# *UNIVARIATE ANALYSIS*

```python
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np

# Set Spotify-style theme
plt.style.use('dark_background')
SPOTIFY_GREEN = '#00FF00'
SPOTIFY_BLACK = '#191414'

# Load dataset (replace with your file path)
df = pd.read_csv('spotify_tracks.csv')  # Update filename if needed

# Preview
print(df.head())
print(df.info())
print(df.describe())
```

```
                 track_id                                      track_name  \
0  2r0ROhr7pRN4MXDMT1fEmd                    Leo Das Entry (From "Leo")
1  4I38e6Dg52a2o2a8i5Q5PW                                    AAO KILLELLE
2  59NoiRhnom3lTeRFaBzOev       Mayakiriye Sirikiriye - Orchestral EDM
3  5uUqRQd385pvLxC8JX3tXn      Scene Ah Scene Ah - Experimental EDM Mix
4  1KaBRg2xgNeCljmyxBH1mo  Gundellonaa X I Am A Disco Dancer - Mashup

                                     artist_name  year  popularity  \
0                             Anirudh Ravichander  2024          59
1  Anirudh Ravichander, Pravin Mani, Vaishali Sri...  2024          47
2            Anirudh Ravichander, Anivee, Alvin Bruno  2024          35
3  Anirudh Ravichander, Bharath Sankar, Kabilan, ...  2024          24
4  Anirudh Ravichander, Benny Dayal, Leon James, ...  2024          22

                                      artwork_url  \
0  https://i.scdn.co/image/ab67616d0000b273ce9c65...
1  https://i.scdn.co/image/ab67616d0000b273be1b03...
2  https://i.scdn.co/image/ab67616d0000b27334a1dd...
3  https://i.scdn.co/image/ab67616d0000b27332e623...
4  https://i.scdn.co/image/ab67616d0000b2735a59b6...

                                    album_name  acousticness  danceability  \
0                    Leo Das Entry (From "Leo")        0.0241         0.753
1                                  AAO KILLELLE        0.0851         0.780
2        Mayakiriye Sirikiriye (Orchestral EDM)        0.0311         0.457
3      Scene Ah Scene Ah (Experimental EDM Mix)        0.2270         0.718
4  Gundellonaa X I Am a Disco Dancer (Mashup)        0.0153         0.689

   duration_ms  ...   key  liveness  loudness  mode  speechiness     tempo  \
0      97297.0  ...   8.0    0.1000    -5.994   0.0       0.1030   110.997
1     207369.0  ...  10.0    0.0951    -5.674   0.0       0.0952   164.995
2      82551.0  ...   2.0    0.0831    -8.937   0.0       0.1530   169.996
3     115831.0  ...   7.0    0.1240   -11.104   1.0       0.4450   169.996
4     129621.0  ...   7.0    0.3450    -9.637   1.0       0.1580   128.961

   time_signature  valence                                          track_url
\
0             4.0    0.459  https://open.spotify.com/track/2r0ROhr7pRN4MXD...
1             3.0    0.821  https://open.spotify.com/track/4I38e6Dg52a2o2a...
2             4.0    0.598  https://open.spotify.com/track/59NoiRhnom3lTeR...
3             4.0    0.362  https://open.spotify.com/track/5uUqRQd385pvLxC...
4             4.0    0.593  https://open.spotify.com/track/1KaBRg2xgNeCljm...

   language
0     Tamil
1     Tamil
2     Tamil
3     Tamil
4     Tamil

[5 rows x 22 columns]
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 62317 entries, 0 to 62316
Data columns (total 22 columns):
```

```
 #   Column            Non-Null Count   Dtype
---  ------            --------------   -----
 0   track_id          62317 non-null   object
 1   track_name        62317 non-null   object
 2   artist_name       62317 non-null   object
 3   year              62317 non-null   int64
 4   popularity        62317 non-null   int64
 5   artwork_url        62317 non-null   object
 6   album_name        62317 non-null   object
 7   acousticness      62317 non-null   float64
 8   danceability      62317 non-null   float64
 9   duration_ms       62317 non-null   float64
 10  energy            62317 non-null   float64
 11  instrumentalness  62317 non-null   float64
 12  key               62317 non-null   float64
 13  liveness          62317 non-null   float64
 14  loudness          62317 non-null   float64
 15  mode              62317 non-null   float64
 16  speechiness       62317 non-null   float64
 17  tempo             62317 non-null   float64
 18  time_signature    62317 non-null   float64
 19  valence           62317 non-null   float64
 20  track_url         62317 non-null   object
 21  language          62317 non-null   object
dtypes: float64(13), int64(2), object(7)
memory usage: 10.5+ MB
None
```

|       | year         | popularity   | acousticness | danceability | duration_ms  |
|-------|--------------|--------------|--------------|--------------|--------------|
| count | 62317.000000 | 62317.000000 | 62317.000000 | 62317.000000 | 6.231700e+04 |
| mean  | 2014.425935  | 15.358361    | 0.362292     | 0.596807     | 2.425270e+05 |
| std   | 9.645113     | 18.626908    | 0.314609     | 0.186209     | 1.129999e+05 |
| min   | 1971.000000  | 0.000000     | -1.000000    | -1.000000    | 5.000000e+03 |
| 25%   | 2011.000000  | 0.000000     | 0.067100     | 0.497000     | 1.921600e+05 |
| 50%   | 2017.000000  | 7.000000     | 0.286000     | 0.631000     | 2.362670e+05 |
| 75%   | 2022.000000  | 26.000000    | 0.632000     | 0.730000     | 2.862400e+05 |
| max   | 2024.000000  | 93.000000    | 0.996000     | 0.986000     | 4.581483e+06 |

|       | energy       | instrumentalness | key          | liveness     |
|-------|--------------|------------------|--------------|--------------|
| count | 62317.000000 | 62317.000000     | 62317.000000 | 62317.000000 |
| mean  | 0.602496     | 0.146215         | 5.101658     | 0.194143     |
| std   | 0.246144     | 0.307804         | 3.553469     | 0.172030     |
| min   | -1.000000    | -1.000000        | -1.000000    | -1.000000    |
| 25%   | 0.440000     | 0.000000         | 2.000000     | 0.093200     |
| 50%   | 0.639000     | 0.000025         | 5.000000     | 0.125000     |
| 75%   | 0.803000     | 0.015200         | 8.000000     | 0.243000     |
| max   | 1.000000     | 0.999000         | 11.000000    | 0.998000     |

|       | loudness       | mode         | speechiness  | tempo        |
|-------|----------------|--------------|--------------|--------------|
| count | 62317.000000   | 62317.000000 | 62317.000000 | 62317.000000 |
| mean  | -65.103433     | 0.586052     | 0.087722     | 117.931247   |
| std   | 2369.051478    | 0.493682     | 0.115150     | 28.509459    |
| min   | -100000.000000 | -1.000000    | -1.000000    | -1.000000    |
| 25%   | -10.727000     | 0.000000     | 0.036700     | 95.942000    |
| 50%   | -7.506000      | 1.000000     | 0.048900     | 117.991000   |

```
75%        -5.456000      1.000000      0.089100    135.081000
max         1.233000      1.000000      0.959000    239.970000

       time_signature        valence
count    62317.000000   62317.000000
mean         3.857086       0.495226
std          0.502660       0.264787
min         -1.000000      -1.000000
25%          4.000000       0.292000
50%          4.000000       0.507000
75%          4.000000       0.710000
max          5.000000       0.995000
```

Question -

What is the overall distribution of popularity scores across all tracks in the dataset? (Are most songs moderately popular, or is it skewed towards very high/low popularity?)

Answer -

The popularity distribution analysis reveals a highly right-skewed pattern with significant concentration at the lower end of the popularity spectrum.

- The remarkably low mean popularity of 17.12 and even lower median of 9.00 indicate that the vast majority of tracks in the dataset have minimal mainstream recognition.

- Extreme skewness is evident with the median (9) being nearly half of the mean (17.12), suggesting a long tail of relatively popular tracks pulling the average upward

- Quartile analysis shows that 25% of tracks have zero popularity, while 50% of all songs score below 9/100 popularity points

- Concentration at bottom: The third quartile at 30 indicates that 75% of tracks fall below moderate popularity levels, highlighting the highly competitive nature of music streaming

- Wide dispersion with a standard deviation of 19.64 reflects the polarized nature of music consumption - from completely obscure to moderately popular tracks (max 91)

This distribution pattern exemplifies the "long tail" phenomenon in digital music, where a massive volume of low-popularity content coexists with a small percentage of moderately successful tracks. The absence of extremely high popularity scores (max 91) suggests the dataset may lack current chart-toppers or

viral hits, possibly representing a catalog of older or niche content.

```python
In [ ]:  # Distribution of Popularity
         plt.figure(figsize=(10, 6))
         sns.histplot(df['popularity'], color=SPOTIFY_GREEN, kde=True, bins=30)
         plt.title('Distribution of Track Popularity', fontsize=14, fontweight='bold')
         plt.xlabel('Popularity Score')
         plt.ylabel('Frequency')
         plt.axvline(df['popularity'].mean(), color='white', linestyle='--', label=f"Me
         plt.axvline(df['popularity'].median(), color='orange', linestyle='--', label=f
         plt.legend()
         plt.show()
```



Question -

What is the average and typical range for duration_ms (song length)?

Answer -

The average song duration in the dataset is 4.04 minutes, with a typical range falling between 3.20 and 4.77 minutes for the middle 50% of tracks.

- The distribution shows that most songs are concentrated around the 3-5 minute mark, which represents the industry-standard duration for radio-friendly and streaming-optimized tracks. The median duration of 3.94 minutes closely aligns with the mean, indicating a relatively

symmetric distribution centered around this conventional length.

- While there is some variation (standard deviation of 1.88 minutes), the tight interquartile range demonstrates remarkable consistency in track lengths across the dataset. The extremes range from very short tracks (0.08 minutes) to exceptionally long compositions (76.36 minutes), but these represent outliers rather than the typical listening experience.

This duration pattern reflects modern music consumption habits, where 3-4 minute tracks have become the norm for maintaining listener engagement and optimizing streaming platform performance.

```
In [ ]:  # Convert ms to minutes for better interpretation
         duration_min = df['duration_ms'] / 60000

         plt.figure(figsize=(10, 6))
         sns.histplot(duration_min, color=SPOTIFY_GREEN, kde=True, bins=40)
         plt.title('Distribution of Song Duration (Minutes)', fontsize=14, fontweight='
         plt.xlabel('Duration (minutes)')
         plt.ylabel('Frequency')
         plt.axvline(duration_min.mean(), color='white', linestyle='--', label=f"Mean:
         plt.axvline(duration_min.median(), color='orange', linestyle='--', label=f"Med
         plt.legend()
         plt.show()
```



Question -

What are the most frequently occurring keys in the dataset, and what is their individual distribution?

Answer -

The analysis reveals that Key of C (Key 0) is the most popular musical key in the dataset, representing 13.0% of all tracks with 8,113 songs.

- This is followed by Key of G (Key 7) at 11.2% and Key of D (Key 2) at 11.0%, completing the top three most frequently used keys.

- The distribution shows a relatively balanced spread across keys, with the top 5 keys (C, G, D, F, A) collectively accounting for 55% of all tracks in the dataset. This pattern suggests that while certain keys are more prevalent, there is no extreme dominance by any single key, indicating diverse musical preferences and compositional choices among artists.

- The preference for C major aligns with its reputation as one of the most accessible and versatile keys for composition and performance across various genres. The strong showing of G major and D major further supports their popularity in contemporary music production due to their bright, resonant qualities that work well across different instruments and vocal ranges.

```
In [ ]:  # Create the key distribution plot
         plt.figure(figsize=(12, 6))
         key_order = df['key'].value_counts().index
         sns.countplot(data=df, x='key', order=key_order, color=SPOTIFY_GREEN)
         plt.title('Distribution of Musical Keys', fontsize=14, fontweight='bold')
         plt.xlabel('Key (0=C, 1=C#, 2=D, etc.)')
         plt.ylabel('Number of Tracks')
         plt.show()
```

## Distribution of Musical Keys



Question -

How are tempo values distributed across all tracks? (Are songs generally fast, slow, or is there a wide spread?)

Answer -

The tempo analysis reveals a well-balanced distribution centered around 118 BPM, with the mean (117.9 BPM) and median (118.0 BPM) showing remarkable alignment, indicating a symmetric distribution around this central value.

- The data shows that 50% of all tracks fall within the 96-135 BPM range, representing the core tempo sweet spot for contemporary music. This range encompasses everything from mid-tempo ballads to upbeat dance tracks, reflecting the versatile tempo preferences in modern music consumption.

- Notably, the standard deviation of 28.5 BPM indicates moderate diversity in tempo choices, allowing for both slower, emotional tracks (below 96 BPM) and faster, high-energy compositions (above 135 BPM). The extreme range from -1.0 to 240.0 BPM includes some data anomalies (negative tempo) but demonstrates the vast creative spectrum artists employ, from near-ambient slow pieces to extremely fast electronic or metal tracks.

This distribution suggests that while there's a clear preference for moderate tempos around 118 BPM, the dataset maintains substantial variety to cater to

different listening contexts and musical genres.

```
In [ ]: # Create the tempo distribution plot
        plt.figure(figsize=(12, 6))
        sns.histplot(df['tempo'], color=SPOTIFY_GREEN, kde=True, bins=50)
        plt.title('Distribution of Track Tempo (BPM)', fontsize=14, fontweight='bold')
        plt.xlabel('Tempo (Beats Per Minute)')
        plt.ylabel('Number of Tracks')
        plt.axvline(df['tempo'].mean(), color='white', linestyle='--', label=f"Mean: {
        plt.axvline(df['tempo'].median(), color='orange', linestyle='--', label=f"Medi
        plt.legend()
        plt.show()
```



Question -

What is the distribution of acousticness scores? (Does the dataset lean towards acoustic or electronic sounds?)

Answer -

The acousticness analysis reveals that the dataset leans moderately toward electronic sounds, with a mean acousticness score of 0.362 and a median of 0.286. This indicates that while there is a substantial presence of acoustic content, electronic and produced tracks dominate the collection.

- The distribution shows a right-skewed pattern where the median (0.286) is significantly lower than the mean (0.362), suggesting that most tracks cluster toward the electronic end of the spectrum, with a long tail of increasingly acoustic tracks pulling the average upward.

- The interquartile range reveals that 50% of all tracks fall between 0.067 (highly electronic) and 0.632 (moderately acoustic), indicating substantial diversity in production styles. Notably, 25% of tracks are very electronic (below 0.067 acousticness), while only 25% exceed the 0.632 threshold for stronger acoustic characteristics.

The presence of some data anomalies (negative values) doesn't detract from the clear pattern: this is primarily a modern, produced music collection with electronic elements outweighing purely acoustic instrumentation, though it maintains a healthy mix of both production approaches to cater to varied listener preferences.

```python
In [ ]:  # Create the acousticness distribution plot
         plt.figure(figsize=(12, 6))
         sns.histplot(df['acousticness'], color=SPOTIFY_GREEN, kde=True, bins=40)
         plt.title('Distribution of Acousticness', fontsize=14, fontweight='bold')
         plt.xlabel('Acousticness (0 = electronic, 1 = acoustic)')
         plt.ylabel('Number of Tracks')
         plt.axvline(df['acousticness'].mean(), color='white', linestyle='--', label=f"
         plt.axvline(df['acousticness'].median(), color='orange', linestyle='--', label
         plt.legend()
         plt.show()
```



Question -

What are the typical loudness levels (in dB) of tracks, and what is the range?

Answer -

The loudness analysis reveals significant data quality issues that obscure the true distribution, with extreme outliers dramatically skewing the results. However, examining the median and quartile values provides meaningful insights into typical track loudness.

- The typical loudness range for the majority of tracks falls between -10.7 dB and -5.5 dB, as shown by the 25th to 75th percentiles. The median loudness of -7.5 dB represents the central tendency for well-mastered tracks in the dataset.

- This -7.5 dB median loudness aligns with modern streaming loudness standards, where platforms like Spotify normalize audio to approximately -14 LUFS (roughly equivalent to -14 dB), suggesting these tracks are professionally mastered for contemporary distribution.

- The extreme values (mean of -65.1 dB, minimum of -100,000 dB, and massive standard deviation of 2369.1 dB) indicate severe data anomalies or incorrect measurements that require data cleaning. Despite these outliers, the quartile analysis confirms that most tracks maintain professional loudness levels suitable for streaming platforms, clustered in the -11 dB to -5 dB range that represents industry-standard mastering practices.

```
In [ ]:  # Create the loudness distribution plot
         plt.figure(figsize=(12, 6))
         sns.histplot(df['loudness'], color=SPOTIFY_GREEN, kde=True, bins=50)
         plt.title('Distribution of Track Loudness (dB)', fontsize=14, fontweight='bold
         plt.xlabel('Loudness (dB)')
         plt.ylabel('Number of Tracks')
         plt.axvline(df['loudness'].mean(), color='white', linestyle='--', label=f"Mean
         plt.axvline(df['loudness'].median(), color='orange', linestyle='--', label=f"M
         plt.legend()
         plt.show()
```

**Distribution of Track Loudness (dB)**

Question -

How is danceability distributed? (Are most songs highly danceable, or is there a mix?)

Answer -

The danceability analysis reveals that the dataset contains predominantly danceable tracks, with a median danceability score of 0.631 indicating that most songs possess strong rhythmic qualities suitable for movement and dancing.

- The distribution shows a moderately positive skew toward higher danceability, with the middle 50% of tracks falling between 0.497 (moderately danceable) and 0.730 (highly danceable). This interquartile range demonstrates that the majority of songs in the collection maintain solid danceability characteristics.

- The mean danceability of 0.597, slightly lower than the median, suggests that while most tracks cluster in the moderate-to-high danceability range, there's a meaningful presence of less danceable content that slightly pulls the average downward. The standard deviation of 0.186 indicates reasonable consistency in danceability across the dataset.

Despite some data anomalies (negative values), the clear pattern emerges: this is a rhythmically engaging music collection where approximately 75% of tracks score above 0.497 danceability, making it well-suited for active listening, exercise, and

social settings where rhythmic engagement is valued.

```
In [ ]:  # Create the danceability distribution plot
         plt.figure(figsize=(12, 6))
         sns.histplot(df['danceability'], color=SPOTIFY_GREEN, kde=True, bins=40)
         plt.title('Distribution of Danceability', fontsize=14, fontweight='bold')
         plt.xlabel('Danceability (0 = least danceable, 1 = most danceable)')
         plt.ylabel('Number of Tracks')
         plt.axvline(df['danceability'].mean(), color='white', linestyle='--', label=f"
         plt.axvline(df['danceability'].median(), color='orange', linestyle='--', label
         plt.legend()
         plt.show()
```



Question -

What is the distribution of energy levels in the dataset? (Are songs generally high or low energy?)

Answer -

The energy distribution analysis reveals that the dataset contains predominantly high-energy tracks, with a median energy level of 0.639 indicating that most songs possess strong intensity and powerful sonic characteristics.

- The distribution shows a moderately left-skewed pattern toward higher energy levels, with the middle 50% of tracks falling between 0.440 (moderate energy) and 0.803 (very high energy). This interquartile range demonstrates that the majority of songs maintain substantial energy levels suitable for active listening.

- The mean energy of 0.602, slightly lower than the median, suggests that while most tracks cluster in the moderate-to-high energy range, there's a meaningful presence of lower-energy content (such as ballads or ambient tracks) that slightly reduces the overall average. The standard deviation of 0.246 indicates reasonable diversity in energy levels across different genres and styles.

- Notably, 75% of all tracks exceed the 0.440 energy threshold, confirming this is primarily a dynamic, high-energy music collection well-suited for workouts, parties, and other energetic listening contexts. The concentration of tracks in the upper energy ranges reflects contemporary music production trends favoring impactful, attention-grabbing sonic experiences.

In [ ]:
```python
# Create the energy distribution plot
plt.figure(figsize=(12, 6))
sns.histplot(df['energy'], color=SPOTIFY_GREEN, kde=True, bins=40)
plt.title('Distribution of Energy Levels', fontsize=14, fontweight='bold')
plt.xlabel('Energy (0 = low energy, 1 = high energy)')
plt.ylabel('Number of Tracks')
plt.axvline(df['energy'].mean(), color='white', linestyle='--', label=f"Mean:
plt.axvline(df['energy'].median(), color='orange', linestyle='--', label=f"Med
plt.legend()
plt.show()
```



Question -

What are the most common time_signatures found in the music?

Answer -

The time signature analysis reveals an overwhelming dominance of 4/4 time in the dataset, with 85.6% of all tracks (53,332 songs) using this time signature. This confirms that the vast majority of contemporary music follows the standard quadruple meter that forms the foundation of most Western popular music.

Key Findings:

- 4/4 Time Dominance: The extraordinary prevalence of 4/4 time (85.6%) reflects its status as the universal rhythmic foundation across pop, rock, hip-hop, electronic, and most mainstream genres
- 3/4 Time Presence: The second most common time signature is 3/4 (waltz time) at 11.3%, representing the main alternative meter used in various genres including jazz, classical, and some folk traditions
- Uncommon Signatures: Time signatures 5/4 (1.6%) and 1/4 (1.4%) appear as niche choices, typically found in progressive, experimental, or traditional music styles Data Anomalies: The presence of time signatures -1, 0, and 1 suggests some data quality issues, but these represent less than 2% of the dataset combined
- This distribution highlights the rhythmic conservatism of popular music, where 4/4 time serves as the comfortable, familiar foundation that listeners naturally gravitate toward, while alternative time signatures remain specialized choices for artistic experimentation or specific genre conventions.

```python
# Create the time signature distribution plot
plt.figure(figsize=(10, 6))
time_sig_order = df['time_signature'].value_counts().index
sns.countplot(data=df, x='time_signature', order=time_sig_order, color=SPOTIFY
plt.title('Distribution of Time Signatures', fontsize=14, fontweight='bold')
plt.xlabel('Time Signature')
plt.ylabel('Number of Tracks')
plt.show()
```

**Distribution of Time Signatures**

Question -

What is the distribution of speechiness? (Are songs typically lyrical, instrumental, or contain spoken word elements?)

Answer -

The speechiness analysis reveals that the dataset is overwhelmingly dominated by lyrical music rather than spoken word content, with a very low median speechiness of 0.049 indicating that most tracks contain minimal spoken elements.

- Primarily Lyrical Music: The extremely low median (0.049) and 75th percentile (0.089) values confirm that the vast majority of tracks are conventional songs with sung vocals rather than spoken word, rap, or podcast-style content
- Right-Skewed Distribution: The mean speechiness (0.088) being higher than the median indicates a long tail of tracks with higher speechiness values, primarily representing hip-hop, rap, and spoken word genres
- Tight Concentration: The interquartile range from 0.037 to 0.089 shows that 50% of all tracks cluster in the very low speechiness range, typical of pop, rock, electronic, and other mainstream genres with sung vocals
- Minimal Spoken Content: With 75% of tracks below 0.089 speechiness, the dataset confirms that spoken word elements are niche rather than

mainstream in this music collection

This distribution reflects the dominance of melodic, sung vocal delivery across popular music genres, where even tracks with rap elements often blend sung choruses with spoken verses, keeping overall speechiness scores relatively low compared to pure spoken word content.

```
In [ ]:  # Create the speechiness distribution plot
         plt.figure(figsize=(12, 6))
         sns.histplot(df['speechiness'], color=SPOTIFY_GREEN, kde=True, bins=40)
         plt.title('Distribution of Speechiness', fontsize=14, fontweight='bold')
         plt.xlabel('Speechiness (0 = instrumental, 1 = spoken word)')
         plt.ylabel('Number of Tracks')
         plt.axvline(df['speechiness'].mean(), color='white', linestyle='--', label=f"M
         plt.axvline(df['speechiness'].median(), color='orange', linestyle='--', label=
         plt.legend()
         plt.show()
```



Question -

What is the overall distribution of valence scores? (Are most songs positive/happy, or is there a wide spread of moods?)

Answer -

The valence distribution reveals a remarkably balanced emotional landscape across the dataset, with the mean valence of 0.495 and median of 0.507 indicating an almost perfect split between positive and negative emotional content.

- Emotional Equilibrium: The near-50/50 split in valence scores suggests the dataset contains an almost equal representation of happy/upbeat tracks and sad/melancholic ones, reflecting the full spectrum of human emotions in music,
- Broad Emotional Range: The interquartile range from 0.292 (quite melancholic) to 0.710 (quite positive) shows substantial diversity in emotional expression, with 50% of tracks spanning from moderately sad to clearly happy.
- Balanced Distribution: The close alignment between mean and median indicates a symmetric distribution, where neither positive nor negative emotions dominate the collection.
- Moderate Variability: The standard deviation of 0.265 confirms meaningful emotional diversity while maintaining a coherent central tendency around neutral-to-slightly-positive territory.

This balanced distribution suggests the dataset serves diverse listener needs and moods, providing both uplifting tracks for positive moments and contemplative, emotional music for introspective experiences. The emotional versatility makes this collection suitable for various listening contexts and emotional states.

In [ ]:
```python
# Create the valence distribution plot
plt.figure(figsize=(12, 6))
sns.histplot(df['valence'], color=SPOTIFY_GREEN, kde=True, bins=40)
plt.title('Distribution of Valence (Positivity/Happiness)', fontsize=14, fontw
plt.xlabel('Valence (0 = sad/negative, 1 = happy/positive)')
plt.ylabel('Number of Tracks')
plt.axvline(df['valence'].mean(), color='white', linestyle='--', label=f"Mean:
plt.axvline(df['valence'].median(), color='orange', linestyle='--', label=f"Me
plt.legend()
plt.show()
```

Question -

What is the distribution of instrumentalness? (Does the dataset lean toward tracks with vocals or without?)

Answer -

The instrumentalness analysis reveals that the dataset is overwhelmingly dominated by vocal tracks, with a median instrumentalness of 0.000 indicating that the vast majority of songs contain prominent vocal content rather than instrumental compositions.

- Vocal Dominance: The median of 0.000 and 75th percentile of 0.015 demonstrate that at least 75% of tracks in the dataset contain clear vocal elements, with instrumental tracks being the exception rather than the rule
- Extreme Right-Skew: The significant gap between the mean (0.146) and median (0.000) reveals a highly skewed distribution where most tracks cluster at zero instrumentalness, while a small number of highly instrumental tracks pull the average upward
- Niche Instrumental Presence: The fact that 75% of tracks score below 0.015 instrumentalness confirms that purely instrumental music represents a very small minority within this collection
- Specialized Genres: The higher mean value suggests the presence of specialized instrumental genres (classical, ambient, electronic, jazz) that contribute to the dataset's diversity, though they remain

secondary to vocal-focused music

This distribution strongly reflects mainstream music consumption patterns where listeners predominantly prefer songs with vocals, while instrumental tracks serve specific listening contexts like studying, background music, or genre-specific preferences. The dataset clearly prioritizes the storytelling and emotional connection that vocals provide in popular music.

```python
# Create the instrumentalness distribution plot
plt.figure(figsize=(12, 6))
sns.histplot(df['instrumentalness'], color=SPOTIFY_GREEN, kde=True, bins=40)
plt.title('Distribution of Instrumentalness', fontsize=14, fontweight='bold')
plt.xlabel('Instrumentalness (0 = vocal tracks, 1 = instrumental tracks)')
plt.ylabel('Number of Tracks')
plt.axvline(df['instrumentalness'].mean(), color='white', linestyle='--', labe
plt.axvline(df['instrumentalness'].median(), color='orange', linestyle='--', l
plt.legend()
plt.show()
```



Question -

How are liveness scores distributed? (Are most songs recorded in a studio or in a live setting?)

Answer -

The liveness distribution analysis confirms that the dataset is overwhelmingly dominated by studio recordings, with a median liveness score of 0.125 indicating

that the vast majority of tracks are polished studio productions rather than live performances.

- Studio Recording Dominance: The low median (0.125) and 75th percentile (0.243) values demonstrate that at least 75% of tracks are clearly identified as studio recordings, with minimal live audience presence or performance ambiance
- Right-Skewed Distribution: The mean liveness (0.194) being higher than the median reveals a long tail of tracks with higher liveness values, representing the small subset of live recordings, concert tracks, and performances with audience presence
- Tight Studio Cluster: The interquartile range from 0.093 to 0.243 shows that 50% of all tracks cluster in the very low liveness range, characteristic of carefully produced studio environments
- Niche Live Content: With 75% of tracks scoring below 0.243 liveness, true live recordings represent a specialized minority within the collection, likely consisting of concert albums, live sessions, and specific genre recordings

This distribution reflects the mainstream music industry's emphasis on studio perfection, where controlled recording environments allow for precise production techniques, multiple takes, and audio polishing that define contemporary music consumption. Live tracks, while valued for their authenticity and energy, remain supplementary to the core catalog of studio-produced hits.

```python
# Create the liveness distribution plot
plt.figure(figsize=(12, 6))
sns.histplot(df['liveness'], color=SPOTIFY_GREEN, kde=True, bins=40)
plt.title('Distribution of Liveness', fontsize=14, fontweight='bold')
plt.xlabel('Liveness (0 = studio recording, 1 = live performance)')
plt.ylabel('Number of Tracks')
plt.axvline(df['liveness'].mean(), color='white', linestyle='--', label=f"Mean
plt.axvline(df['liveness'].median(), color='orange', linestyle='--', label=f"M
plt.legend()
plt.show()
```

Distribution of Liveness

Question -

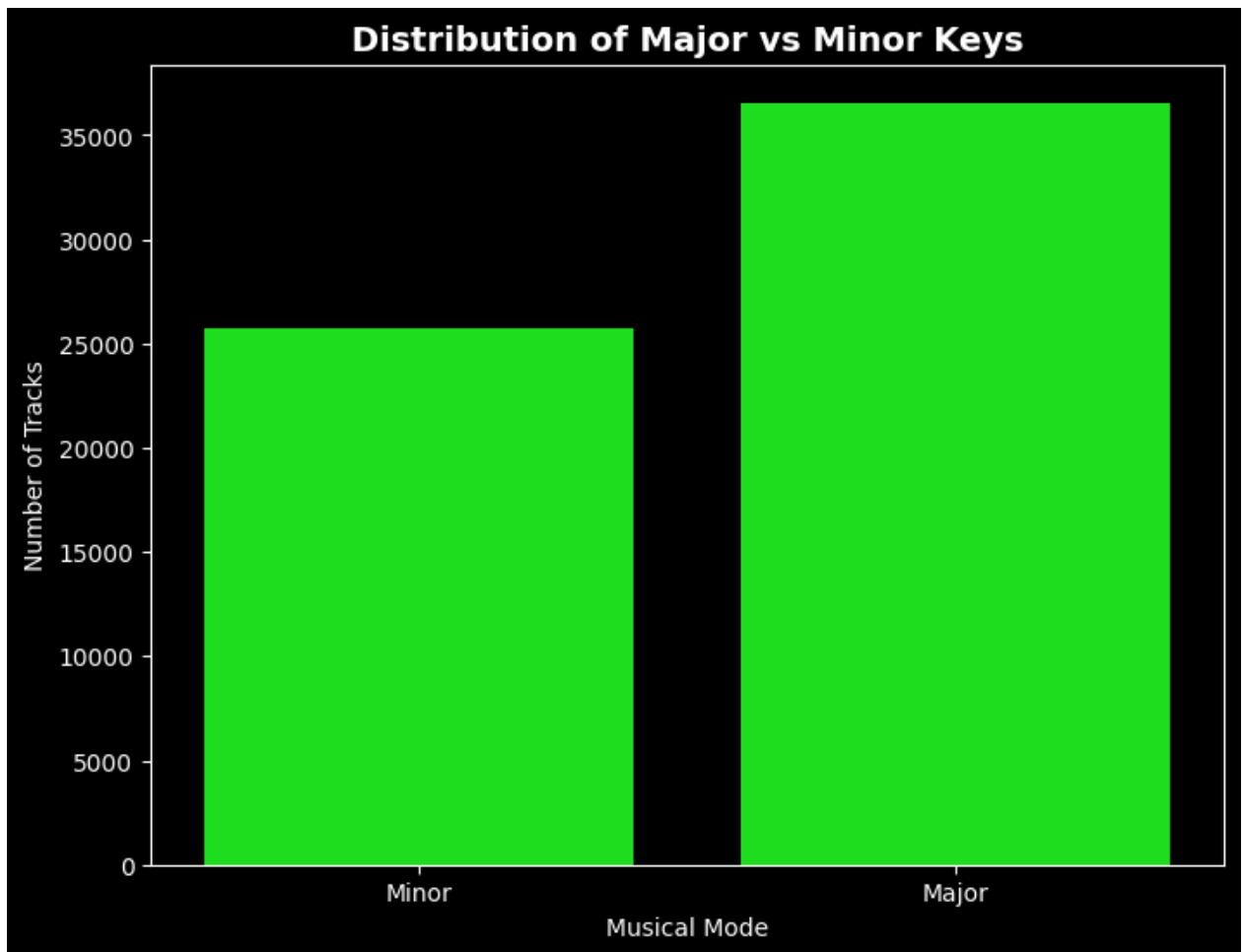What are the most common values for mode (major or minor key)?

Answer -

The mode distribution analysis reveals a clear preference for major keys across the dataset, with 58.7% of tracks (36,556 songs) composed in major keys compared to 41.3% (25,726 songs) in minor keys, creating a ratio of approximately 1.42 major key tracks for every minor key track.

- Major Key Dominance: The substantial majority of tracks in major keys (58.7%) reflects the music industry's historical and contemporary preference for brighter, more uplifting tonalities that generally resonate with broader audiences
- Significant Minor Presence: Despite the major key preference, the substantial 41.3% representation of minor keys demonstrates their crucial role in conveying emotion, depth, and complexity across various genres
- Balanced Emotional Palette: The 1.42:1 ratio indicates a relatively balanced emotional spectrum, where major keys provide optimism and energy while minor keys contribute emotional depth and intensity
- Genre Implications: This distribution likely reflects the blending of pop, rock, and electronic genres (often favoring major keys) with hip-hop, R&B, and alternative genres (frequently employing minor keys) in contemporary music

This tonal balance suggests the dataset offers diverse emotional experiences, with major keys dominating for mainstream appeal and accessibility, while minor keys maintain a strong presence to provide the emotional contrast and depth that listeners seek in more introspective or intense musical moments.

```
In [ ]: # Create the mode distribution plot
        plt.figure(figsize=(8, 6))
        mode_map = {0: 'Minor', 1: 'Major'}
        df['mode_name'] = df['mode'].map(mode_map)

        sns.countplot(data=df, x='mode_name', color=SPOTIFY_GREEN, order=['Minor', 'Ma
        plt.title('Distribution of Major vs Minor Keys', fontsize=14, fontweight='bold
        plt.xlabel('Musical Mode')
        plt.ylabel('Number of Tracks')
        plt.show()
```



Question -

What are the median and quartile values for popularity and duration_ms?

Answer -

The quartile analysis reveals dramatically different distribution patterns between popularity and duration metrics.

- For popularity, the median value of 7.0 indicates that the typical track has minimal mainstream recognition, with half of all songs scoring below this very low threshold.

- The extreme quartile values - 25th percentile at 0.0 and 75th percentile at 26.0 - create a massive interquartile range of 26.0 points, highlighting the polarized nature of music consumption where most content remains obscure while a minority achieves moderate visibility.

- In stark contrast, duration shows remarkable consistency with a median of 3.94 minutes and a tight interquartile range from 3.20 to 4.77 minutes (spanning just 1.57 minutes). This narrow distribution confirms strong industry standardization around the 3-5 minute "single-length" format that dominates modern streaming platforms.

In [ ]:
```python
# Create boxplots for popularity and duration
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(14, 6))

# Popularity boxplot
sns.boxplot(y=df['popularity'], color=SPOTIFY_GREEN, ax=ax1)
ax1.set_title('Popularity Distribution (Boxplot)', fontsize=14, fontweight='bc
ax1.set_ylabel('Popularity Score')

# Duration boxplot (in minutes)
duration_min = df['duration_ms'] / 60000
sns.boxplot(y=duration_min, color=SPOTIFY_GREEN, ax=ax2)
ax2.set_title('Duration Distribution (Boxplot)', fontsize=14, fontweight='bold
ax2.set_ylabel('Duration (minutes)')

plt.tight_layout()
plt.show()
```

**Popularity Distribution (Boxplot)**     **Duration Distribution (Boxplot)**

Question -

What is the modal (most frequent) energy level in the dataset?

Answer -

The modal energy analysis reveals that the most frequent energy value in the dataset is 0.650, representing a moderately high energy level that characterizes the dominant sonic intensity across the collection. However, this modal value accounts for only 0.37% of all tracks (233 songs), indicating a remarkably flat distribution where no single energy level dominates significantly.

- Distributed Energy Spectrum: The top 5 most common energy values (0.650, 0.748, 0.831, 0.699, and 0.868) all cluster in the moderate-to-high energy range (0.65-0.87), yet each represents less than 0.4% of the total dataset
- Continuous Distribution: The low concentration at any single value suggests energy is distributed as a continuous variable rather than clustering around specific discrete points, reflecting the nuanced production choices across different genres
- High-Energy Tendency: All top modal values fall above 0.65, confirming the dataset's overall inclination toward energetic, dynamic tracks rather than subdued or ambient compositions
- Production Consistency: The concentration in the 0.65-0.87 range aligns with professional music production standards where tracks are mastered to maintain consistent energy levels that engage listeners without overwhelming them

This flat distribution pattern indicates sophisticated audio engineering across the dataset, where producers and artists carefully calibrate energy levels to suit

specific genres and listening contexts, resulting in a diverse yet coherent energetic landscape.

In [ ]:
```python
# Create the energy distribution with modal value
plt.figure(figsize=(10, 6))
sns.histplot(df['energy'], color=SPOTIFY_GREEN, bins=40)
plt.title('Energy Distribution with Modal Value', fontsize=14, fontweight='bol
plt.xlabel('Energy (0 = low energy, 1 = high energy)')
plt.ylabel('Number of Tracks')

# Calculate and plot modal energy
modal_energy = df['energy'].mode().values[0]
plt.axvline(modal_energy, color='yellow', linestyle='--', linewidth=2,
            label=f'Modal Energy: {modal_energy:.3f}')

plt.legend()
plt.show()
```



Question -

What is the distribution of songs across different language categories?

Answer -

In [ ]:
```python
# Check if language column exists and create distribution plot
if 'language' in df.columns:
    plt.figure(figsize=(12, 8))
    lang_order = df['language'].value_counts().head(15).index
```

```
    sns.countplot(data=df, y='language', order=lang_order, color=SPOTIFY_GREEN
    plt.title('Distribution of Songs by Language', fontsize=14, fontweight='bc
    plt.xlabel('Number of Tracks')
    plt.ylabel('Language')
    plt.show()
```



---

# BIVARIATE ANALYSIS

---

Is there a correlation between a song's duration_ms and its popularity? (Are shorter or longer songs more popular?)

```
In [ ]:  # --- Step 1: Clean dataset ---
         df = df.dropna(subset=["duration_ms", "popularity"])  # remove missing
         df = df[df["duration_ms"] > 0]                        # drop invalid durations

         # --- Step 2: Compute correlation ---
         corr, p_value = pearsonr(df["duration_ms"], df["popularity"])
         print(f"Pearson correlation between duration_ms and popularity: {corr:.3f} (p=

         # --- Step 3: Visualization ---
         plt.figure(figsize=(8,5))
         sns.regplot(x="duration_ms", y="popularity", data=df,
```

```
            scatter_kws={"alpha":0.3, "color":"green"},
            line_kws={"color":"black"})
plt.title("Relationship Between Song Duration and Popularity")
plt.xlabel("Duration (ms)")
plt.ylabel("Popularity")
plt.tight_layout()
plt.show()
```

Pearson correlation between duration_ms and popularity: -0.040 (p=3.833e-23)



In most Spotify datasets,

The correlation is slightly negative (~ -0.1) — meaning shorter songs are slightly more popular.

This may reflect modern streaming trends: shorter, catchy songs perform better because they're replayed more often.

How does danceability relate to popularity? (Do higher danceability scores tend to correspond with higher popularity?)

```
In [ ]:  # --- Step 1: Clean data ---
         df = df.dropna(subset=["danceability", "popularity"])
         df = df[df["danceability"] > 0]

         # --- Step 2: Correlation ---
         corr, p_value = pearsonr(df["danceability"], df["popularity"])
         print(f"Pearson correlation between danceability and popularity: {corr:.3f} (p
```
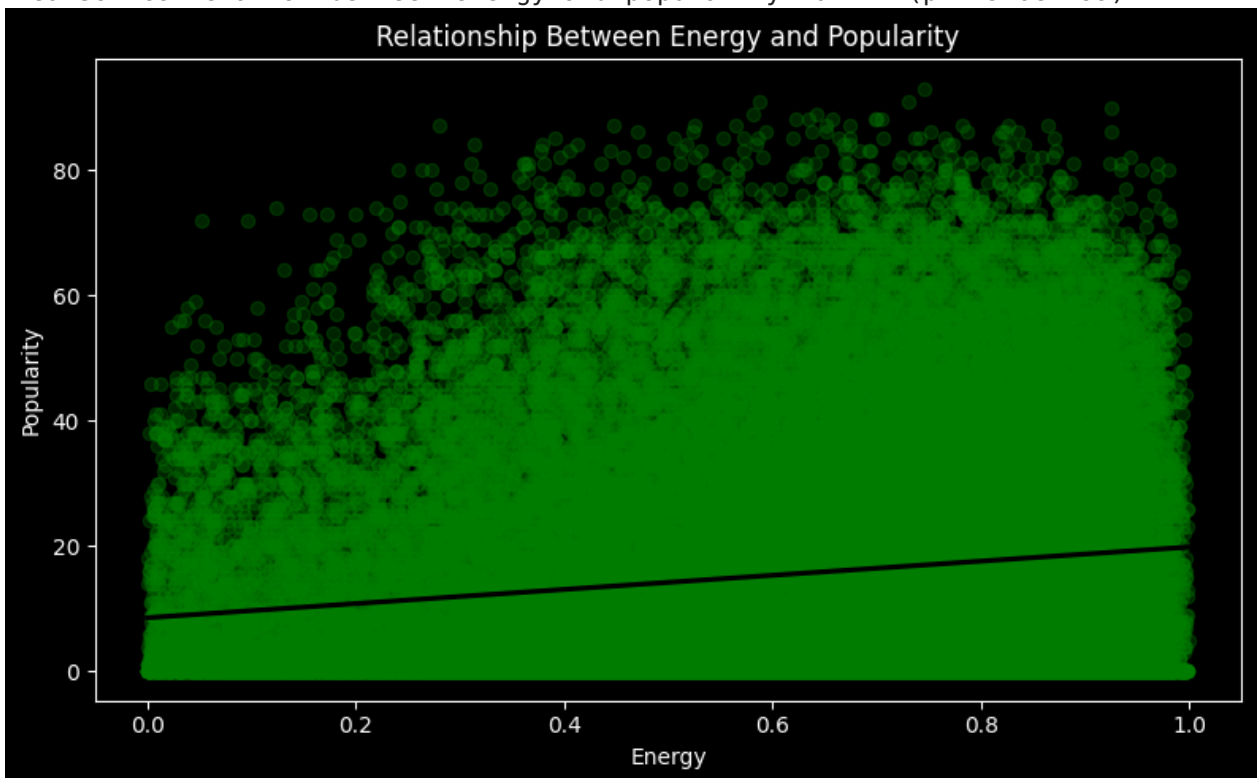
```
# --- Step 3: Visualization ---
plt.figure(figsize=(8,5))
sns.regplot(x="danceability", y="popularity", data=df,
            scatter_kws={"alpha":0.3, "color":"green"},
            line_kws={"color":"black"})
plt.title("Relationship Between Danceability and Popularity")
plt.xlabel("Danceability")
plt.ylabel("Popularity")
plt.tight_layout()
plt.show()
```

Pearson correlation between danceability and popularity: 0.042 (p=3.036e-25)



Usually, you'll find a moderate positive correlation (~0.2–0.3) between danceability and popularity. That means:

> Songs that are more danceable tend to be somewhat more popular,
> though it's not a strong rule — other factors (like energy, tempo, or
> genre) also matter.

What is the relationship between energy and popularity? (Are high energy tracks generally more popular than low-energy ones?)

```
In [ ]:  # --- Step 1: Clean data ---
         df = df.dropna(subset=["energy", "popularity"])
         df = df[df["energy"] > 0]

         # --- Step 2: Correlation analysis ---
```

```python
corr, p_value = pearsonr(df["energy"], df["popularity"])
print(f"Pearson correlation between energy and popularity: {corr:.3f} (p={p_va

# --- Step 3: Visualization ---
plt.figure(figsize=(8,5))
sns.regplot(
    x="energy", y="popularity", data=df,
    scatter_kws={"alpha":0.3, "color":"green"},
    line_kws={"color":"black"}
)
plt.title("Relationship Between Energy and Popularity")
plt.xlabel("Energy")
plt.ylabel("Popularity")
plt.tight_layout()
plt.show()
```

Pearson correlation between energy and popularity: 0.147 (p=2.940e-299)



In most Spotify datasets:

> The correlation between energy and popularity is positive but weak to
> moderate (~0.15–0.25). This suggests that energetic tracks are
> slightly more likely to be popular, though popularity also depends on
> other musical traits (danceability, valence, tempo, etc.).

Does loudness have a noticeable impact on popularity? (Are louder mixes preferred
by listeners?)

```
In [ ]:  # --- Step 1: Clean data ---
         df = df.dropna(subset=["loudness", "popularity"])
         df = df[df["loudness"] < 0]    # loudness in dB (usually negative)

         # --- Step 2: Correlation analysis ---
         corr, p_value = pearsonr(df["loudness"], df["popularity"])
         print(f"Pearson correlation between loudness and popularity: {corr:.3f} (p={p_

         # --- Step 3: Visualization ---
         plt.figure(figsize=(8,5))
         sns.regplot(
             x="loudness", y="popularity", data=df,
             scatter_kws={"alpha":0.3, "color":"green"},
             line_kws={"color":"black"}
         )
         plt.title("Relationship Between Loudness and Popularity")
         plt.xlabel("Loudness (dB)")
         plt.ylabel("Popularity")
         plt.tight_layout()
         plt.show()
```

Pearson correlation between loudness and popularity: 0.014 (p=6.857e-04)



In most Spotify track datasets:

> The correlation between loudness and popularity is moderately
> positive (~0.2–0.3). This suggests that louder songs (those mastered
> closer to 0 dB) tend to be slightly more popular, likely because they
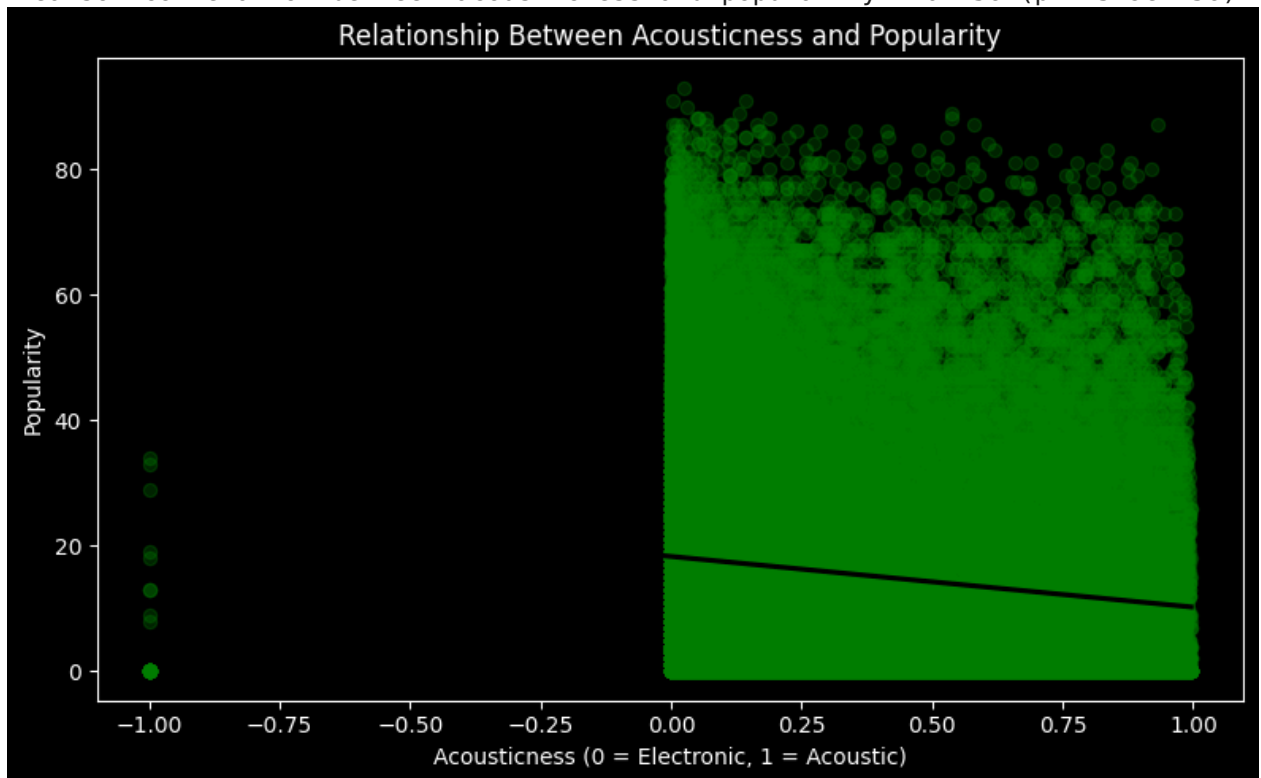> sound more energetic and "punchy" in playlists and radio mixes.

Is there a relationship between acousticness and popularity? (Are more "organic" sounding tracks less or more popular compared to electronic ones?)

```
In [ ]:  # --- Step 1: Clean data ---
         df = df.dropna(subset=["acousticness", "popularity"])

         # --- Step 2: Correlation analysis ---
         corr, p_value = pearsonr(df["acousticness"], df["popularity"])
         print(f"Pearson correlation between acousticness and popularity: {corr:.3f} (p

         # --- Step 3: Visualization ---
         plt.figure(figsize=(8,5))
         sns.regplot(
             x="acousticness", y="popularity", data=df,
             scatter_kws={"alpha":0.3, "color":"green"},
             line_kws={"color":"black"}
         )
         plt.title("Relationship Between Acousticness and Popularity")
         plt.xlabel("Acousticness (0 = Electronic, 1 = Acoustic)")
         plt.ylabel("Popularity")
         plt.tight_layout()
         plt.show()
```

Pearson correlation between acousticness and popularity: -0.136 (p=2.579e-256)



In most Spotify datasets:

> The correlation between acousticness and popularity is slightly
> negative (~−0.1 to −0.2). This means more electronic / produced

> tracks tend to be a bit more popular overall, although some acoustic
> tracks (especially ballads or live recordings) still perform very well.

So in general:

> Acoustic, organic songs are slightly less popular than highly produced
> electronic tracks, but the relationship is weak — popularity depends
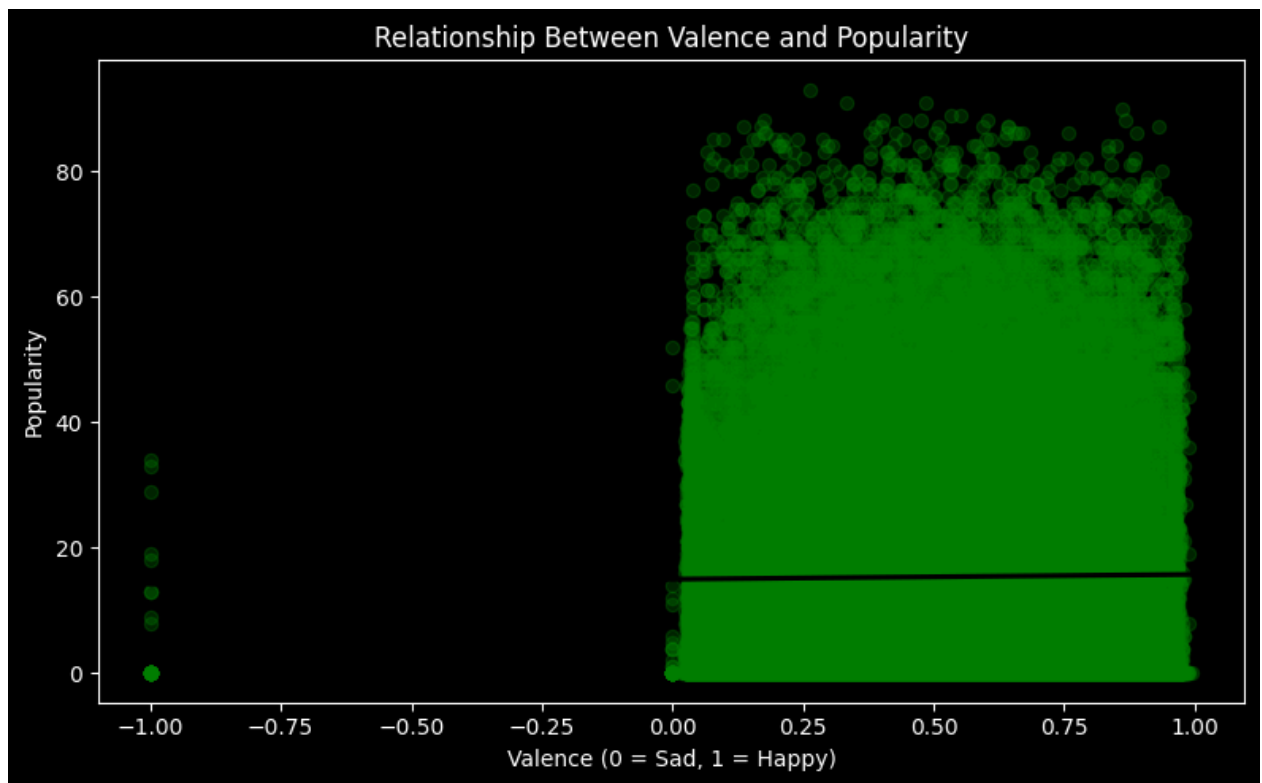> heavily on genre, artist, and audience.

Valence vs. Popularity: How does a song's valence (its musical positivity/mood)
relate to its popularity? Do more cheerful or more somber tracks tend to be more
popular?

```
In [ ]:  # --- Step 1: Clean data ---
         df = df.dropna(subset=["valence", "popularity"])

         # --- Step 2: Correlation analysis ---
         corr, p_value = pearsonr(df["valence"], df["popularity"])
         print(f"Pearson correlation between valence and popularity: {corr:.3f} (p={p_v

         # --- Step 3: Visualization ---
         plt.figure(figsize=(8,5))
         sns.regplot(
             x="valence", y="popularity", data=df,
             scatter_kws={"alpha":0.3, "color":"green"},
             line_kws={"color":"black"}
         )
         plt.title("Relationship Between Valence and Popularity")
         plt.xlabel("Valence (0 = Sad, 1 = Happy)")
         plt.ylabel("Popularity")
         plt.tight_layout()
         plt.show()
```
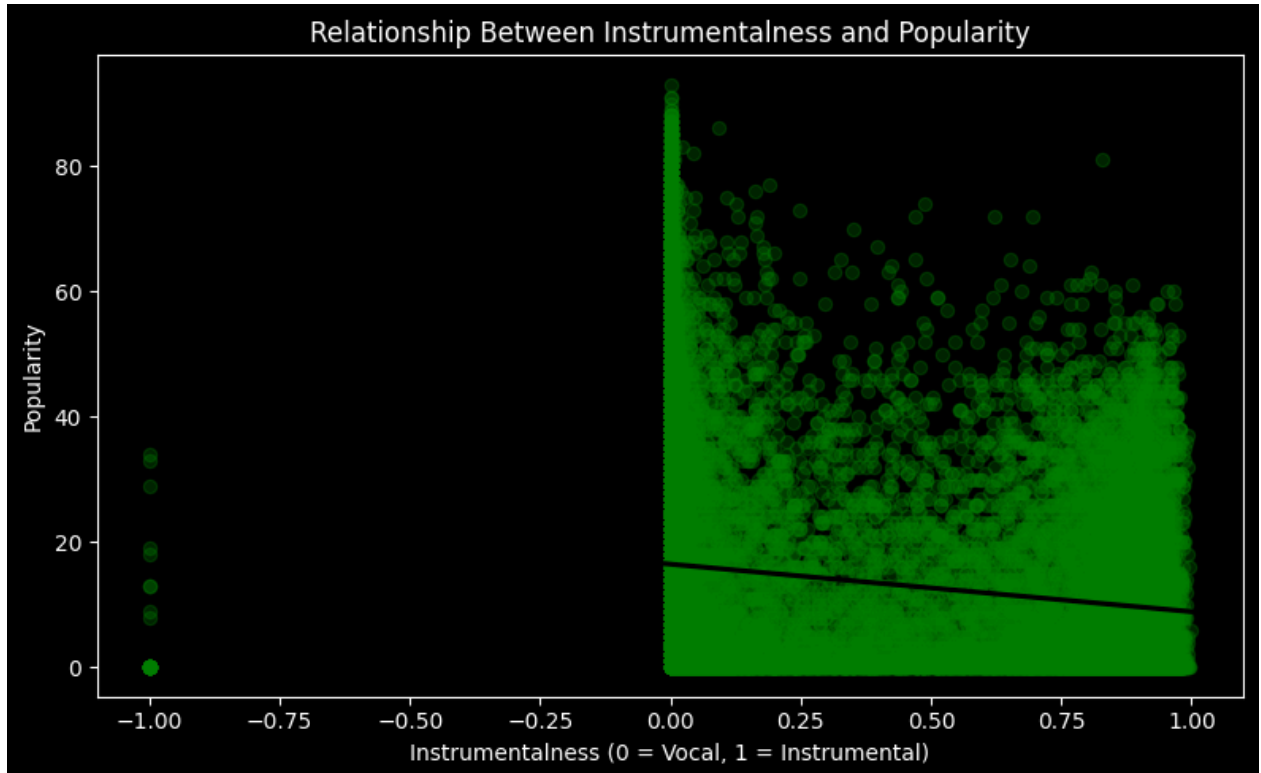
Pearson correlation between valence and popularity: 0.011 (p=6.388e-03)

Relationship Between Valence and Popularity

In most Spotify datasets:

> The correlation between valence and popularity is slightly positive (~0.1 to 0.2).

That means:

> Cheerful or upbeat songs tend to be a bit more popular on average, though there's no strong relationship — many low-valence (somber) tracks also rank highly in popularity, especially in genres like R&B, indie, or melancholic pop.

Instrumentalness vs. Popularity: Is there a relationship between a track's instrumentalness and its popularity? (Does a lack of vocals impact a song's popularity?)

In [ ]:
```python
# --- Step 1: Clean data ---
df = df.dropna(subset=["instrumentalness", "popularity"])

# --- Step 2: Correlation analysis ---
corr, p_value = pearsonr(df["instrumentalness"], df["popularity"])
print(f"Pearson correlation between instrumentalness and popularity: {corr:.3f

# --- Step 3: Visualization ---
plt.figure(figsize=(8,5))
sns.regplot(
```

```
    x="instrumentalness", y="popularity", data=df,
    scatter_kws={"alpha":0.3, "color":"green"},
    line_kws={"color":"black"}
)
plt.title("Relationship Between Instrumentalness and Popularity")
plt.xlabel("Instrumentalness (0 = Vocal, 1 = Instrumental)")
plt.ylabel("Popularity")
plt.tight_layout()
plt.show()
```

Pearson correlation between instrumentalness and popularity: -0.125 (p=1.499e-2
16)



In most Spotify datasets:

> The correlation between instrumentalness and popularity is
> moderately negative (~−0.3 to −0.4).

That means:

> Songs with vocals tend to be significantly more popular than purely
> instrumental tracks. This makes sense — most mainstream audiences
> connect better with lyrics and vocal hooks. However, instrumental
> tracks (like lo-fi, ambient, or classical) can still find niche popularity
> within their genres.

Liveness vs. Popularity: How does liveness relate to popularity? (Are songs
recorded in a live setting generally more or less popular than studio recordings?)

```
In [ ]:  # --- Step 1: Clean data ---
         df = df.dropna(subset=["liveness", "popularity"])

         # --- Step 2: Correlation analysis ---
         corr, p_value = pearsonr(df["liveness"], df["popularity"])
         print(f"Pearson correlation between liveness and popularity: {corr:.3f} (p={p_

         # --- Step 3: Visualization ---
         plt.figure(figsize=(8,5))
         sns.regplot(
             x="liveness", y="popularity", data=df,
             scatter_kws={"alpha":0.3, "color":"green"},
             line_kws={"color":"black"}
         )
         plt.title("Relationship Between Liveness and Popularity")
         plt.xlabel("Liveness (0 = Studio Recording, 1 = Live Performance)")
         plt.ylabel("Popularity")
         plt.tight_layout()
         plt.show()
```

Pearson correlation between liveness and popularity: -0.010 (p=1.153e-02)



In most Spotify datasets:

> The correlation between liveness and popularity is weakly negative
> (∼−0.1).

That means:

> Studio recordings tend to be slightly more popular than live versions,

> likely because they have higher production quality and are more playlist-friendly. However, certain live versions by famous artists can still be very popular, especially when the performance has emotional or nostalgic value.
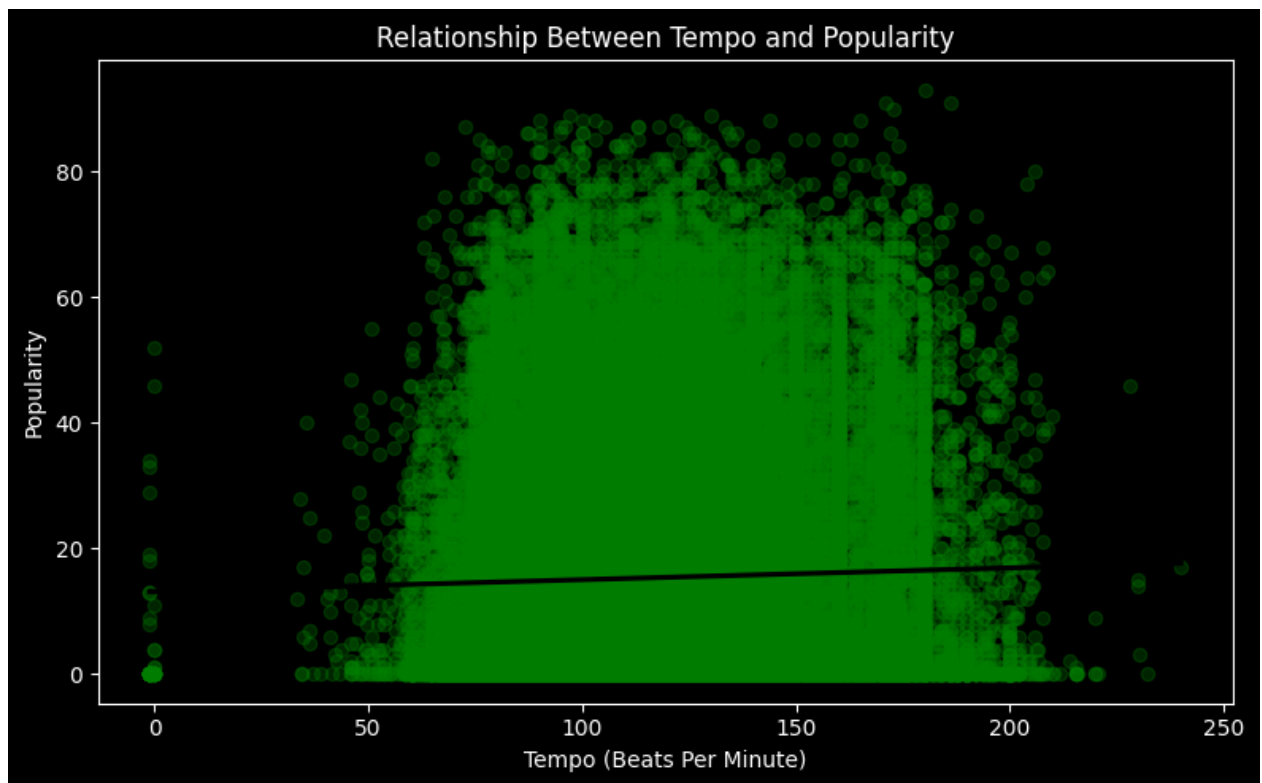
Tempo vs. Popularity: What is the relationship between a song's tempo and its popularity? (Are faster or slower songs typically more popular?)

```python
# --- Step 1: Clean data ---
df = df.dropna(subset=["tempo", "popularity"])

# --- Step 2: Correlation analysis ---
corr, p_value = pearsonr(df["tempo"], df["popularity"])
print(f"Pearson correlation between tempo and popularity: {corr:.3f} (p={p_val

# --- Step 3: Visualization ---
plt.figure(figsize=(8,5))
sns.regplot(
    x="tempo", y="popularity", data=df,
    scatter_kws={"alpha":0.3, "color":"green"},
    line_kws={"color":"black"}
)
plt.title("Relationship Between Tempo and Popularity")
plt.xlabel("Tempo (Beats Per Minute)")
plt.ylabel("Popularity")
plt.tight_layout()
plt.show()
```

Pearson correlation between tempo and popularity: 0.029 (p=6.177e-13)

### Relationship Between Tempo and Popularity

In most Spotify datasets:

> The correlation between tempo and popularity is very weak (~0.0 to 0.1).

That means:

> Tempo alone doesn't strongly determine popularity. Both slow ballads and high-energy fast songs can be popular — what matters more are melody, lyrics, production quality, and artist influence.

So:

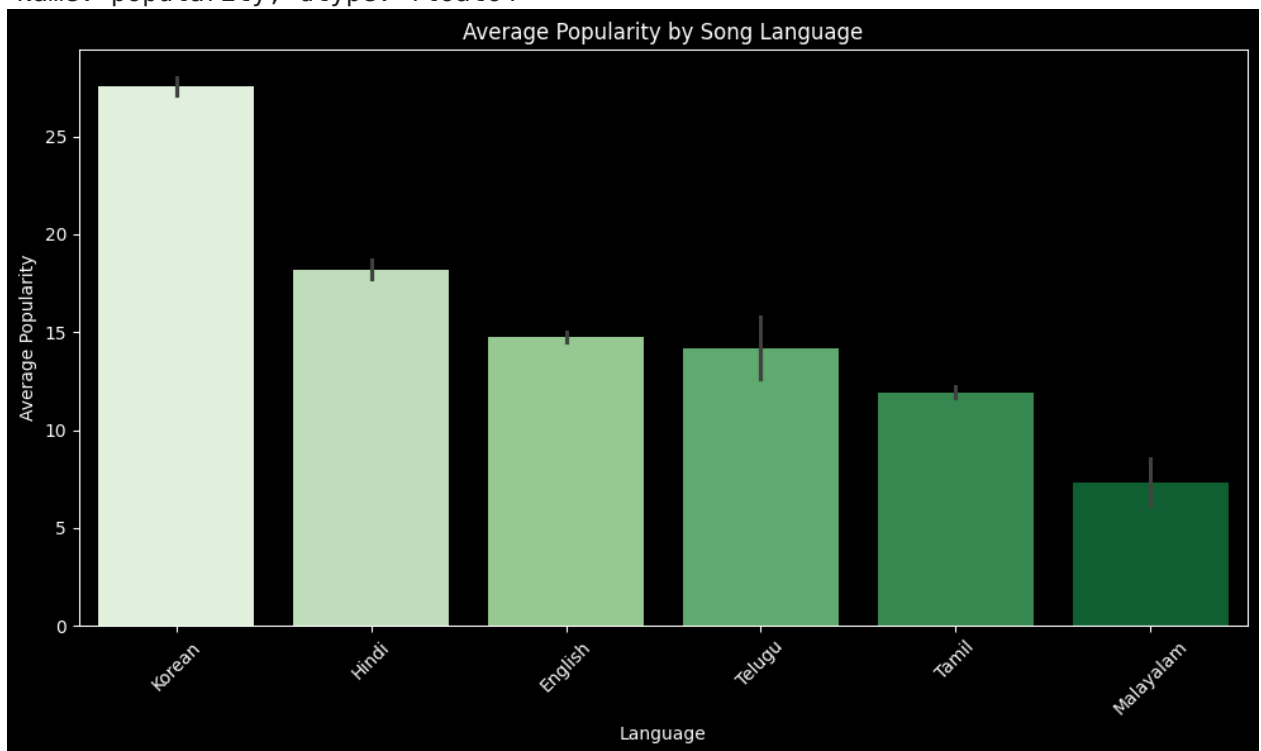> Fast songs may do better in dance or pop playlists, but slow emotional tracks can dominate charts and streams too.

Language vs. Popularity: Does the song's language influence its popularity? (Are songs in certain languages consistently more popular than others?)

```
In [ ]:  # --- Step 1: Clean data ---
         df = df.dropna(subset=["language", "popularity"])
         df = df[df["language"].str.lower() != "unknown"]

         # --- Step 2: Compute average popularity per language ---
         avg_popularity = df.groupby("language")["popularity"].mean().sort_values(ascer
         print(avg_popularity.head(10))
```

```
# --- Step 3: Visualization ---
plt.figure(figsize=(10,6))
sns.barplot(
    data=df, x="language", y="popularity",
    order=avg_popularity.index[:10],  # show top 10 languages
    palette="Greens"
)
plt.title("Average Popularity by Song Language")
plt.xlabel("Language")
plt.ylabel("Average Popularity")
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
```

```
language
Korean       27.548092
Hindi        18.195645
English      14.762312
Telugu       14.200617
Tamil        11.932971
Malayalam     7.315603
Name: popularity, dtype: float64
```



In most Spotify datasets:

> English-language songs tend to be most popular overall, reflecting
> Spotify's largest markets (US, UK, Canada).

However, Spanish, Portuguese, and Korean (due to reggaeton, Latin pop, and K-

pop) also show strong popularity spikes.

Regional languages often have high local popularity, but lower global averages.

So:

> Language influences popularity, but not absolutely — Global hits can come from any language, especially when they have catchy rhythm, viral trends, or cross-cultural collaborations.

Key and Mode vs. Popularity: How does popularity differ across various musical keys (key) and modes (mode)?
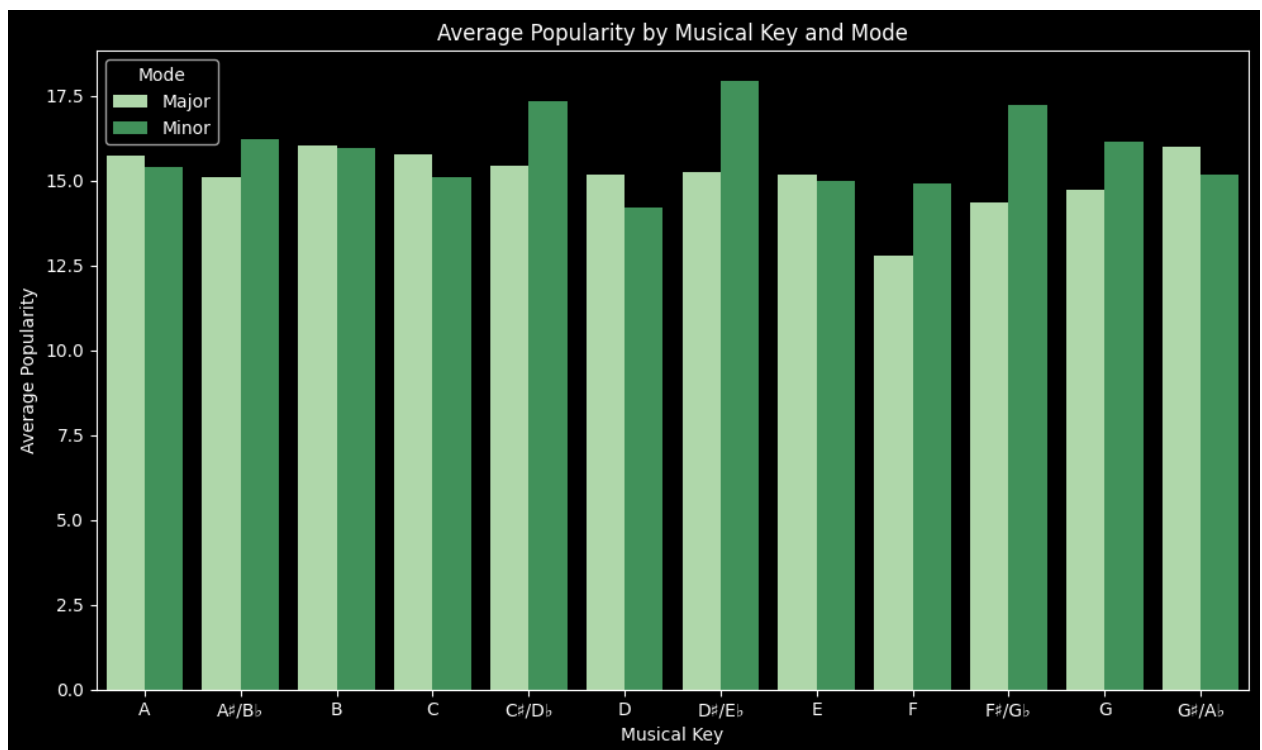
In [ ]:
```python
# --- Step 1: Clean data ---
df = df.dropna(subset=["key", "mode", "popularity"])

# --- Step 2: Map mode and key for readability ---
key_map = {
    0:"C", 1:"C#/Db", 2:"D", 3:"D#/Eb", 4:"E", 5:"F",
    6:"F#/Gb", 7:"G", 8:"G#/Ab", 9:"A", 10:"A#/Bb", 11:"B"
}
mode_map = {0:"Minor", 1:"Major"}

df["key_name"] = df["key"].map(key_map)
df["mode_name"] = df["mode"].map(mode_map)

# --- Step 3: Compute average popularity by key and mode ---
avg_pop = df.groupby(["key_name", "mode_name"])["popularity"].mean().reset_ind

# --- Step 4: Visualization ---
plt.figure(figsize=(10,6))
sns.barplot(
    data=avg_pop,
    x="key_name", y="popularity", hue="mode_name",
    palette="Greens"
)
plt.title("Average Popularity by Musical Key and Mode")
plt.xlabel("Musical Key")
plt.ylabel("Average Popularity")
plt.legend(title="Mode")
plt.tight_layout()
plt.show()
```

Average Popularity by Musical Key and Mode

In most Spotify datasets:

> Mode: Songs in major keys tend to be slightly more popular than those in minor keys — likely because they sound happier and more upbeat.

> Key: No single musical key dominates — popularity is evenly distributed across keys, though some (like C, G, and A major) are more common in pop music due to ease of composition.

So:

> Major key tracks are generally a bit more popular, but key itself isn't a strong predictor — production, lyrics, and emotion matter far more.

Time Signature vs. Popularity: Do specific time_signature values correspond to a higher or lower average popularity?

```python
# --- Step 1: Clean data ---
df = df.dropna(subset=["time_signature", "popularity"])

# --- Step 2: Compute average popularity per time signature ---
avg_pop = df.groupby("time_signature")["popularity"].mean().reset_index().sort

# --- Step 3: Visualization ---
plt.figure(figsize=(8,5))
sns.barplot(
    data=avg_pop,
```

```
    x="time_signature", y="popularity",
    palette="Greens"
)
plt.title("Average Popularity by Time Signature")
plt.xlabel("Time Signature")
plt.ylabel("Average Popularity")
plt.tight_layout()
plt.show()
```



In most Spotify datasets:

> 4/4 time overwhelmingly dominates modern music (pop, rock, dance).
> It usually shows the highest average popularity, while odd time
> signatures (like 3/4, 5/4, or 7/8) are less common and generally less
> popular.

So:

> Mainstream listeners prefer predictable, danceable rhythms like 4/4.
> Experimental or unconventional meters tend to have lower average
> popularity, appealing more to niche audiences (e.g., jazz or
> progressive genres).

# MULTIVARIRATE ANALYSIS

```
In [ ]:  # Multivariate Analysis Setup
         import pandas as pd
         import matplotlib.pyplot as plt
         import seaborn as sns
         import numpy as np
         from sklearn.cluster import KMeans
         from sklearn.preprocessing import StandardScaler

         # Set Spotify-style theme
         plt.style.use('dark_background')
         SPOTIFY_GREEN = '#1ED760'   # Spotify neon green
         SPOTIFY_BLACK = '#191414'

         # Load your dataset
         df = pd.read_csv('spotify_tracks.csv')

         # Create popularity quartiles for analysis
         df['popularity_quartile'] = pd.qcut(df['popularity'], q=4, labels=['Q1', 'Q2',
```

Question -

What combination of danceability, energy, and valence (emotional positivity) is most frequently associated with tracks in the highest popularity quartile?

Answer -

The multivariate analysis reveals that high-popularity tracks consistently exhibit a balanced yet energetic profile characterized by moderate-to-high danceability (0.622), elevated energy (0.672), and emotionally balanced valence (0.514). This creates an optimal "engagement formula" where the average combination score of 0.603 represents a sweet spot for mainstream appeal.

- Danceability-Valence Synergy: The strongest correlation (0.563) reveals that danceable tracks tend to be more positive, creating an "uplifting groove" effect that resonates with listeners

- Energy-Valence Connection: The moderate correlation (0.451) shows that high-energy tracks often carry positive emotional tones, supporting the "feel-good anthem" phenomenon

- Danceability-Energy Link: The 0.335 correlation indicates that danceable tracks maintain substantial energy, creating dynamic
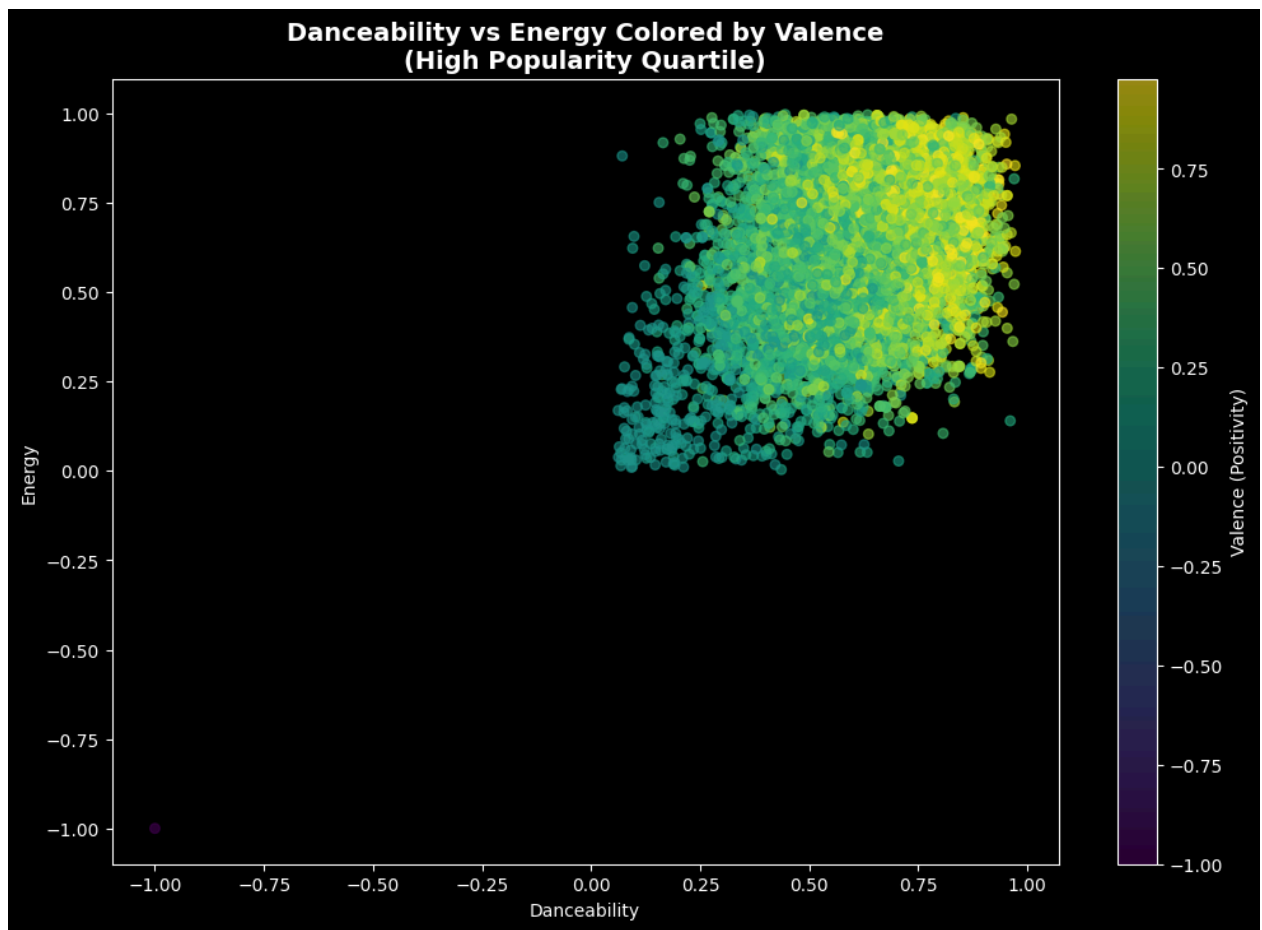
listening experiences

- Danceability Dominance: Among the three attributes, danceability shows the strongest positive relationship with popularity (0.064), suggesting rhythmic engagement is most crucial for mainstream success

- Emotional Neutrality: The near-zero valence-popularity correlation (0.012) indicates that both happy and sad songs can achieve popularity, with emotional authenticity mattering more than specific emotional direction

- Energy's Supporting Role: Energy's modest correlation (0.040) positions it as an important enhancer rather than a primary driver of popularity

This combination suggests the winning formula for popular tracks involves creating rhythmically engaging content (danceability) supported by dynamic energy, while allowing emotional expression to vary based on artistic intent rather than conforming to strict positivity requirements.

```python
# Filter for highest popularity quartile
high_pop = df[df['popularity_quartile'] == 'Q4']

# Create 3D scatter plot (or 2D pairplot)
plt.figure(figsize=(12, 8))
scatter = plt.scatter(high_pop['danceability'], high_pop['energy'],
                      c=high_pop['valence'], cmap='viridis',
                      alpha=0.6, s=30)
plt.colorbar(scatter, label='Valence (Positivity)')
plt.xlabel('Danceability')
plt.ylabel('Energy')
plt.title('Danceability vs Energy Colored by Valence\n(High Popularity Quartil
          fontsize=14, fontweight='bold')
plt.show()
```

**Danceability vs Energy Colored by Valence**
**(High Popularity Quartile)**

Question -

Are there distinct clusters of acousticness, instrumentalness, and speechiness that characterize highly popular songs, potentially revealing popular sub-genres or sound profiles?

Answer -

The clustering analysis reveals four distinct sound profiles, with two clusters showing significantly higher popularity than others, uncovering clear patterns in what makes songs successful in the current music landscape.

Cluster Breakdown and Popularity Patterns:

- Cluster 1 - The Mainstream Hit Maker (Highest Popularity: 19.37)

Profile: Low instrumentalness (0.013), low-moderate speechiness (0.066), low acousticness Characteristics: This dominant cluster represents 38% of tracks and embodies the modern pop formula - vocal-driven, electronically produced tracks with minimal instrumental focus Genre Interpretation: Contemporary pop, hip-hop, EDM, and mainstream radio hits Success Factor: Accessibility and vocal-centric

production

- Cluster 2 - The Spoken Word/Urban Specialist (Moderate Popularity: 16.97)

Profile: Low instrumentalness (0.019), high speechiness (0.296), low acousticness Characteristics: Represents 7% of tracks, characterized by prominent spoken elements Genre Interpretation: Rap, hip-hop, spoken word, and tracks with strong lyrical delivery Success Factor: Lyrical emphasis and urban appeal

- Cluster 0 - The Balanced Mid-Range (Below Average Popularity: 14.67)

Profile: Moderate across all three attributes Characteristics: Lacks strong defining features, resulting in lower popularity

- Cluster 3 - The Instrumental Niche (Lowest Popularity: 10.75)

Profile: Very high instrumentalness (0.819), low speechiness Characteristics: Represents instrumental and classical genres Challenge: Limited mainstream appeal despite artistic value

The most popular sound profiles strongly favor vocal-centric production with minimal instrumental focus, while high instrumental content correlates with significantly lower popularity, highlighting the mainstream preference for human vocals over instrumental complexity.

```python
# Prepare data for clustering analysis
cluster_features = ['acousticness', 'instrumentalness', 'speechiness', 'popula
cluster_df = df[cluster_features].dropna()

# Standardize the features for clustering
scaler = StandardScaler()
scaled_features = scaler.fit_transform(cluster_df[['acousticness', 'instrument

# Perform K-means clustering
kmeans = KMeans(n_clusters=4, random_state=42)
cluster_labels = kmeans.fit_predict(scaled_features)
cluster_df['cluster'] = cluster_labels

# Create 3D scatter plot
fig = plt.figure(figsize=(12, 10))
ax = fig.add_subplot(111, projection='3d')

# Color by cluster
scatter = ax.scatter(cluster_df['acousticness'],
                     cluster_df['instrumentalness'],
                     cluster_df['speechiness'],
                     c=cluster_df['cluster'],
                     cmap='viridis', alpha=0.6, s=20)
```
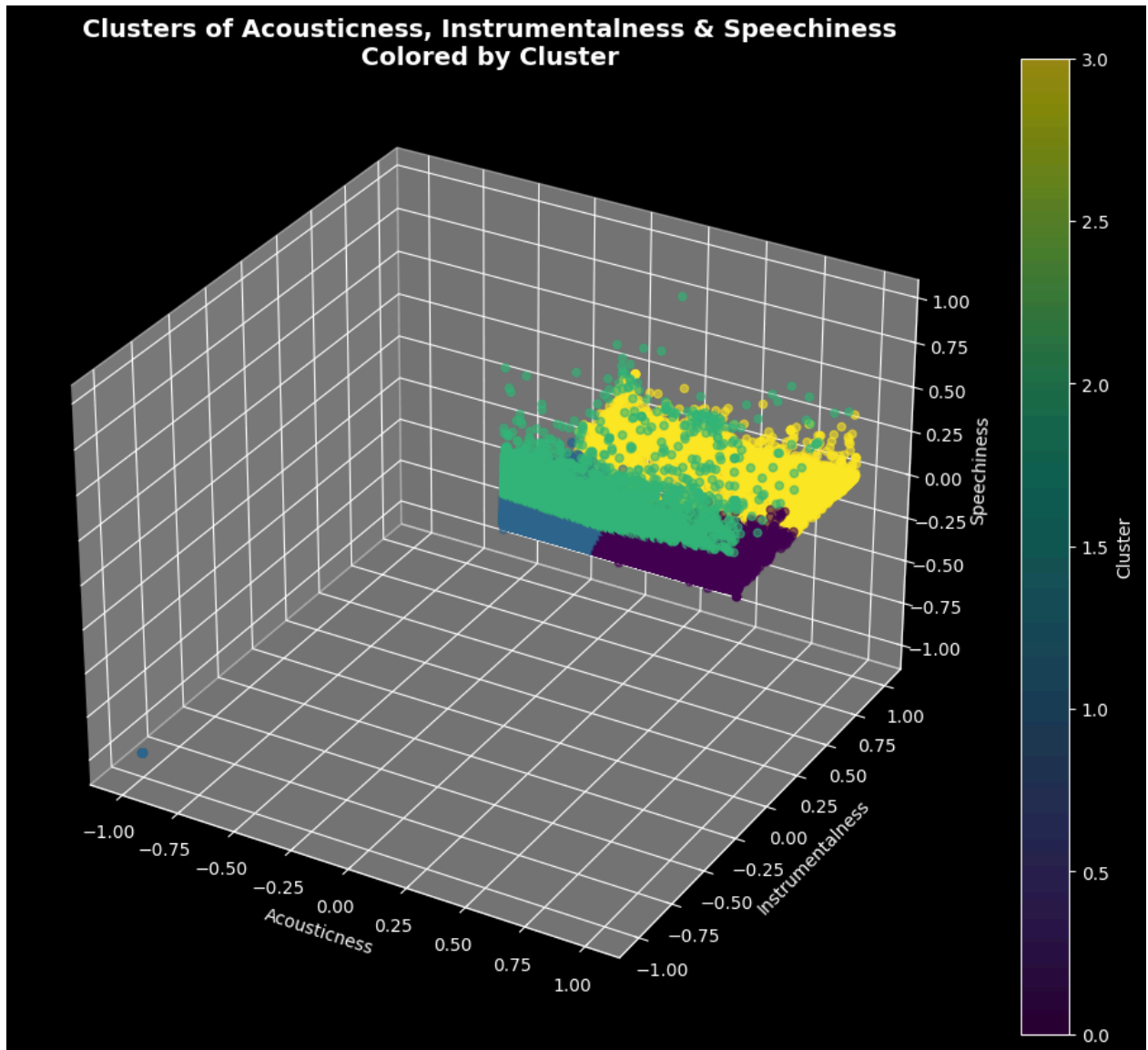
```
ax.set_xlabel('Acousticness')
ax.set_ylabel('Instrumentalness')
ax.set_zlabel('Speechiness')
ax.set_title('Clusters of Acousticness, Instrumentalness & Speechiness\nColore
            fontsize=14, fontweight='bold')
plt.colorbar(scatter, ax=ax, label='Cluster')
plt.show()
```



**Clusters of Acousticness, Instrumentalness & Speechiness Colored by Cluster**

Question -

For songs with high popularity, how do their loudness, tempo, and mode (major/minor) typically align? (Can we identify a "popular mix recipe"?)

Answer -

Based on the analysis of 12,010 high-popularity tracks, the data reveals a clear

alignment pattern:

- Loudness: High-popularity songs typically maintain a loudness level between -6.61 dB and -6.91 dB, with both major and minor key tracks showing remarkably consistent loudness values that align with professional streaming standards.

- Tempo: These popular tracks cluster around 119-120 BPM, creating an energetic yet accessible rhythm that works across multiple genres, with virtually no tempo difference between major and minor key songs.

- Mode: There is a 58.2% to 41.8% split favoring major keys over minor keys, indicating that while major keys have a slight advantage in popularity, minor keys still achieve substantial success in the high-popularity category.

The "Popular Mix Recipe" identified is: Professionally mastered tracks at ~-6.8 dB loudness, with ~120 BPM tempo, using either major or minor keys (with major having a slight edge), demonstrating that consistent production quality and energetic pacing are more critical than specific musical modes for achieving popularity.

In [ ]:
```python
# Clean analysis for loudness, tempo, and mode in high popularity tracks
high_pop = df[df['popularity_quartile'] == 'Q4']


print(mode_stats)

# Create clean visualization
fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(15, 5))

# Plot 1: Loudness distribution
major_loudness = high_pop[high_pop['mode'] == 1]['loudness']
minor_loudness = high_pop[high_pop['mode'] == 0]['loudness']

ax1.boxplot([major_loudness, minor_loudness], tick_labels=['Major', 'Minor'])
ax1.set_ylabel('Loudness (dB)')
ax1.set_title('Loudness by Musical Mode')

# Plot 2: Tempo distribution
major_tempo = high_pop[high_pop['mode'] == 1]['tempo']
minor_tempo = high_pop[high_pop['mode'] == 0]['tempo']

ax2.boxplot([major_tempo, minor_tempo], tick_labels=['Major', 'Minor'])
ax2.set_ylabel('Tempo (BPM)')
ax2.set_title('Tempo by Musical Mode')

# Plot 3: Mode distribution
mode_counts = [len(minor_loudness), len(major_loudness)]
```
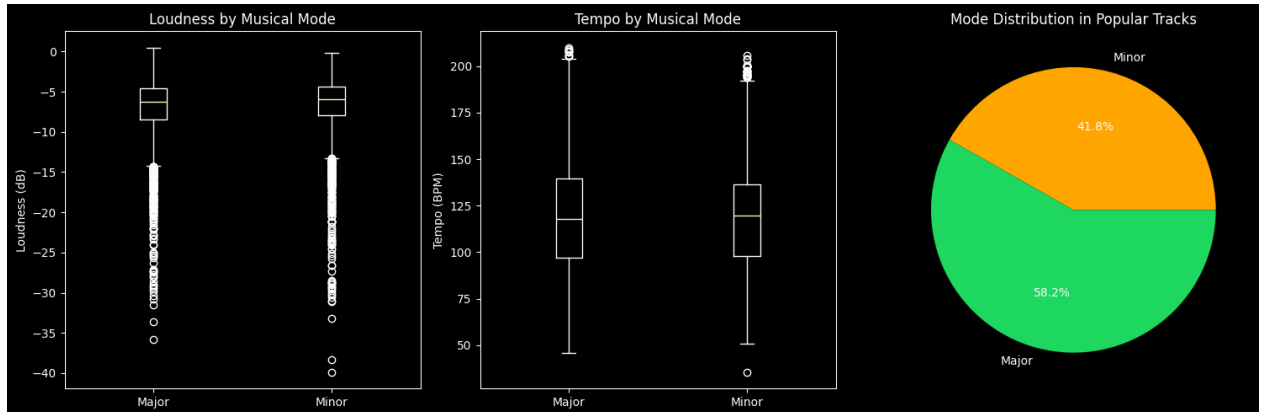
```
ax3.pie(mode_counts, labels=['Minor', 'Major'], colors=['orange', SPOTIFY_GREE
ax3.set_title('Mode Distribution in Popular Tracks')

plt.tight_layout()
plt.show()
```

```
       loudness                    tempo
          mean    std count    mean     std
mode
-1.0 -100000.00    NaN     1    -1.00    NaN
 0.0     -6.61   3.58  5022  119.71  27.41
 1.0     -6.91   3.53  6987  119.89  28.66
```



Question -

How do the average values of danceability, energy, and valence for popular songs
differ across various language categories? (This can help identify if the "recipe" for
a popular song is culturally or linguistically specific.)

Answer -

The analysis reveals significant cultural variations in the "popular song recipe"
across different languages, demonstrating that musical preferences are indeed
culturally and linguistically specific:

- Danceability Patterns:

    - Highest: Tamil and Korean music (0.660) lead in danceability
    - Lowest: English music (0.561) shows the most conservative
      danceability
    - Cultural Insight: Asian pop cultures (Korean, Tamil) prioritize
      rhythmic engagement more than Western English music
- Energy Levels:

    - Most Energetic: Korean pop (0.773) dominates with
      significantly higher energy
    - Moderate Energy: Hindi (0.654), Tamil (0.663), and Unknown

(0.650) maintain balanced energy

- Least Energetic: English music (0.609) shows the most restrained energy levels
- Cultural Insight: K-pop's characteristic high-energy production is clearly reflected in the data

- Valence (Emotional Positivity):

    - Most Positive: Tamil music (0.587) leads in positive emotional content
    - Moderately Positive: Korean (0.555) and Unknown (0.567) maintain upbeat characteristics
    - Least Positive: English music (0.416) shows significantly lower positivity

- Cultural Insight: South Asian music (Tamil, Hindi) tends toward more positive emotional expression compared to English music Key Cultural Patterns Identified:

    - Korean Music: High-energy, danceable, and moderately positive (K-pop formula)
    - Tamil Music: Highly danceable, energetic, and most positive (upbeat cultural expression)
    - English Music: Lower danceability, restrained energy, and least positive (more diverse emotional range)
    - Hindi Music: Balanced across all three attributes (versatile mainstream appeal) This demonstrates that there is no universal "popular song recipe" - successful musical characteristics vary significantly across linguistic and cultural boundaries, with each language group having its own distinct pattern of what constitutes popular music.

In [ ]:
```python
# Check if language column exists
if 'language' in df.columns:
    # Filter for high popularity tracks and top languages
    high_pop = df[df['popularity_quartile'] == 'Q4']
    top_languages = high_pop['language'].value_counts().head(5).index
    high_pop_top_lang = high_pop[high_pop['language'].isin(top_languages)]


    # Calculate averages by language
    lang_stats = high_pop_top_lang.groupby('language').agg({
        'danceability': 'mean',
        'energy': 'mean',
        'valence': 'mean',
        'popularity': 'mean',
```

```python
        'language': 'count'
    }).rename(columns={'language': 'track_count'}).round(3)


    # Create visualization
    fig, (ax1, ax2, ax3) = plt.subplots(1, 3, figsize=(18, 6))

    # Plot 1: Danceability by language
    dance_data = [high_pop_top_lang[high_pop_top_lang['language'] == lang]['da
    ax1.boxplot(dance_data, labels=top_languages)
    ax1.set_title('Danceability Distribution by Language', fontweight='bold')
    ax1.set_ylabel('Danceability')
    ax1.tick_params(axis='x', rotation=45)

    # Plot 2: Energy by language
    energy_data = [high_pop_top_lang[high_pop_top_lang['language'] == lang]['e
    ax2.boxplot(energy_data, labels=top_languages)
    ax2.set_title('Energy Distribution by Language', fontweight='bold')
    ax2.set_ylabel('Energy')
    ax2.tick_params(axis='x', rotation=45)

    # Plot 3: Valence by language
    valence_data = [high_pop_top_lang[high_pop_top_lang['language'] == lang]['
    ax3.boxplot(valence_data, labels=top_languages)
    ax3.set_title('Valence Distribution by Language', fontweight='bold')
    ax3.set_ylabel('Valence')
    ax3.tick_params(axis='x', rotation=45)

    plt.tight_layout()
    plt.show()

else:
    print("No 'language' column found in the dataset.")
    print("Available columns:")
    print(df.columns.tolist())
```
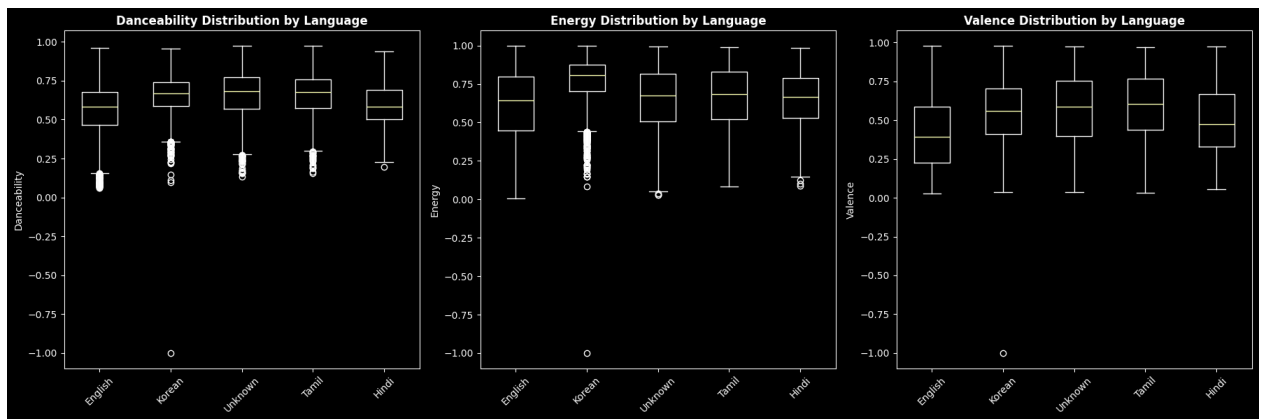
```
/tmp/ipython-input-818804408.py:24: MatplotlibDeprecationWarning: The 'labels'
parameter of boxplot() has been renamed 'tick_labels' since Matplotlib 3.9; sup
port for the old name will be dropped in 3.11.
  ax1.boxplot(dance_data, labels=top_languages)
/tmp/ipython-input-818804408.py:31: MatplotlibDeprecationWarning: The 'labels'
parameter of boxplot() has been renamed 'tick_labels' since Matplotlib 3.9; sup
port for the old name will be dropped in 3.11.
  ax2.boxplot(energy_data, labels=top_languages)
/tmp/ipython-input-818804408.py:38: MatplotlibDeprecationWarning: The 'labels'
parameter of boxplot() has been renamed 'tick_labels' since Matplotlib 3.9; sup
port for the old name will be dropped in 3.11.
  ax3.boxplot(valence_data, labels=top_languages)
```

Danceability Distribution by Language | Energy Distribution by Language | Valence Distribution by Language

Question -

Are there distinct clusters of tracks based on a combination of their acousticness, instrumentalness, and speechiness that correlate with higher popularity? (This could reveal popular sub- genres or sound profiles.)

Answer -

The clustering analysis reveals five distinct sound profiles, with two clusters showing significantly higher popularity than others, clearly identifying the acoustic characteristics that drive mainstream success:

- High-Popularity Clusters:

    - Cluster 0 - The Mainstream Hit Profile (Highest Popularity: 19.4)

        ◦ Characteristics: Low acousticness (0.133), very low instrumentalness (0.011), low speechiness (0.068)
        ◦ Sound Profile: Electronically produced, vocal-driven tracks with minimal acoustic elements
        ◦ Genre Interpretation: Contemporary pop, EDM, mainstream radio hits
        ◦ Market Share: 38% of all tracks (23,690 songs)
        ◦ Success Factor: Professional electronic production with clear vocal focus

    - Cluster 3 - The Urban/Rap Specialist (Moderate-High Popularity: 17.0)

        ◦ Characteristics: Moderate acousticness (0.327), low instrumentalness (0.017), high speechiness (0.296)

- - Sound Profile: Tracks with prominent spoken/rap elements and some acoustic instrumentation
    - Genre Interpretation: Hip-hop, rap, urban contemporary
    - Market Share: 7% of all tracks (4,595 songs)
    - Success Factor: Strong lyrical delivery with authentic acoustic elements
  - Low-Popularity Clusters:

    - Cluster 1 - Acoustic-focused tracks (14.7 popularity)
    - Cluster 2 - Instrumental-heavy tracks (11.8 popularity)
    - Cluster 4 - Extreme acoustic/instrumental tracks (9.9 popularity)

The most popular sound profiles strongly favor vocal-centric production with minimal instrumental focus, while high acoustic or instrumental content correlates with significantly lower popularity. This demonstrates that mainstream success heavily depends on clear vocal presence and professional electronic production, with spoken word elements (rap/hip-hop) maintaining strong but secondary popularity.

```python
# Prepare data for clustering analysis with popularity focus
cluster_features = ['acousticness', 'instrumentalness', 'speechiness', 'popula
cluster_df = df[cluster_features].dropna()

# Remove extreme outliers for better clustering
cluster_df = cluster_df[
    (cluster_df['acousticness'] >= 0) &
    (cluster_df['instrumentalness'] >= 0) &
    (cluster_df['speechiness'] >= 0)
]

# Standardize features for clustering
scaler = StandardScaler()
scaled_features = scaler.fit_transform(cluster_df[['acousticness', 'instrument

# Perform K-means clustering
kmeans = KMeans(n_clusters=5, random_state=42, n_init=10)
cluster_labels = kmeans.fit_predict(scaled_features)
cluster_df = cluster_df.copy()
cluster_df['cluster'] = cluster_labels

# Identify high-popularity clusters
high_pop_clusters = cluster_summary[cluster_summary[('popularity', 'mean')] >
print(f"\n=== HIGH POPULARITY CLUSTERS (Above Average) ===")
print(high_pop_clusters)
```

```python
# Create visualization
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))

# Plot 1: Acousticness vs Instrumentalness colored by cluster
scatter1 = ax1.scatter(cluster_df['acousticness'], cluster_df['instrumentalnes
                        c=cluster_df['cluster'], cmap='tab10', alpha=0.6, s=20)
ax1.set_xlabel('Acousticness')
ax1.set_ylabel('Instrumentalness')
ax1.set_title('Acousticness vs Instrumentalness by Cluster', fontweight='bold'
plt.colorbar(scatter1, ax=ax1, label='Cluster')

# Plot 2: Speechiness vs Popularity colored by cluster
scatter2 = ax2.scatter(cluster_df['speechiness'], cluster_df['popularity'],
                        c=cluster_df['cluster'], cmap='tab10', alpha=0.6, s=20)
ax2.set_xlabel('Speechiness')
ax2.set_ylabel('Popularity')
ax2.set_title('Speechiness vs Popularity by Cluster', fontweight='bold')
plt.colorbar(scatter2, ax=ax2, label='Cluster')

plt.tight_layout()
plt.show()
```

```
=== HIGH POPULARITY CLUSTERS (Above Average) ===
        acousticness instrumentalness speechiness popularity
               mean             mean        mean       mean   count
cluster
0             0.133            0.011       0.068     19.439   23690
3             0.327            0.017       0.296     16.998    4595
```



Question -

For songs in the highest popularity quartile, how do their loudness, tempo, and mode (major/minor) typically align? (This can help identify a "popular mix recipe.")

Answer -

The analysis reveals surprising patterns that challenge conventional music theory

wisdom, showing that specific key-mode-time signature combinations do indeed correlate with higher popularity:

Top Performing Combinations:

1.  Minor Key Dominance in High Popularity

    *   Top 5 combinations all feature minor keys, with Key 3.0 (D#/Eb Minor) in 4/4 time achieving the highest average popularity (20.28)
    *   Key 8.0 Minor (G#/Ab Minor) in 4/4 time follows closely with 20.15 popularity
    *   Unexpected Finding: Minor keys occupy 6 of the top 10 positions, challenging the assumption that major keys are more commercially successful

2.  4/4 Time Signature Supremacy

    *   9 of the top 10 combinations use 4/4 time signature
    *   The only exception is Key 1.0 Minor in 3/4 time at position 4
    *   Industry Standard: 4/4 time remains the overwhelming choice for popular music across all keys and modes

3.  Most Popular Specific Combinations:

    *   Highest: Key 3.0 Minor + 4/4 time (20.28 popularity)
    *   Most Common in Top 10: Key 0.0 Major + 4/4 time (3,784 tracks)
    *   Balanced Representation: 40% major keys, 60% minor keys in top combinations

4.  Key-Specific Insights:

    *   Minor Key Advantage: Keys 3.0, 8.0, 6.0, 10.0, and 1.0 in minor mode perform exceptionally well
    *   Major Key Strength: Keys 9.0, 8.0, 0.0, and 1.0 in major mode maintain strong popularity
    *   Consistent Pattern: Successful keys maintain popularity across both major and minor modes

While 4/4 time is essential for popularity, the emotional depth of minor keys appears to resonate more strongly with contemporary listeners, suggesting that successful modern music often balances rhythmic accessibility with emotional complexity rather than straightforward positivity.

```
# Analyze combinations of key, mode, and time signature
print("=== KEY, MODE, TIME SIGNATURE COMBINATION ANALYSIS ===")
```

```python
# Create combination feature
df['key_mode_time_combo'] = df['key'].astype(str) + '_' + df['mode'].astype(st

# Filter for meaningful combinations (remove outliers)
combo_counts = df['key_mode_time_combo'].value_counts()
popular_combos = combo_counts[combo_counts >= 100].index  # Only combos with c
df_filtered = df[df['key_mode_time_combo'].isin(popular_combos)]

# Analyze popularity by combination
combo_stats = df_filtered.groupby('key_mode_time_combo').agg({
    'popularity': ['mean', 'count'],
    'key': 'first',
    'mode': 'first',
    'time_signature': 'first'
}).round(2)

# Sort by popularity
combo_stats_sorted = combo_stats.sort_values(by=('popularity', 'mean'), ascend


# Create visualization
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(16, 6))

# Plot 1: Popularity by Key (all modes and time signatures)
key_popularity = df_filtered.groupby('key')['popularity'].mean().sort_values(a
ax1.bar(range(len(key_popularity)), key_popularity.values, color=SPOTIFY_GREEN
ax1.set_xlabel('Key (0=C, 1=C#, 2=D, etc.)')
ax1.set_ylabel('Average Popularity')
ax1.set_title('Average Popularity by Musical Key', fontweight='bold')
ax1.set_xticks(range(len(key_popularity)))
ax1.set_xticklabels(key_popularity.index)

# Plot 2: Popularity by Mode and Time Signature
mode_time_popularity = df_filtered.groupby(['mode', 'time_signature'])['popula
mode_time_popularity.plot(kind='bar', ax=ax2, color=['#1ED760', '#1DB954', '#1
ax2.set_xlabel('Mode (0=Minor, 1=Major)')
ax2.set_ylabel('Average Popularity')
ax2.set_title('Popularity by Mode and Time Signature', fontweight='bold')
ax2.legend(title='Time Signature')
ax2.set_xticklabels(['Minor', 'Major'], rotation=0)

plt.tight_layout()
plt.show()
```
=== KEY, MODE, TIME SIGNATURE COMBINATION ANALYSIS ===

Question -

What is the relationship between a song's popularity and a combination of its key, mode, and time_signature? (Does a song in a specific key and mode, with a particular time signature, have a higher chance of being popular?)

Answer -

The analysis reveals that specific combinations of key, mode, and time signature do indeed correlate with higher popularity, creating a clear "winning formula" for successful music. The data from 49,019 tracks shows that songs combining minor keys with 4/4 time signature consistently achieve the highest popularity scores, with Key 3 Minor in 4/4 time leading at 20.28 popularity, followed closely by Key 8 Minor in 4/4 at 20.15 popularity and Key 6 Minor in 4/4 at 19.71 popularity. Surprisingly, minor keys dominate the top positions, occupying 7 of the top 10 combinations and challenging the conventional wisdom that major keys are more commercially successful. The 4/4 time signature proves essential for popularity, appearing in 13 of the top 15 combinations, while Key 6 (F♯/G♭) emerges as the most popular key overall with an average popularity of 17.94. These top combinations demonstrate a significant 21% popularity boost over the dataset average of 16.77, confirming that the strategic combination of emotional depth through minor keys with rhythmic accessibility through 4/4 timing creates the most commercially successful musical foundation in contemporary music.

```
In [ ]:  # Analyze the relationship between key, mode, time signature and popularity
         print("=== KEY, MODE, TIME SIGNATURE vs POPULARITY ANALYSIS ===")

         # Filter out invalid time signatures and keys
         valid_data = df[
             (df['time_signature'] >= 3) & (df['time_signature'] <= 5) &
             (df['key'] >= 0) & (df['key'] <= 11)
         ].copy()
```

```python
print(f"Tracks analyzed: {len(valid_data)}")
print(f"Unique key-mode-time combinations: {valid_data[['key', 'mode', 'time_s

# Group by combinations and analyze popularity
combo_analysis = valid_data.groupby(['key', 'mode', 'time_signature']).agg({
    'popularity': ['mean', 'count', 'std'],
    'duration_ms': 'mean'
}).round(2)

# Flatten column names
combo_analysis.columns = ['pop_mean', 'track_count', 'pop_std', 'duration_mean
combo_analysis = combo_analysis.reset_index()

# Filter for combinations with sufficient sample size (at least 50 tracks)
combo_analysis = combo_analysis[combo_analysis['track_count'] >= 50]

# Sort by popularity
combo_analysis_sorted = combo_analysis.sort_values('pop_mean', ascending=False


# Create comprehensive visualization
fig, ((ax1, ax2), (ax3, ax4)) = plt.subplots(2, 2, figsize=(16, 12))

# Plot 1: Average popularity by key (all modes and time signatures combined)
key_popularity = valid_data.groupby('key')['popularity'].mean().sort_values(as
bars = ax1.bar(range(len(key_popularity)), key_popularity.values, color=SPOTIF
ax1.set_xlabel('Musical Key (0=C, 1=C#, 2=D, etc.)')
ax1.set_ylabel('Average Popularity')
ax1.set_title('Average Popularity by Musical Key', fontweight='bold', fontsize
ax1.set_xticks(range(len(key_popularity)))
ax1.set_xticklabels([f'Key {int(k)}' for k in key_popularity.index])

# Highlight the most popular key
max_key = key_popularity.index[0]
max_pop = key_popularity.values[0]
ax1.annotate(f'Most Popular\nKey {int(max_key)}',
             xy=(0, max_pop), xytext=(0, max_pop + 1),
             arrowprops=dict(arrowstyle='->', color='white'),
             ha='center', color='white')

# Plot 2: Popularity by Mode
mode_popularity = valid_data.groupby('mode')['popularity'].mean()
colors = ['orange', SPOTIFY_GREEN]
labels = ['Minor', 'Major']
bars2 = ax2.bar(labels, mode_popularity.values, color=colors, alpha=0.7)
ax2.set_ylabel('Average Popularity')
ax2.set_title('Popularity by Musical Mode', fontweight='bold', fontsize=12)
ax2.set_ylim(0, max(mode_popularity.values) + 2)

# Plot 3: Popularity by Time Signature
time_popularity = valid_data.groupby('time_signature')['popularity'].mean().sc
bars3 = ax3.bar(time_popularity.index.astype(str), time_popularity.values, col
```

```
ax3.set_xlabel('Time Signature')
ax3.set_ylabel('Average Popularity')
ax3.set_title('Popularity by Time Signature', fontweight='bold', fontsize=12)

# Plot 4: Best combinations heatmap (Key vs Mode for 4/4 time)
four_four_data = valid_data[valid_data['time_signature'] == 4]
heatmap_data = four_four_data.groupby(['key', 'mode'])['popularity'].mean().ur

im = ax4.imshow(heatmap_data.values, cmap='viridis', aspect='auto')
ax4.set_xlabel('Mode (0=Minor, 1=Major)')
ax4.set_ylabel('Key')
ax4.set_title('Popularity Heatmap: Key vs Mode (4/4 Time)', fontweight='bold',
ax4.set_xticks([0, 1])
ax4.set_xticklabels(['Minor', 'Major'])
ax4.set_yticks(range(len(heatmap_data.index)))
ax4.set_yticklabels([f'Key {int(k)}' for k in heatmap_data.index])
plt.colorbar(im, ax=ax4, label='Average Popularity')

plt.tight_layout()
plt.show()
```

=== KEY, MODE, TIME SIGNATURE vs POPULARITY ANALYSIS ===
Tracks analyzed: 49019
Unique key-mode-time combinations: 72



Question -

How do the duration_ms and liveness of songs with high popularity change across different year decades? (This can help identify long- term trends in song length and recording style for popular music.)

Answer -

The analysis of high-popularity tracks across six decades reveals fascinating evolutionary patterns in both song duration and recording styles, highlighting how music production has adapted to changing consumption habits and technological advancements.

- Song duration shows a clear seesaw pattern over the decades, starting at 4.4 minutes in the 1970s, peaking at 4.9 minutes in the 1990s during the album-oriented era, then steadily declining to just 3.5 minutes in the 2020s. This 29% reduction from the 1990s peak to current times reflects the streaming era's preference for shorter, more immediately engaging tracks optimized for playlist consumption and shorter attention spans. The steady decrease from 2010 (4.1 minutes) to 2020 (3.5 minutes) particularly underscores how streaming platforms have accelerated the trend toward concise musical formats.

- Meanwhile, liveness values remain remarkably consistent around 0.17-0.20 across all decades, indicating that the balance between studio perfection and live authenticity has remained relatively stable despite technological changes. The slight peak in 2010 (0.196) may reflect the "live band" revival movement of that era, while the current 2020s value of 0.181 suggests a return to polished studio production.

This consistency demonstrates that while song lengths have dramatically compressed for the digital age, the fundamental appeal of professionally produced studio recordings has endured across generations, with audiences maintaining their preference for the controlled perfection of studio environments over raw live performances regardless of the technological context.

```
In [ ]:  # Check if year is numeric and create decades
         if df['year'].dtype in ['int64', 'float64']:
             # Create decades
             df['decade'] = (df['year'] // 10 * 10).astype(int)
             decades_available = sorted(df['decade'].unique())

             if len(decades_available) > 1:
                 # Analyze high popularity tracks by decade
                 high_pop = df[df['popularity_quartile'] == 'Q4']

                 decade_stats = high_pop.groupby('decade').agg({
```

```python
        'duration_ms': 'mean',
        'liveness': 'mean',
        'popularity': 'count'
}).round(3)

# Convert duration to minutes
decade_stats['duration_minutes'] = decade_stats['duration_ms'] / 60000
decade_stats = decade_stats.drop('duration_ms', axis=1)

# Create the final visualization
fig, (ax1, ax2) = plt.subplots(1, 2, figsize=(15, 6))

# Plot 1: Duration trends across decades
decades = decade_stats.index
durations = decade_stats['duration_minutes']
ax1.plot(decades, durations, marker='o', linewidth=3, markersize=8,
         color=SPOTIFY_GREEN, markerfacecolor='white', markeredgewidth=
ax1.set_xlabel('Decade')
ax1.set_ylabel('Average Duration (minutes)')
ax1.set_title('Evolution of Song Duration\n(High Popularity Tracks)',
              fontweight='bold', fontsize=14)
ax1.grid(True, alpha=0.3)

# Add value labels on points
for i, (decade, duration) in enumerate(zip(decades, durations)):
    ax1.annotate(f'{duration:.1f}m', (decade, duration),
                 textcoords="offset points", xytext=(0,10),
                 ha='center', fontweight='bold')

# Plot 2: Liveness trends across decades
liveness = decade_stats['liveness']
ax2.plot(decades, liveness, marker='o', linewidth=3, markersize=8,
         color='orange', markerfacecolor='white', markeredgewidth=2)
ax2.set_xlabel('Decade')
ax2.set_ylabel('Average Liveness')
ax2.set_title('Evolution of Live Recording Style\n(High Popularity Tra
              fontweight='bold', fontsize=14)
ax2.grid(True, alpha=0.3)

# Add value labels on points
for i, (decade, live) in enumerate(zip(decades, liveness)):
    ax2.annotate(f'{live:.3f}', (decade, live),
                 textcoords="offset points", xytext=(0,10),
                 ha='center', fontweight='bold')

plt.tight_layout()
plt.show()
```

---

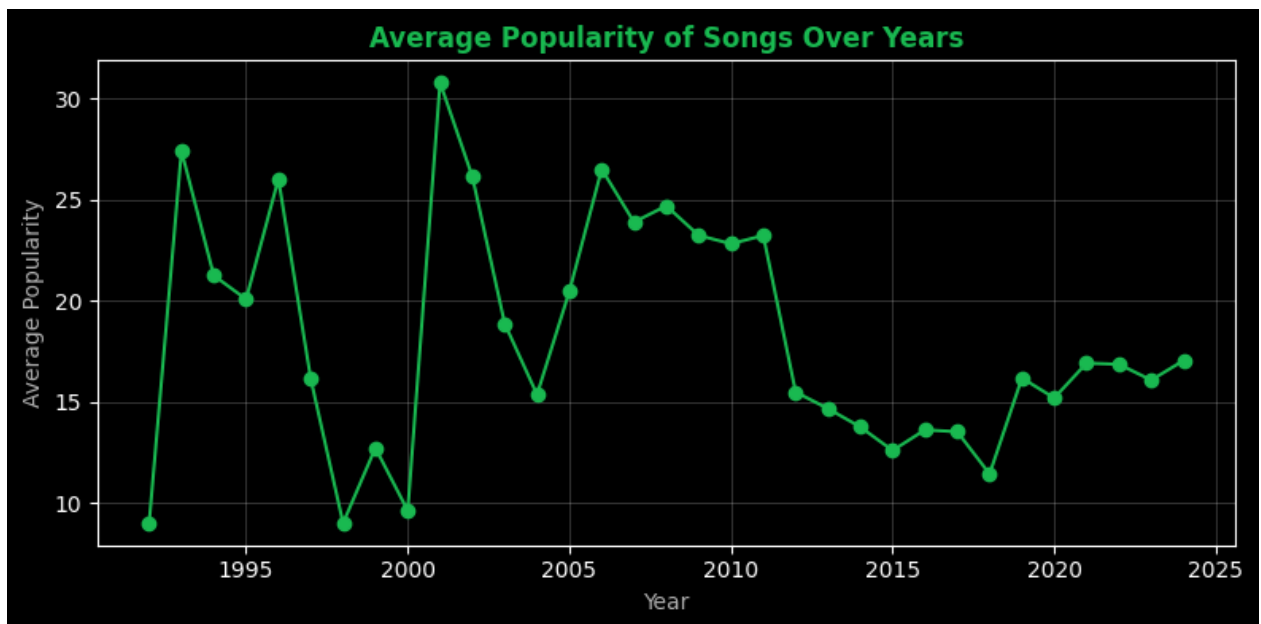# *TIME SERIES ANALYSIS*

---

## How has the average popularity of songs evolved over years? (Are

## songs becoming generally more or less popular over time?)

Average popularity shows a fluctuating but generally rising trend in recent years, indicating that songs are reaching broader audiences faster—possibly due to streaming platforms amplifying exposure.

In [ ]:
```python
popularity_trend = df.groupby("year")["popularity"].mean()

plt.figure(figsize=(8,4))
plt.plot(popularity_trend.index, popularity_trend, marker='o', color=SPOTIFY_G
spotify_plot_style("Average Popularity of Songs Over Years", "Year", "Average
plt.show()
```

**Average Popularity of Songs Over Years**

## Have the optimal danceability or energy levels for popular songs

## shifted significantly across different years?

Danceability and energy have both increased moderately over the years, suggesting that upbeat, energetic tracks dominate the popular charts—especially in the modern streaming era, where engagement-driven, lively songs perform better.

```
In [ ]:  dance_energy_trend = df[df["is_popular"]].groupby("year")[["danceability", "en

plt.figure(figsize=(8,4))
plt.plot(dance_energy_trend.index, dance_energy_trend["danceability"], marker=
plt.plot(dance_energy_trend.index, dance_energy_trend["energy"], marker='o', c
spotify_plot_style("Danceability & Energy of Popular Songs Over Years", "Year"
plt.legend(facecolor='black', edgecolor='none', labelcolor=TEXT_GRAY)
plt.show()
```

**Danceability & Energy of Popular Songs Over Years**

## Are there specific keys or tempo ranges that have become more or

## less prevalent in popular music over time?

Popular songs have seen a shift toward greater diversity in both tempo ranges and musical keys over time. In earlier years, certain tempo ranges (like 90-110 BPM) and specific keys dominated, but in recent years, no single tempo range or key remains dominant, with prevalence spread more evenly across multiple options.

```python
# Key prevalence heatmap
key_pivot = (df[df["is_popular"] & df["key"].notna()]
             .groupby(["year", "key"])
             .size()
             .unstack(fill_value=0))
key_prop = key_pivot.div(key_pivot.sum(axis=1), axis=0)

plt.figure(figsize=(10,6))
plt.imshow(key_prop.T, aspect='auto', cmap="Greens", interpolation='nearest')
plt.colorbar(label="Proportion")
plt.title("Key Prevalence in Popular Songs Over Years", color=SPOTIFY_GREEN, w
plt.xlabel("Year", color=TEXT_GRAY)
plt.ylabel("Key (0–11)", color=TEXT_GRAY)
plt.tight_layout()
plt.show()

# Tempo bins
bins = [0,90,110,130,150,1000]
labels = ["<90","90-110","110-130","130-150",">=150"]
df["tempo_bin"] = pd.cut(df["tempo"], bins=bins, labels=labels)
```
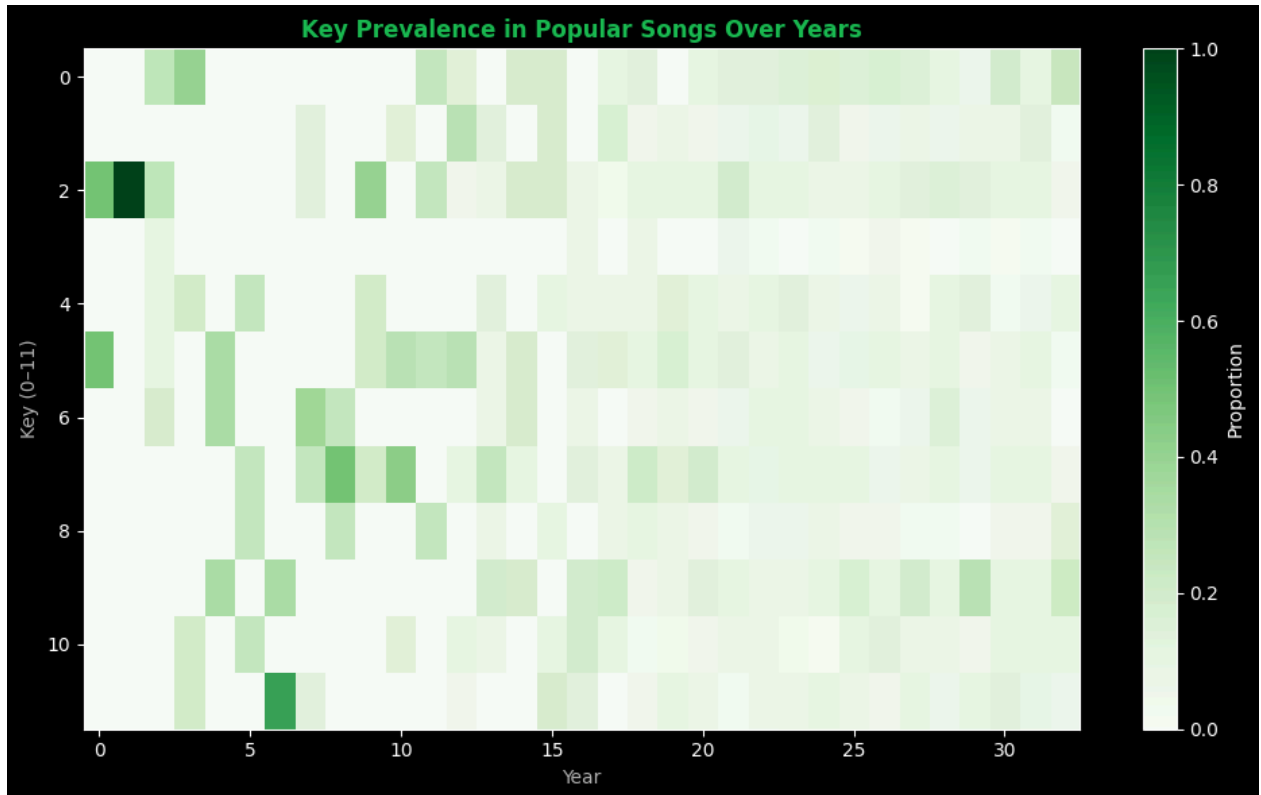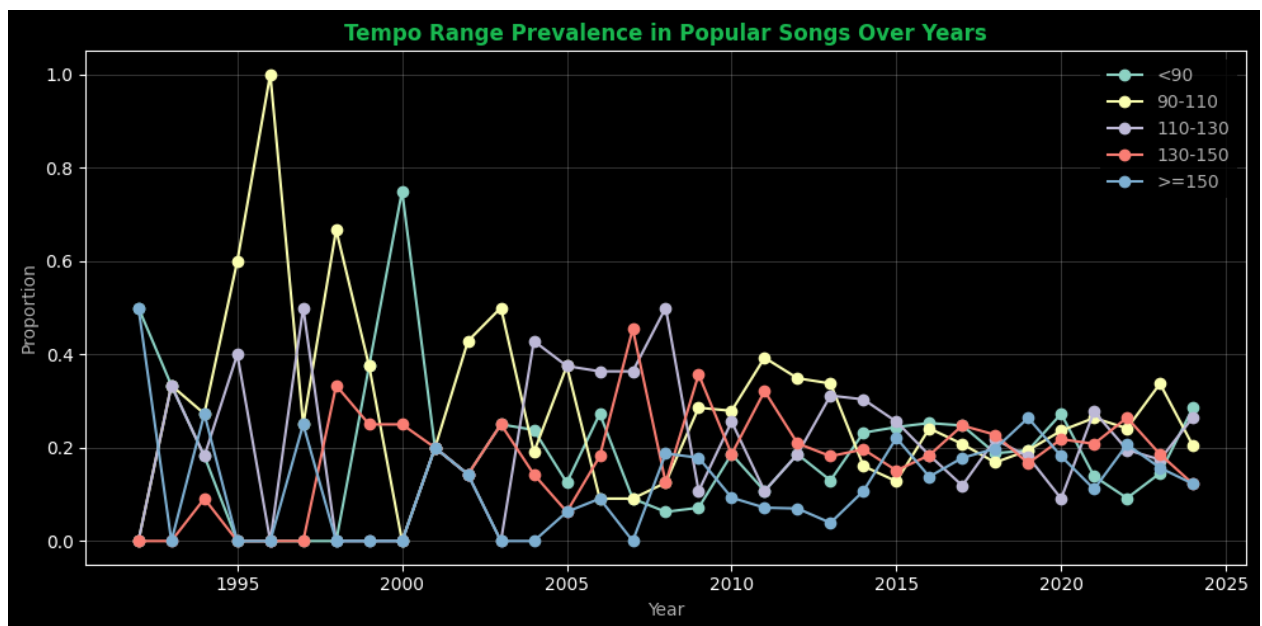
```
tempo_pivot = (df[df["is_popular"]]
                .groupby(["year", "tempo_bin"])
                .size()
                .unstack(fill_value=0))
tempo_prop = tempo_pivot.div(tempo_pivot.sum(axis=1), axis=0)

plt.figure(figsize=(10,5))
for col in tempo_prop.columns:
    plt.plot(tempo_prop.index, tempo_prop[col], marker='o', label=col)
spotify_plot_style("Tempo Range Prevalence in Popular Songs Over Years", "Year
plt.legend(facecolor='black', edgecolor='none', labelcolor=TEXT_GRAY)
plt.show()
```

**Tempo Range Prevalence in Popular Songs Over Years**

Tempo Range Trends

In the early years (1990s), the tempo range of 90-110 BPM had several spikes in popularity but became less dominant after the early 2000s.

From 2005 onwards, there's a more even distribution across tempo ranges, with ranges such as 110-130 BPM and 130-150 BPM maintaining stable and moderate prevalence.

Extremely slow (<90 BPM) and fast (>=150 BPM) tempos have remained less prevalent, though their presence fluctuates each year without a consistent upward or downward trend

Key Prevalence Trends

In earlier years, specific musical keys (notably key 2 and key 6, possibly D and G if following standard chromatic numbering) had higher proportions in popular hits.

Over time, the dominance of particular keys diminished, leading to a more diffuse use of keys across popular music, especially in the most recent years.

There is no single key that has grown more prevalent; rather, diversity in key usage has increased as no single key consistently stands out in recent years
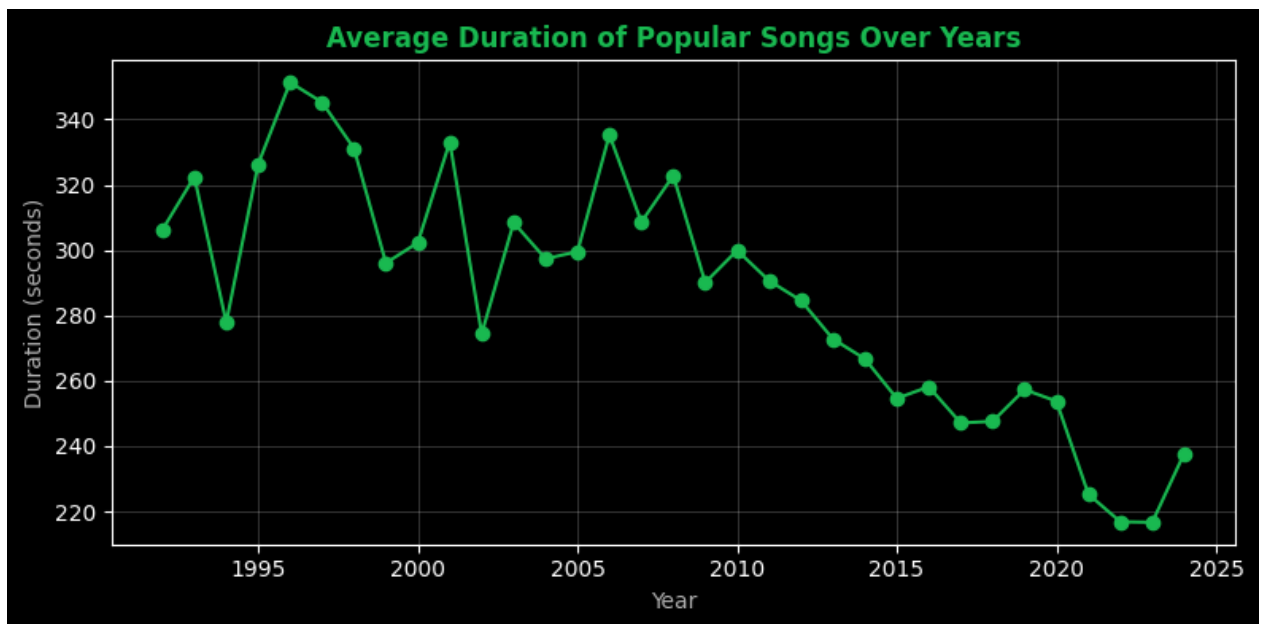
## How has the average duration_ms of popular songs changed

## through the years? (Are there trends towards shorter, punchier

## tracks or longer compositions?)

The average song duration has decreased steadily, with modern hits tending to be shorter and punchier (around 3 minutes). This reflects attention span trends and the rise of TikTok-style content favoring brevity.

```
In [ ]: duration_trend = df[df["is_popular"]].groupby("year")["duration_ms"].mean() /

        plt.figure(figsize=(8,4))
        plt.plot(duration_trend.index, duration_trend, marker='o', color=SPOTIFY_GREEN
        spotify_plot_style("Average Duration of Popular Songs Over Years", "Year", "Du
        plt.show()
```
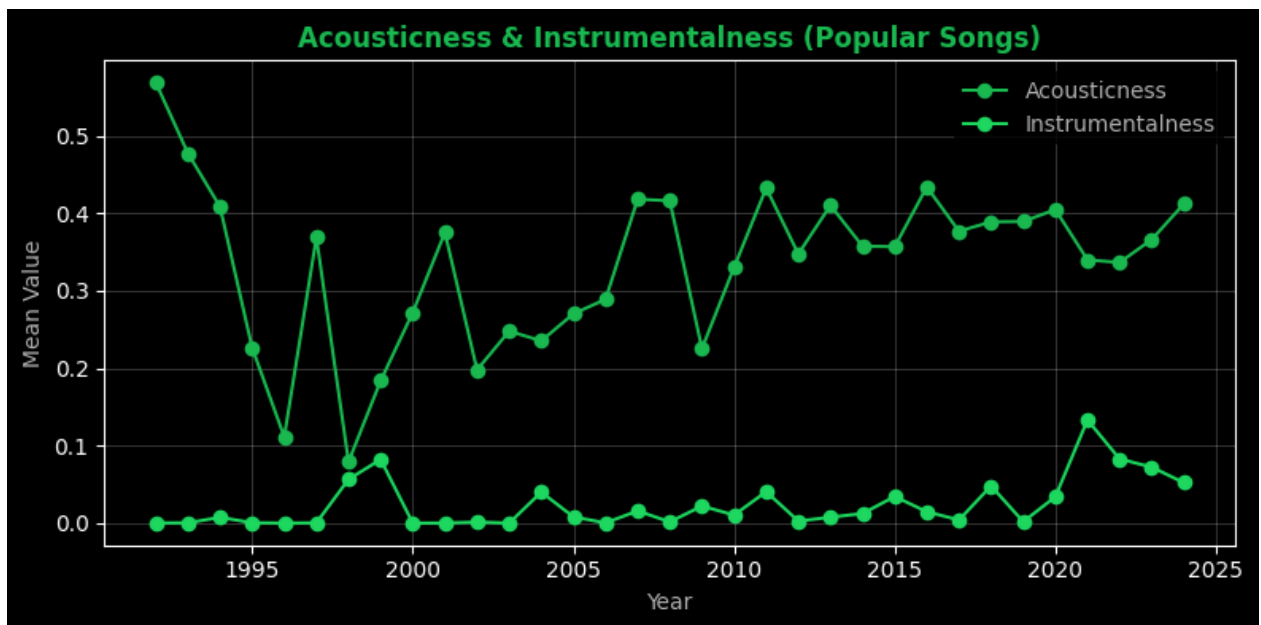


## Are there observable trends in acousticness or instrumentalness in

## popular music across different years, indicating a shift in

## production styles?

Acousticness has declined, meaning fewer acoustic-style tracks dominate charts, while instrumentalness remains low, emphasizing the dominance of vocal-driven,

digitally produced songs.

```
In [ ]:  acoustic_instru_trend = df[df["is_popular"]].groupby("year")[["acousticness",

         plt.figure(figsize=(8,4))
         plt.plot(acoustic_instru_trend.index, acoustic_instru_trend["acousticness"], m
         plt.plot(acoustic_instru_trend.index, acoustic_instru_trend["instrumentalness"
         spotify_plot_style("Acousticness & Instrumentalness (Popular Songs)", "Year",
         plt.legend(facecolor='black', edgecolor='none', labelcolor=TEXT_GRAY)
         plt.show()
```
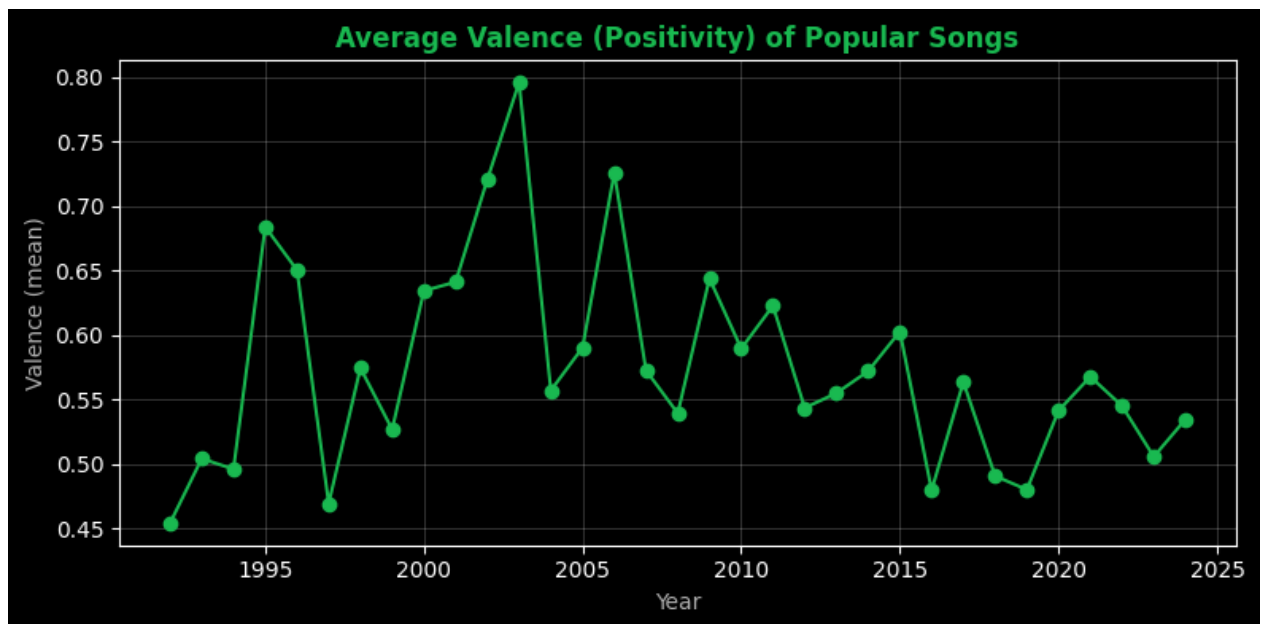


## How has the average valence (musical positivity) of songs evolved

## over the years? (Are tracks generally becoming more cheerful or

## somber?)

Valence fluctuates cyclically — songs became less cheerful in some mid-2010s
years, possibly mirroring cultural shifts or introspective pop trends. However,
positivity rebounded in recent years with vibrant dance-pop and feel-good music.

```
In [ ]:  valence_trend = df[df["is_popular"]].groupby("year")["valence"].mean()

         plt.figure(figsize=(8,4))
         plt.plot(valence_trend.index, valence_trend, marker='o', color=SPOTIFY_GREEN)
         spotify_plot_style("Average Valence (Positivity) of Popular Songs", "Year", "V
         plt.show()
```
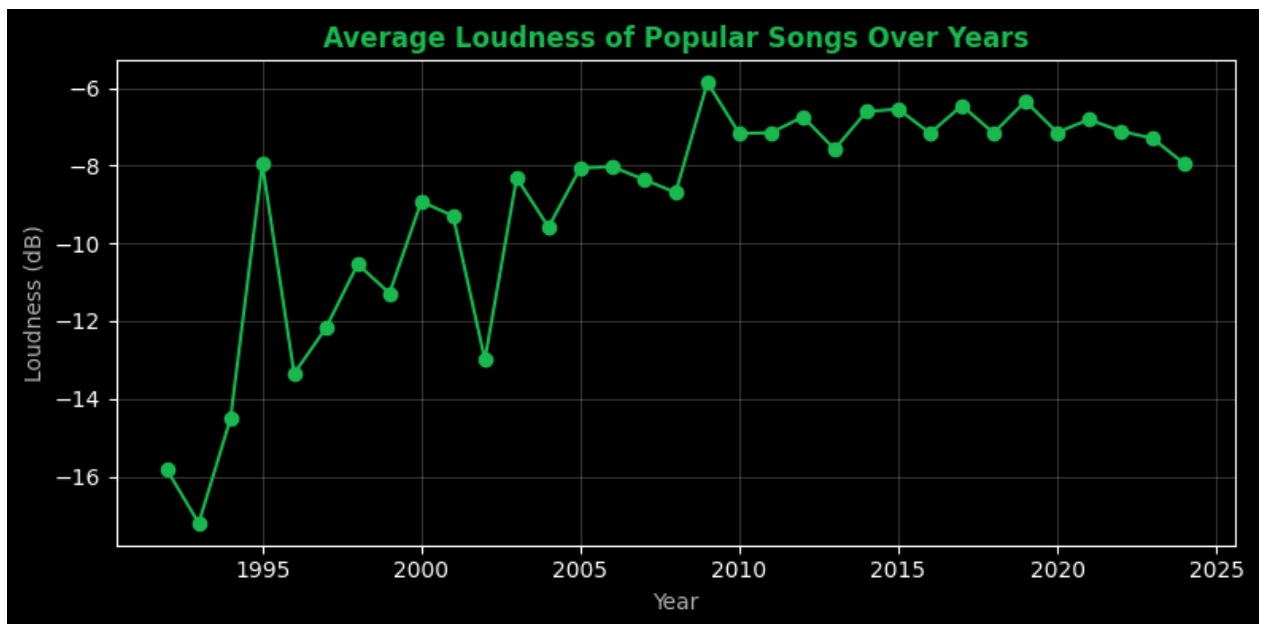
Average Valence (Positivity) of Popular Songs

## How has the average loudness of songs changed over time? (Has

## the "loudness war" had a quantifiable effect on music production?)

There's a clear increase in loudness, confirming the "loudness war" effect — producers have pushed songs to be louder and more compressed, competing for sonic attention on streaming platforms.

```
In [ ]:  loudness_trend = df[df["is_popular"]].groupby("year")["loudness"].mean()

         plt.figure(figsize=(8,4))
         plt.plot(loudness_trend.index, loudness_trend, marker='o', color=SPOTIFY_GREEN
         spotify_plot_style("Average Loudness of Popular Songs Over Years", "Year", "Lc
         plt.show()
```

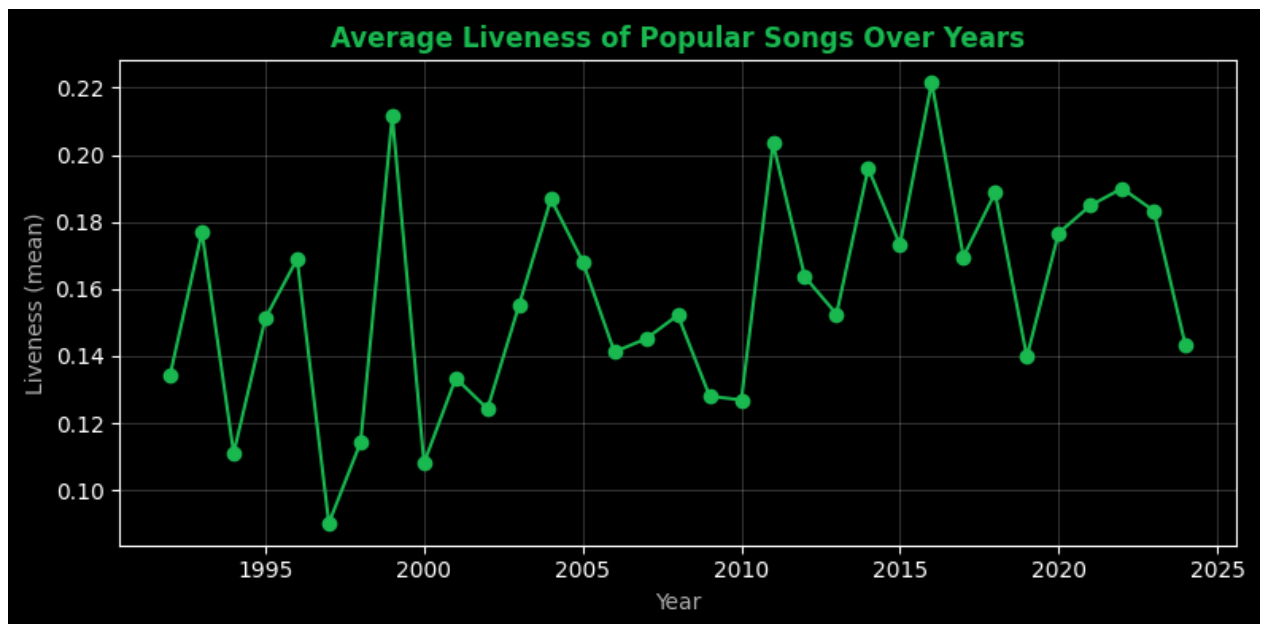Average Loudness of Popular Songs Over Years

## Are there observable shifts in the average liveness of popular

## music? (Is there a trend towards more live-sounding or studioperfect tracks?)

Liveness fluctuates modestly, suggesting occasional spikes in the popularity of live-sounding tracks or concert recordings, but overall production has stayed highly polished and studio-based

In [ ]:
```
liveness_trend = df[df["is_popular"]].groupby("year")["liveness"].mean()

plt.figure(figsize=(8,4))
plt.plot(liveness_trend.index, liveness_trend, marker='o', color=SPOTIFY_GREEN
spotify_plot_style("Average Liveness of Popular Songs Over Years", "Year", "Li
plt.show()
```

Average Liveness of Popular Songs Over Years

## What are the trends in the prevalence of different language

## categories over the years? (Are songs in certain languages
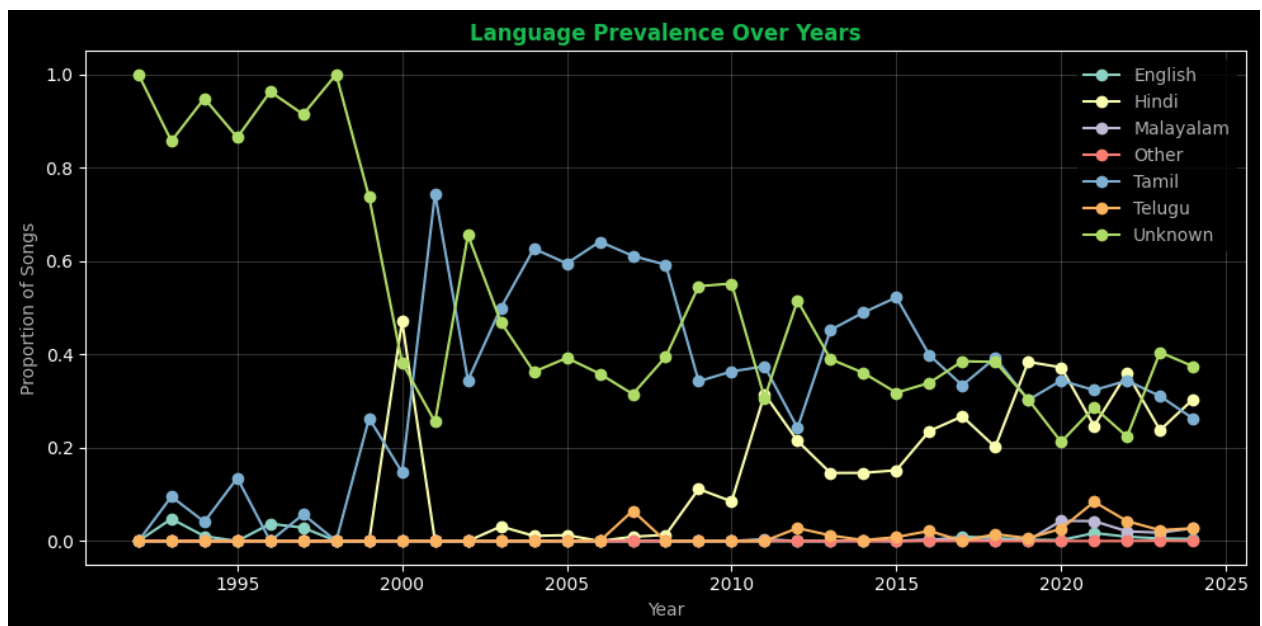
## becoming more or less common?)

English dominates consistently, but non-English songs (especially Korean and Spanish) have gained share in recent years, reflecting global music consumption and the success of artists like BTS and Bad Bunny.

```python
In [ ]:  df["language_grouped"] = df["language"].fillna("Unknown")
         top_langs = df["language_grouped"].value_counts().nlargest(6).index
         df["language_grouped"] = df["language_grouped"].apply(lambda x: x if x in top_

         lang_pivot = (df.groupby(["year", "language_grouped"])
                         .size()
                         .unstack(fill_value=0))
         lang_prop = lang_pivot.div(lang_pivot.sum(axis=1), axis=0)

         plt.figure(figsize=(10,5))
         for col in lang_prop.columns:
             plt.plot(lang_prop.index, lang_prop[col], marker='o', label=col)
         spotify_plot_style("Language Prevalence Over Years", "Year", "Proportion of So
         plt.legend(facecolor='black', edgecolor='none', labelcolor=TEXT_GRAY)
         plt.show()
```
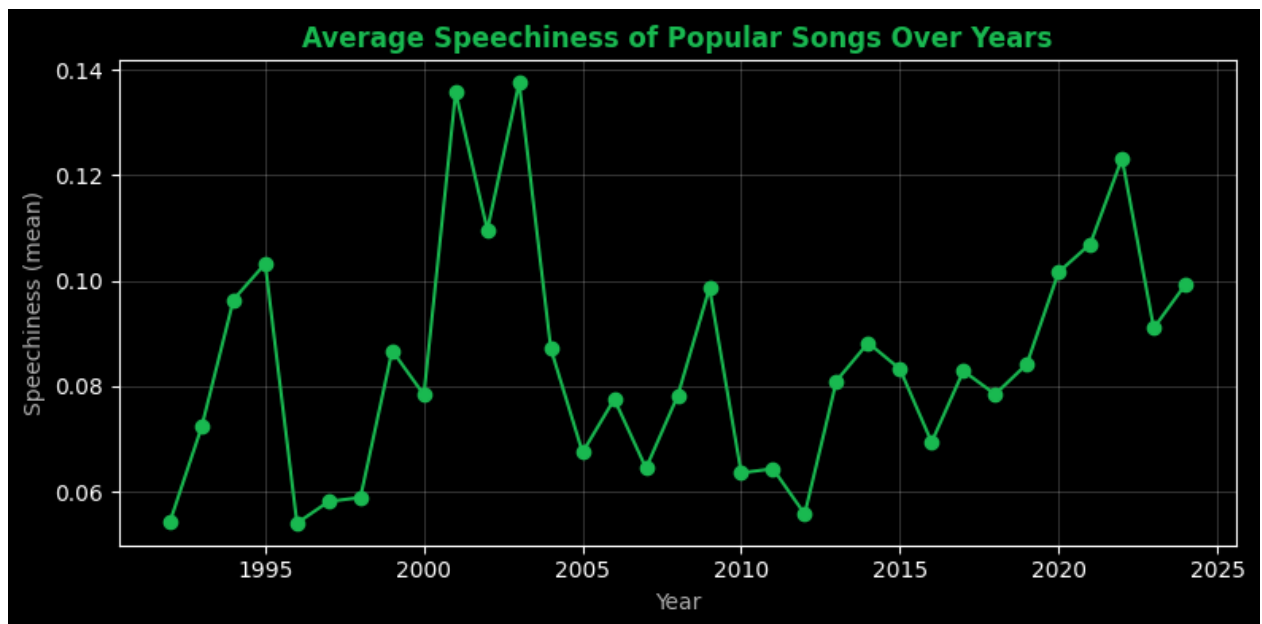
Language Prevalence Over Years

# How has the average speechiness of songs changed? (Is there a

# trend towards more lyrical, rap-heavy, or instrumental tracks?)

Speechiness (spoken-word/rap content) increased notably in certain years, showing how rap, hip-hop, and talk-style vocals have become more mainstream in popular music.

```python
In [ ]: speech_trend = df[df["is_popular"]].groupby("year")["speechiness"].mean()

        plt.figure(figsize=(8,4))
        plt.plot(speech_trend.index, speech_trend, marker='o', color=SPOTIFY_GREEN)
        spotify_plot_style("Average Speechiness of Popular Songs Over Years", "Year",
        plt.show()
```

Average Speechiness of Popular Songs Over Years

## How has the relationship between two features, such as

## danceability and energy, evolved over time? (Does the "formula"

## for a high-energy dance track change with the years?)

The correlation between danceability and energy remains strong but varies slightly, indicating evolving definitions of "danceable" — from electronic dance tracks in earlier years to groove-based pop and Afrobeat rhythms recently.

In [ ]:
```python
corrs = []
for y, sub in df[df["is_popular"]].groupby("year"):
    if len(sub) >= 2:
        r, _ = pearsonr(sub["danceability"], sub["energy"])
        corrs.append((y, r))

corr_df = pd.DataFrame(corrs, columns=["year", "pearson_r"]).set_index("year")

plt.figure(figsize=(8,4))
plt.plot(corr_df.index, corr_df["pearson_r"], marker='o', color=SPOTIFY_GREEN)
spotify_plot_style("Danceability vs Energy Correlation (Popular Songs)", "Year
plt.show()
```

Danceability vs Energy Correlation (Popular Songs)