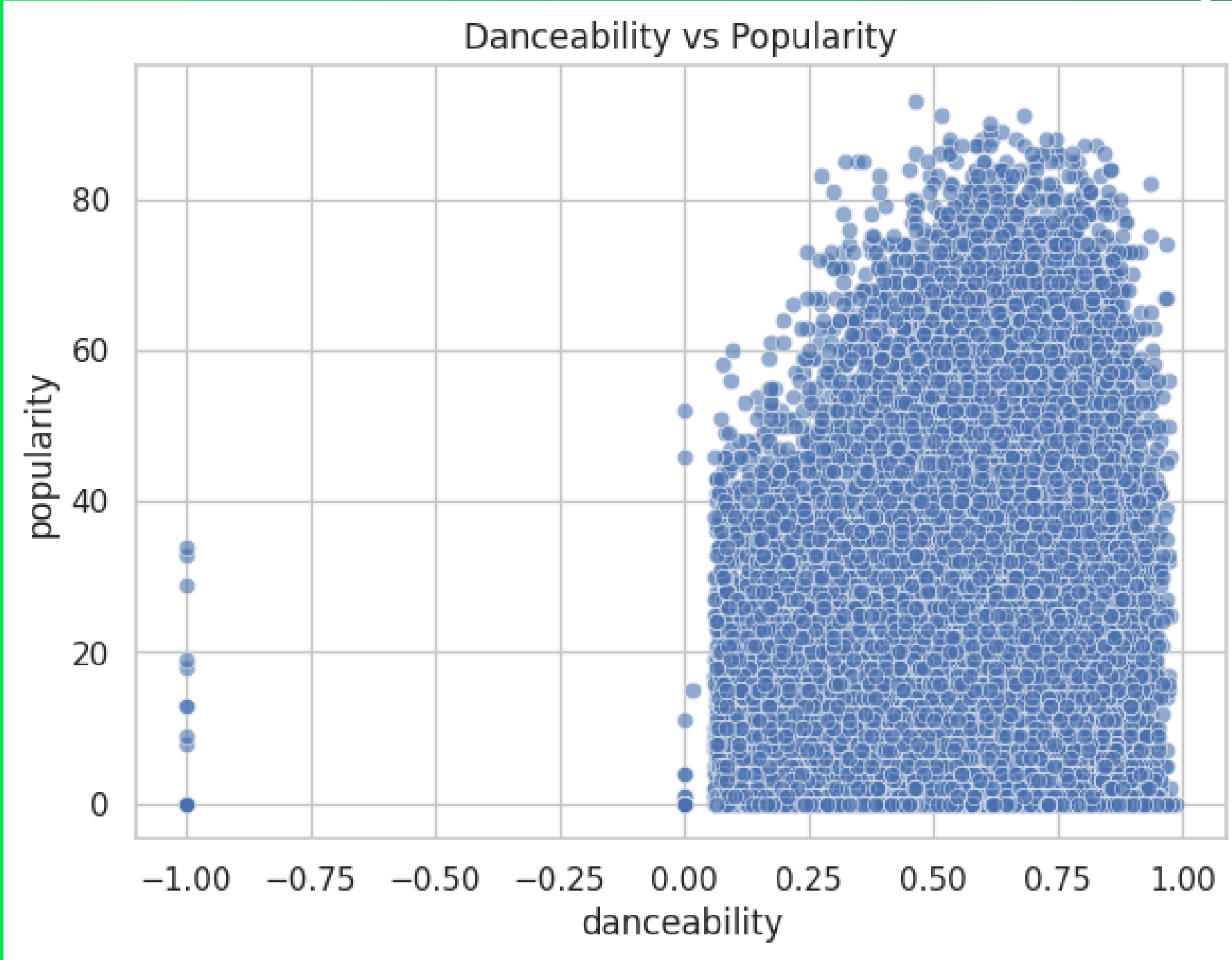Tamil 20.3%

English 37.6%

Unknown 20.9%

# Univariate Analysis (Results)

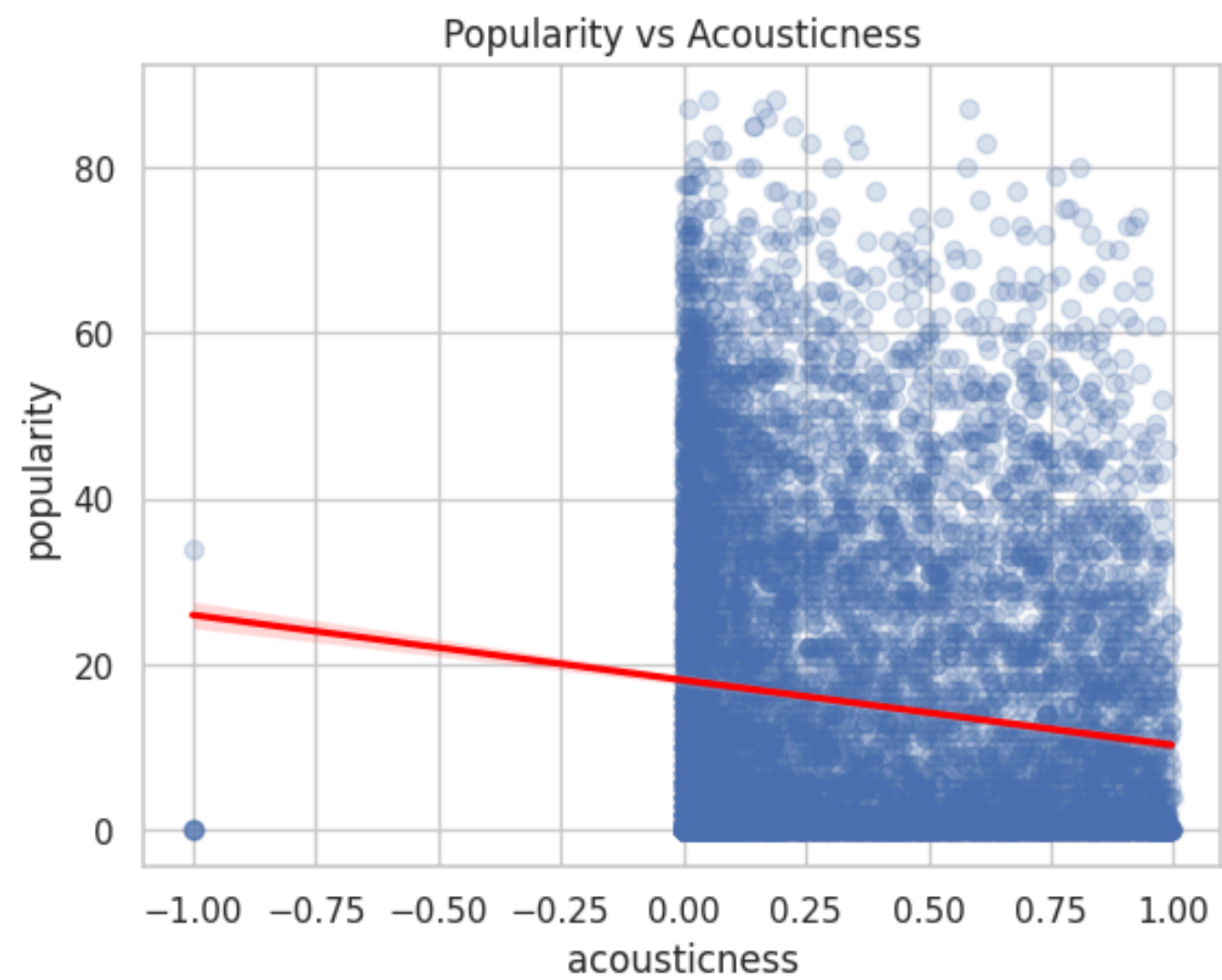Initial analysis of individual features revealed several key trends:
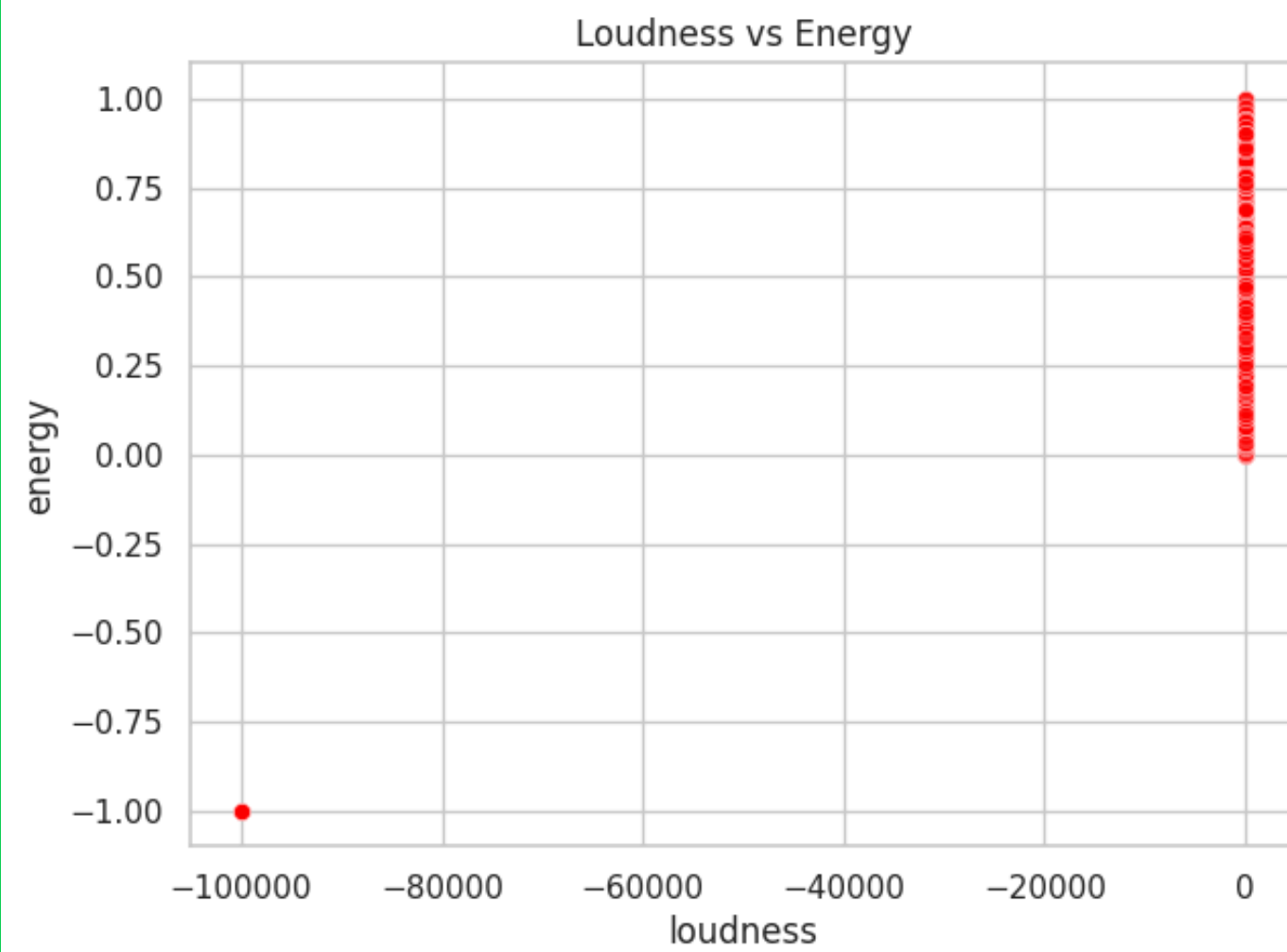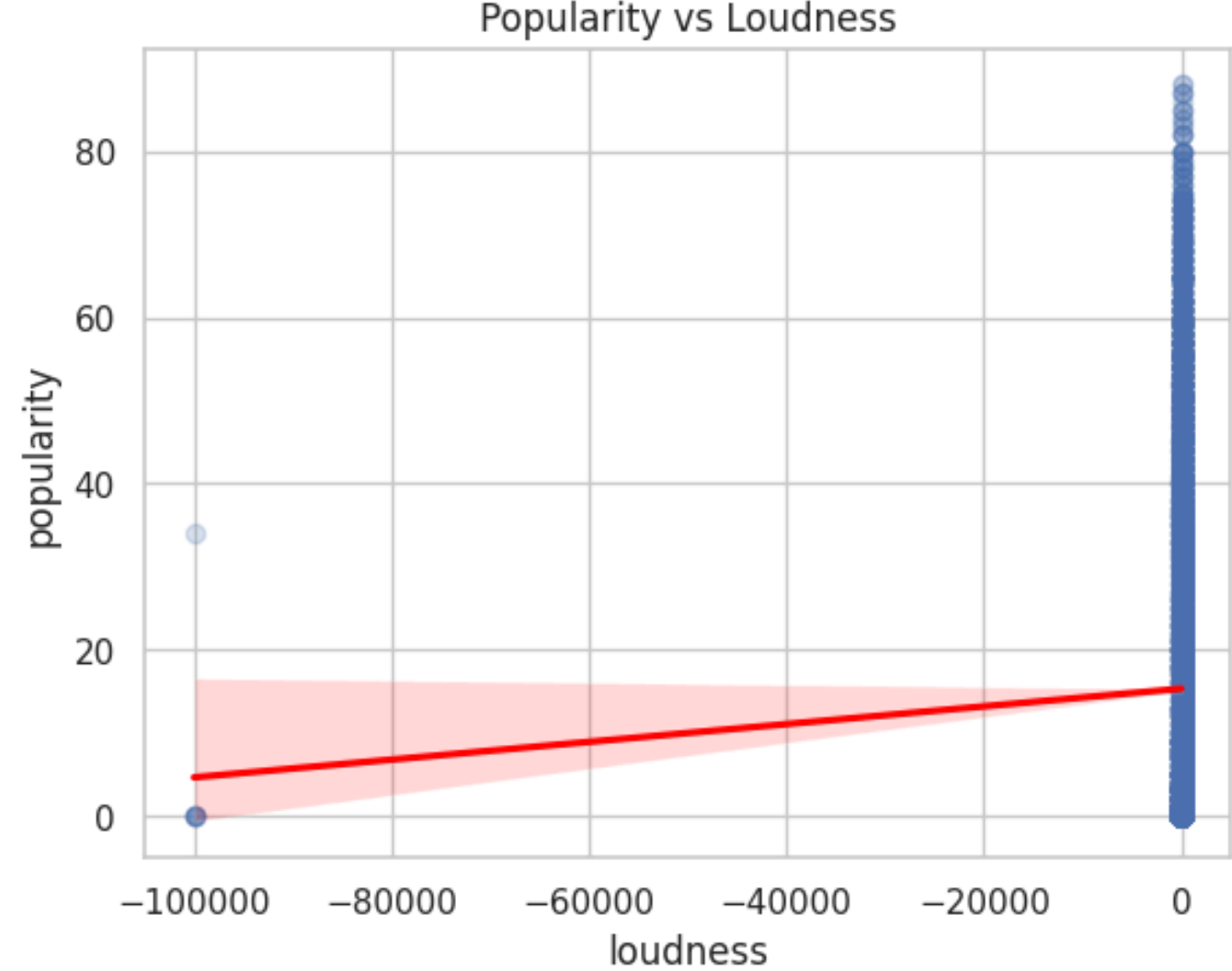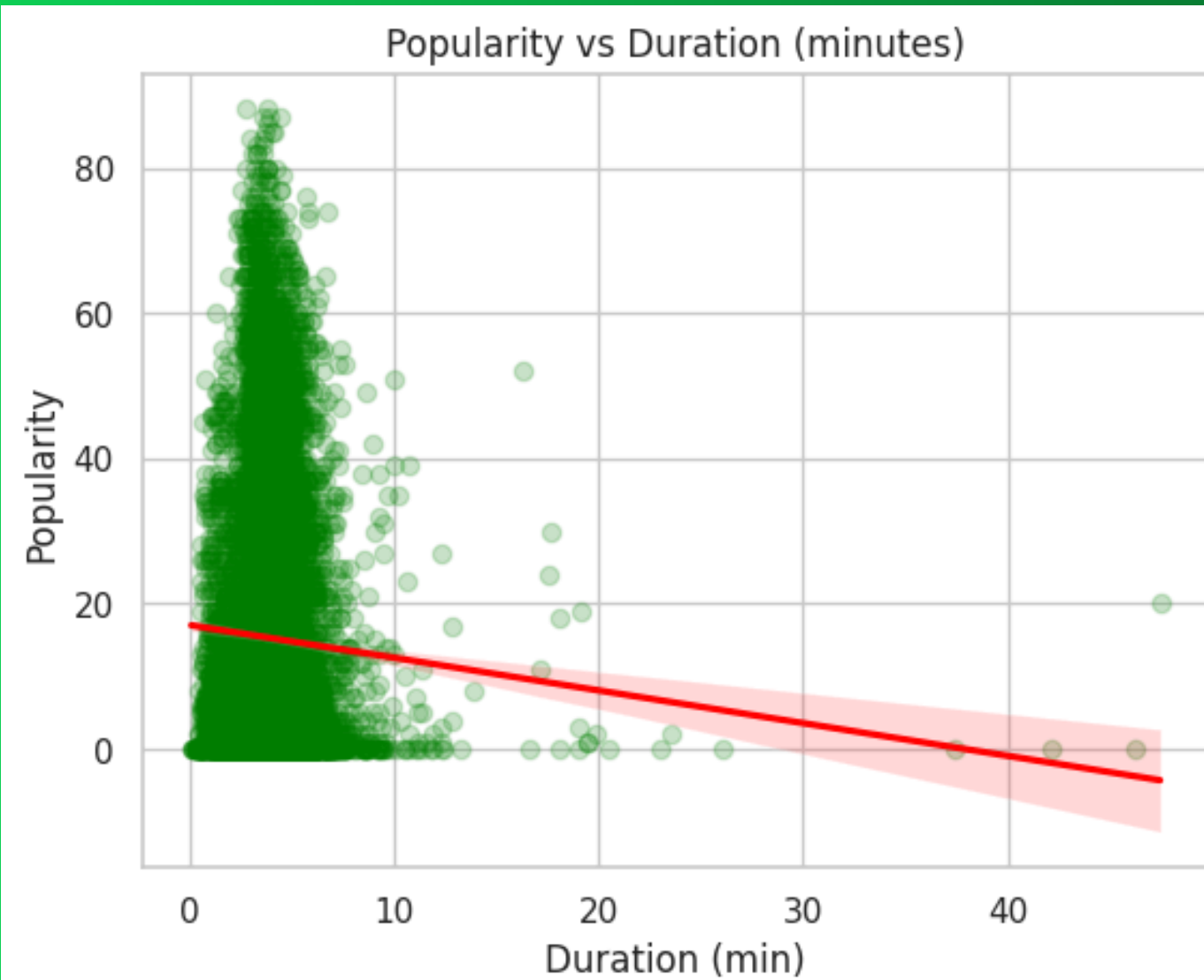
- **Popularity:** Most tracks cluster in the 40–70 range, indicating moderate popularity.
- **Duration:** The average song duration is approximately 3.5 minutes.
- **Audio Features:** A majority of the songs in the dataset are characterized by high danceability and energy.
- **Musical Structure:** The most common time signature is 4/4, and the majority of tracks are in a major key.

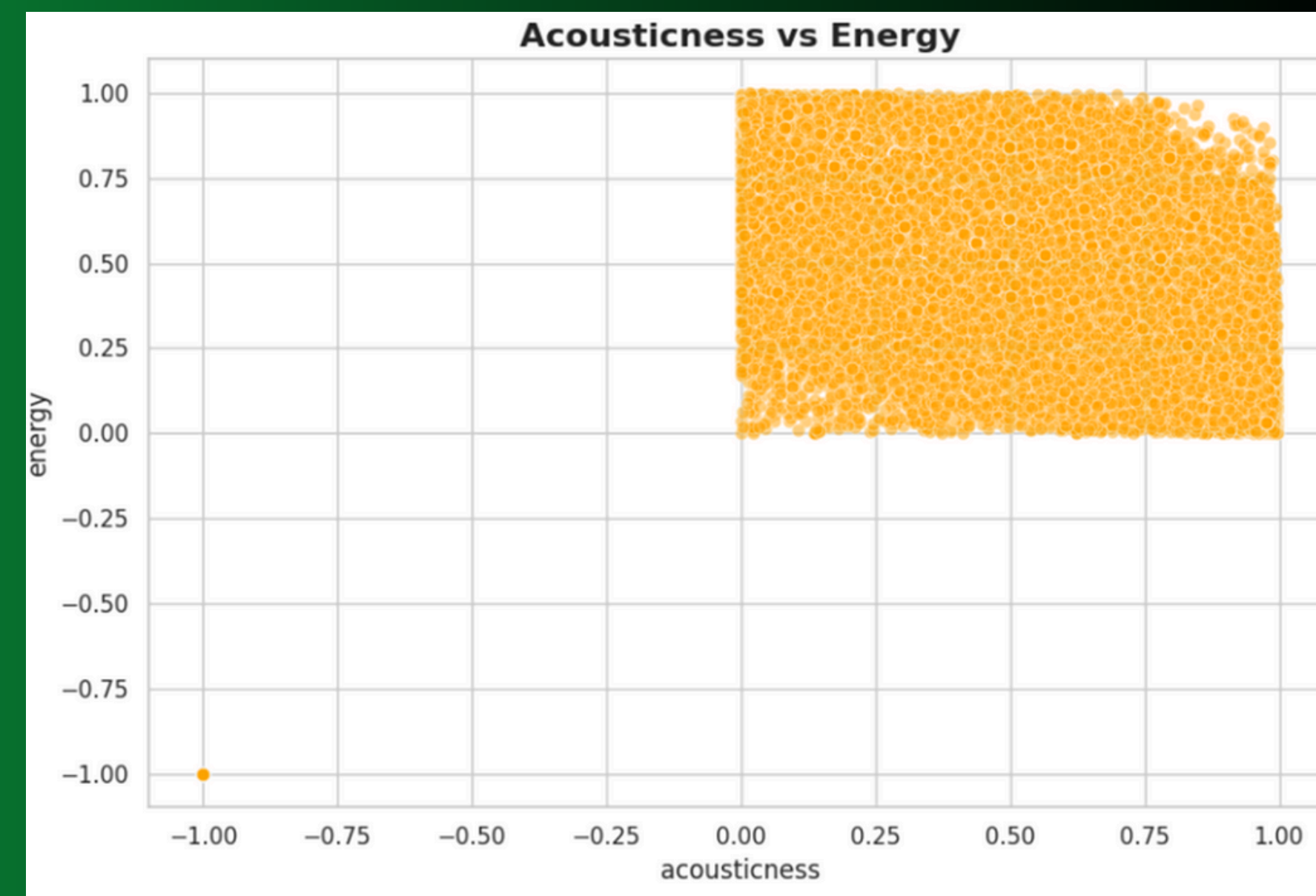# Bivariate Analysis(Numerical vs. Numerical



Both show a slight positive correlation. High danceability and energy help a song's popularity, but they are not the sole drivers. A negative trend is observed in Popularity vs. AcousticnessThe market favors produced, modern sounds..

## Popularity vs Duration (minutes)

## Popularity vs Loudness

## Loudness vs Energy

## Popularity vs Acousticness

Loudness vs Energy — Acousticness vs Energy

There is a very strong positive correlation.. High-energy tracks tend to be more positive/happy (high valence), while low-energy tracks are often associated with lower (sad or calm) valence.

# Bivariate Analysis(Numerical vs. Numerical )



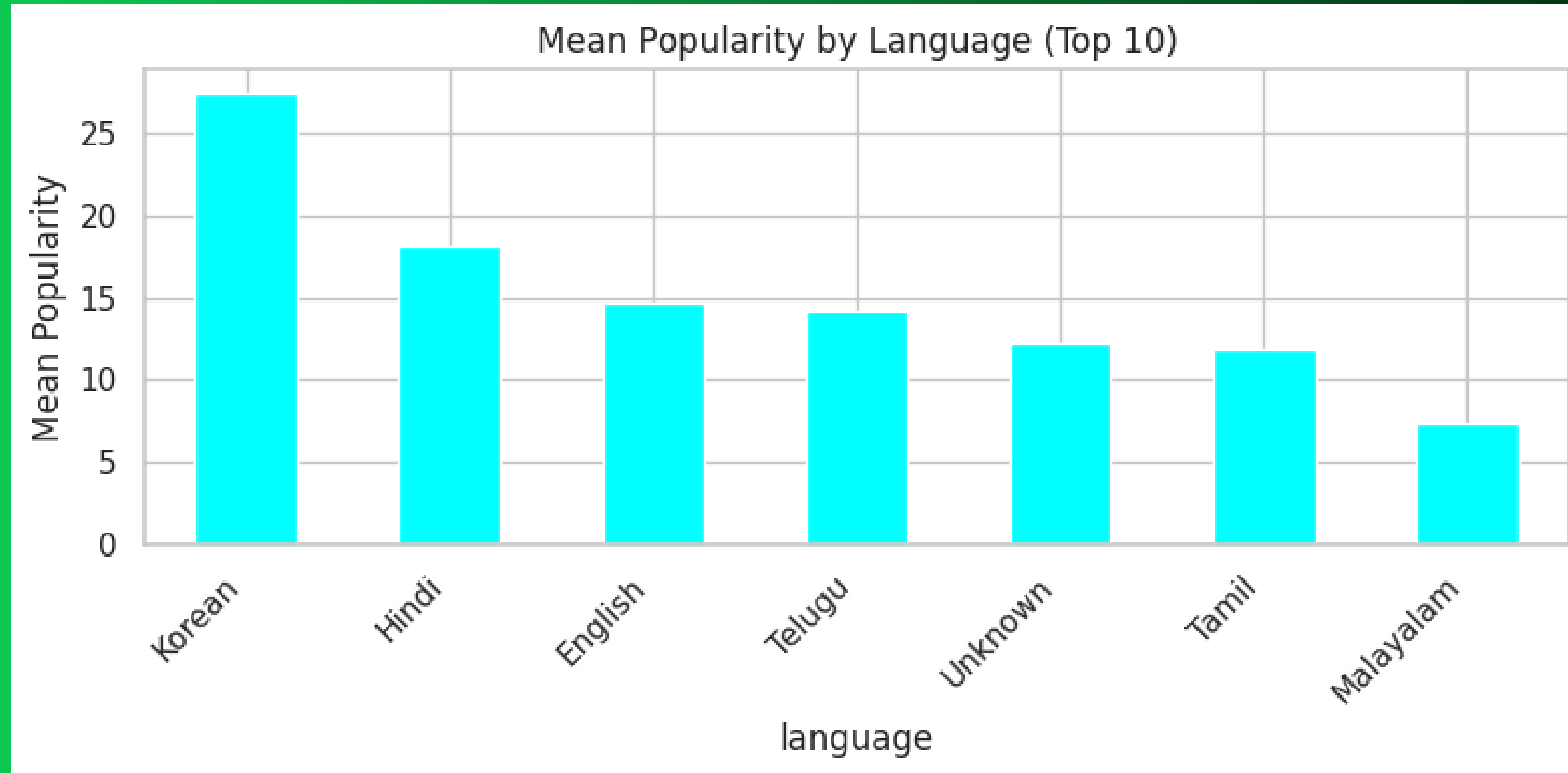Mean Popularity by Language (Top 10)

# Overview: The Correlation Heatmap



Correlation matrix

This map quickly shows all relationships. I highlights the strong links (like Energy-Loudness) and confirms the weak-but-important links to Popularity (e.g., slight positive from Danceability).

# Bivariate & Correlation Analysis

**Examining the relationships between features provided deeper insights:**

- **A strong positive correlation exists between danceability and energy.**
- **Popularity shows a slight positive correlation with both energy and danceability.**
- **Loudness and energy are highly correlated, which is an expected relationship.**
- **A negative trend was observed between acousticness and popularity, suggesting that acoustic songs tend to be less popular on average in this dataset.**

# Multivariate analysis



Top 6 danceability-energy-valence combos in top popularity quartile



PCA projection of acoustic/instrumental/speechiness clusters

# Multivariate analysis summary:

1.Feature Synergy (The Sonic Sweet Spot) Popularity is driven by a balanced blend of features (e.g., high Energy, mid-range Valence, high Loudness) in the Top 5%of tracks. Takeaway: Avoid extremes; synergy matters more than maxing out one single audio metric.

2.The Modern Mix (Efficiency & Impact) :  Time Trend 1: Average song duration is decreasing over the years. Time Trend 2: Average song loudness is increasing over the years.
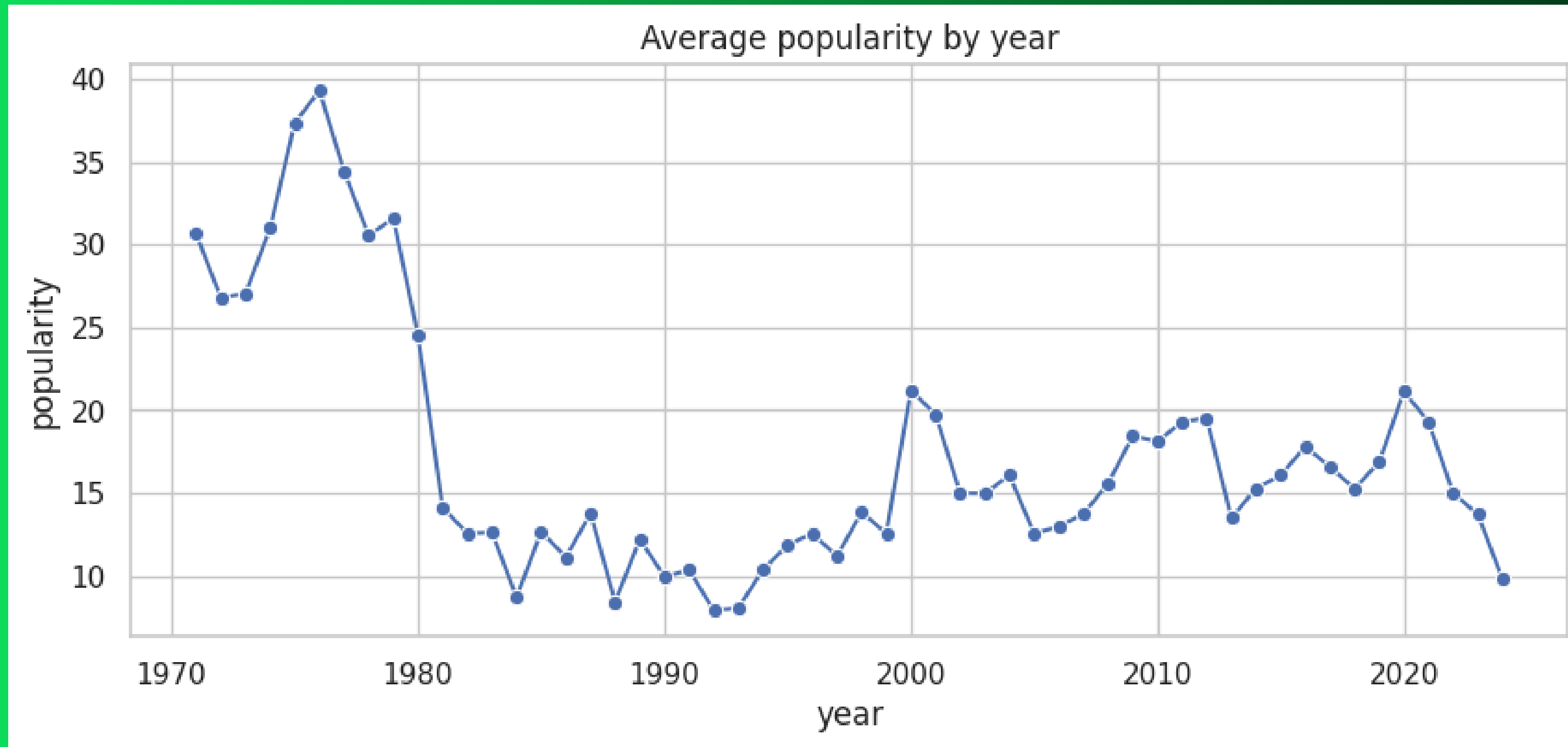
3. Takeaway: Modern hits must be short, concise, and competitively loud for streaming success.

Pillar 3: Global Reach (Market Expansion)  Time Trend 3: The popularity of Non-English music (especially Spanish and Korean tracks) is rapidly accelerating.
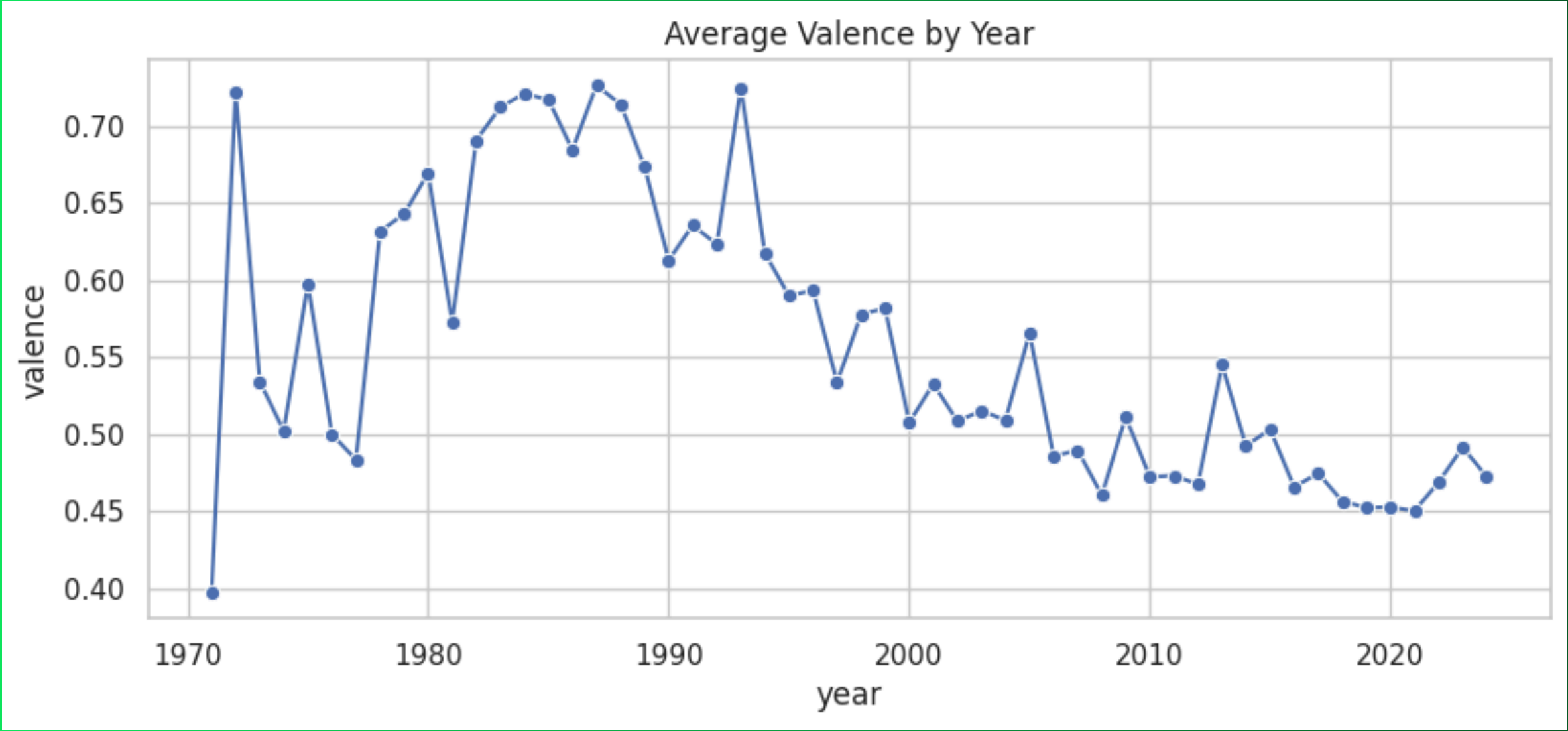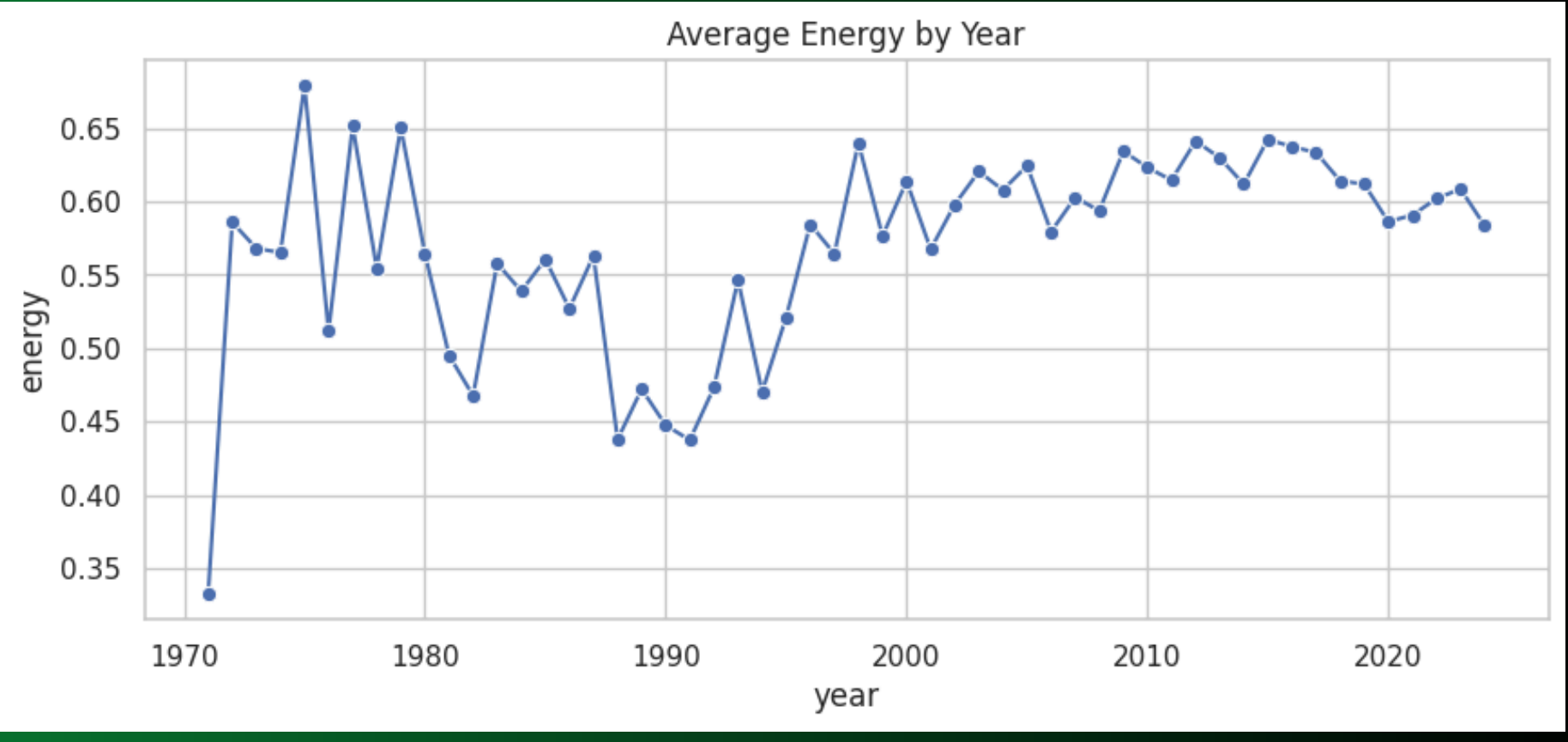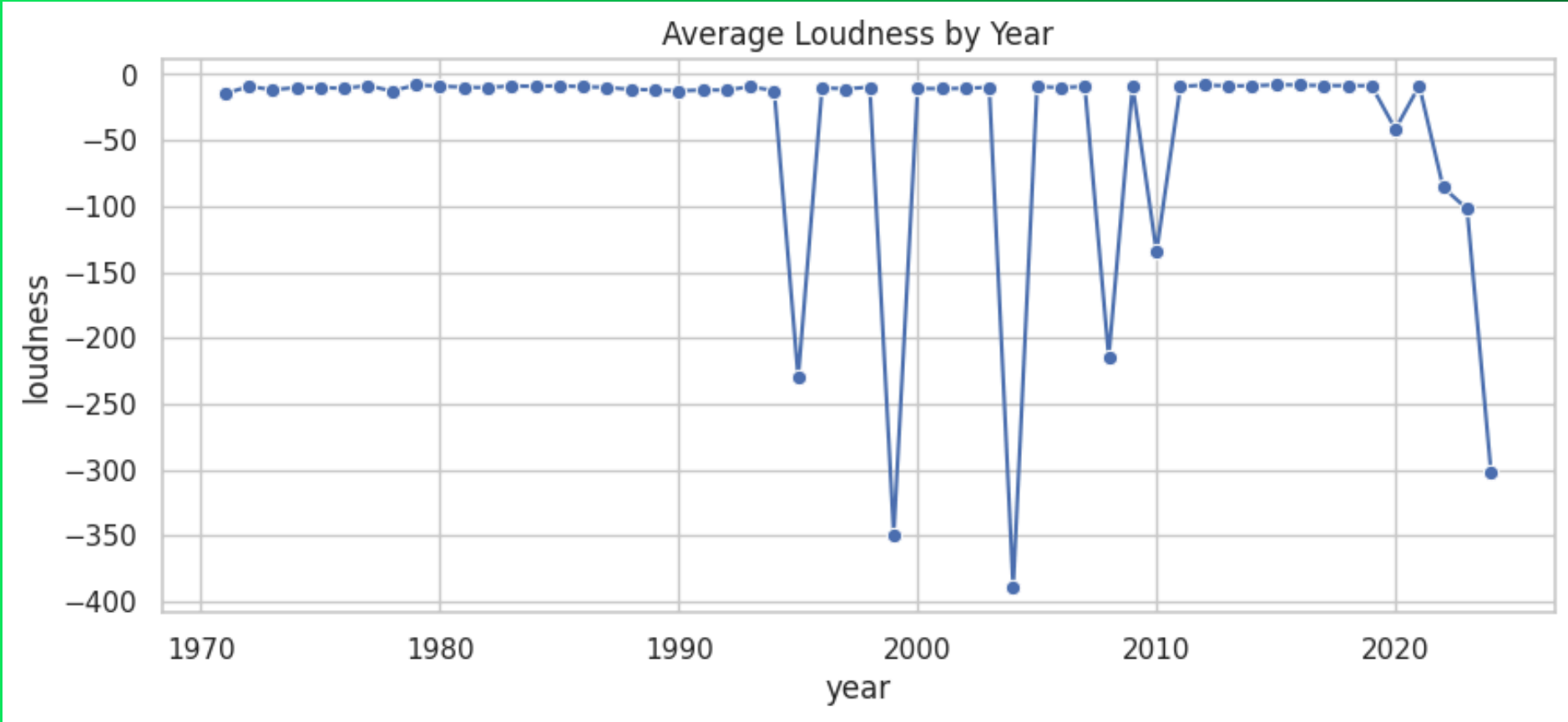
Takeaway: Capitalize on the global market shift by exploring multilingual content and expanding regional focus.

# Time Series Trends



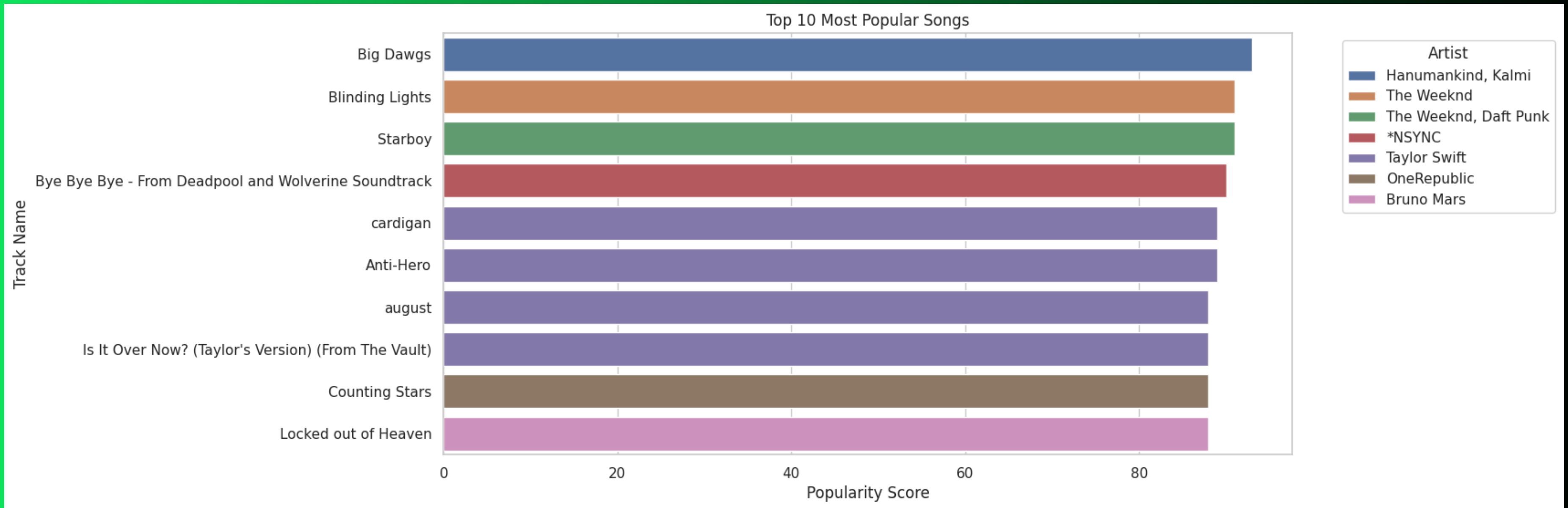Average popularity by year

Newer songs (streaming era) are consistently more popular.

Average Loudness by Year


Average Energy by Year


Average Valence by Year
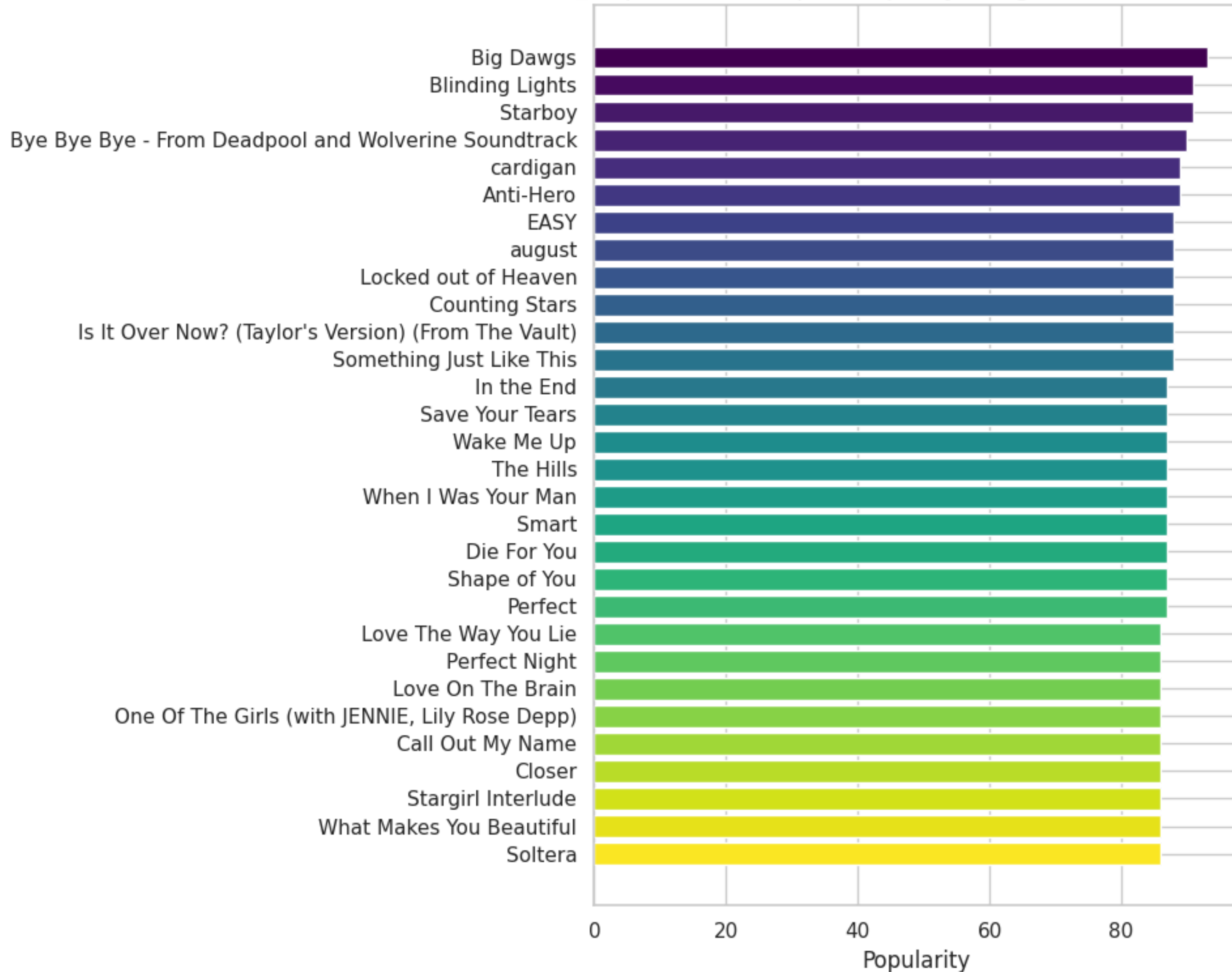
# summary of Time Series Trends

1. The Modern Mix is Shorter and Louder
   - Trend: The average song duration has decreased, while average loudness has increased consistently over the years.
   - Insight: Modern hit songs are designed for maximum impact and efficiency. The industry favors tracks that are loud and concise to hold listener attention on streaming platforms.
   - Visuals: Side-by-Side Line Plots of duration vs. year and loudness vs. year.
2. The Market is Going Global
   - Trend: The average popularity of tracks in languages like Spanish and Korean has seen a rapid, significant rise in recent years.
   - Insight: The music market is truly globalizing. Non-English music is no longer niche and represents a major growth area.
   - Visuals: Line Plot of popularity of Spanish/Korean tracks vs. year

Top 1% Most Popular Spotify Songs (Color Gradient)

# Outlier Analysis Result:

- Focus: The analysis centers on niche tracks defined by extreme scores on Speechiness or Instrumentalness.
- Core Principle (Pillar 4): Success for these tracks depends on Niche Clarity—they must fully commit to their segment.
- Speechiness Outliers (High Lyrical Niche):
- The market is strictly divided; there's little room for medium-level speech.
- Tracks must be either primarily sung or primarily lyrical/spoken.
- Instrumentalness Outliers (Pure Instrumental Niche):
- Instrumental tracks are rare in the popular dataset.
- To succeed, they must compensate by maximizing secondary features like Energy and Valence to create a strong, clear emotional impact.
- Takeaway: Outliers must specialize and optimize their remaining features instead of trying to hit the "average pop song" sweet spot.

# Thank you