

# Data Description

The dataset includes detailed audio features and metadata for songs, which are instrumental in analyzing and predicting song popularity. Each feature has been standardized and described based on Spotify's Web API and other sources. The key attributes represent both audio characteristics and external metadata, allowing for a comprehensive understanding of song attributes that contribute to their popularity. Below is a summary of the key features in the dataset:

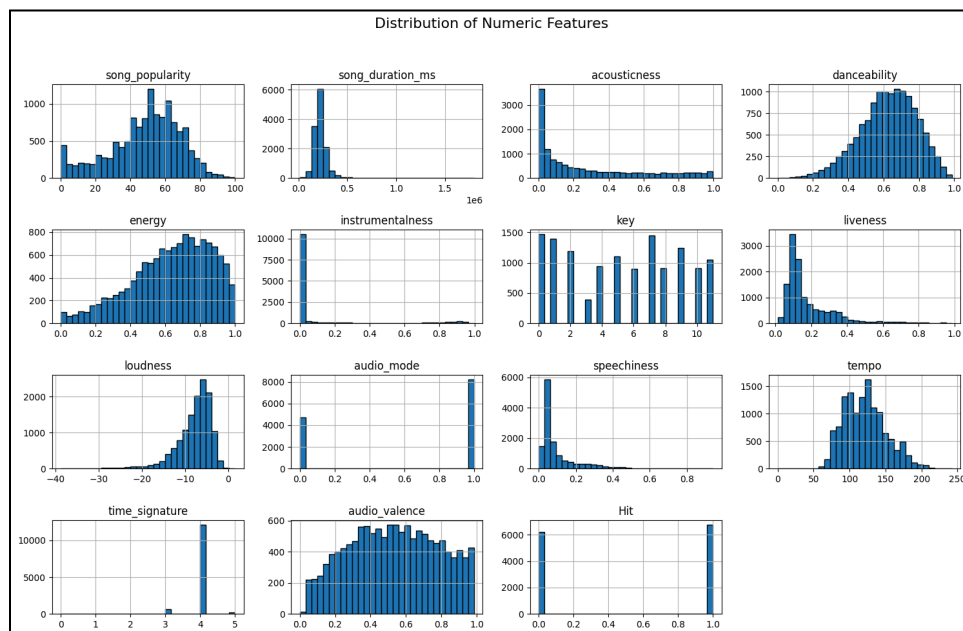
Feature	Description
Acousticness	A confidence measure (0.0 to 1.0) indicates whether the track is acoustic. Higher values indicate higher confidence in the track's acoustic nature.
Danceability	Measures how suitable a track is for dancing based on tempo, rhythm stability, beat strength, and regularity. Ranges from 0.0 (least danceable) to 1.0 (most danceable).
Duration (ms)	The duration of the track is in milliseconds.
Energy	A perceptual measure (0.0 to 1.0) represents a track's intensity and activity. High-energy tracks often feel fast and loud (e.g., rock/metal), while classical tracks score lower.
Instrumentalness	Predicts whether a track contains no vocals. Values closer to 1.0 indicate a higher likelihood of being instrumental (e.g., classical, soundtracks).
Liveness	Detects the presence of an audience in the recording. Higher values (e.g., above 0.8) indicate the track was likely performed live.

<b>Loudness</b>	The overall loudness of a track is measured in decibels (dB). Values range from -60 dB to 0 dB, with higher values indicating louder tracks.
<b>Speechiness</b>	Measures the presence of spoken words in a track. Values closer to 1.0 indicate higher speech content (e.g., podcasts or spoken word tracks).
<b>Tempo</b>	The estimated speed of the track in beats per minute (BPM).
<b>Valence</b>	A measure (0.0 to 1.0) describing the positivity of the track's mood. Higher values indicate happier tracks, while lower values suggest sadness or negativity.
<b>Song Popularity</b>	Represents the popularity of the song on a scale from 0 to 100, with higher values indicating greater popularity.
<b>Key</b>	Defines the tonal center of the song, influencing harmony and mood.
<b>Audio Mode</b>	Indicates whether the song has a major or minor tonality, crucial for emotional context.
<b>tine_signature</b>	Describes the rhythmic structure, essential for analyzing beat patterns and suitability for dancing or other activities.

# Exploratory Data Analysis

The exploratory data analysis (EDA) revealed key insights into the dataset, highlighting the distribution of numeric features, correlations between variables, and characteristics of hit versus non-hit songs. The analysis included visualizations such as histograms, box plots, and correlation matrices to explore relationships and trends in the data.

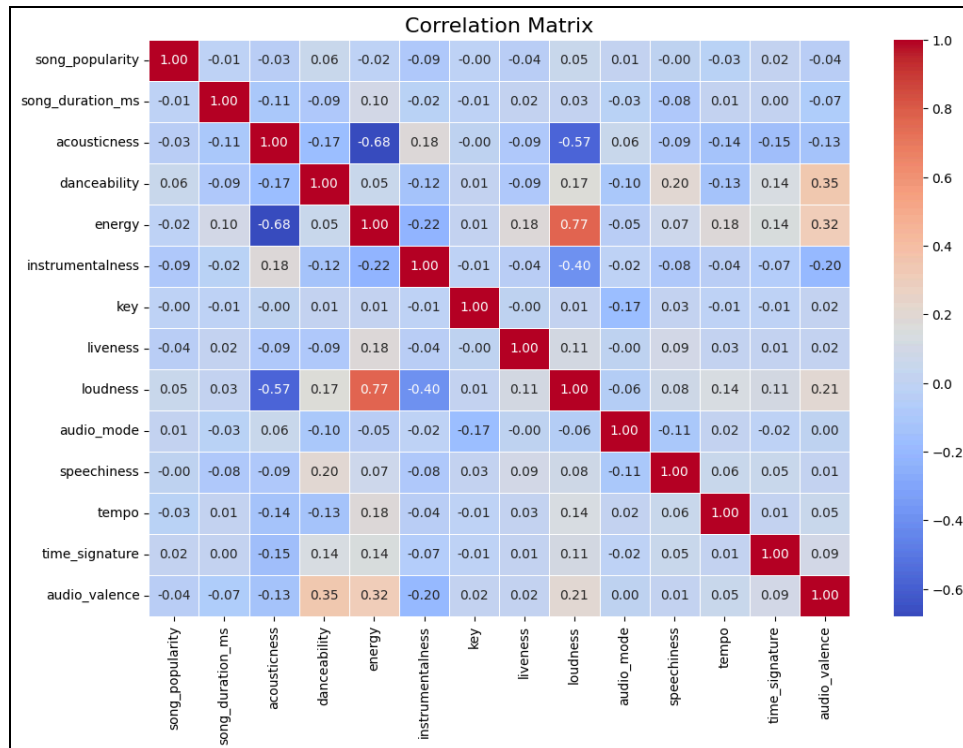
## Distribution of Numeric Features:



**Observation:** Most features like **danceability**, **energy**, and **tempo** are distributed around typical values (e.g., danceability peaks around 0.6, and tempo centers around 120–130 BPM). Features like **instrumentalness** and **liveness** are skewed, with most songs having low values, indicating they are likely studio-recorded and contain vocals.

**Insight:** The dataset contains a variety of songs, with many showing characteristics typical of popular, high-energy music.

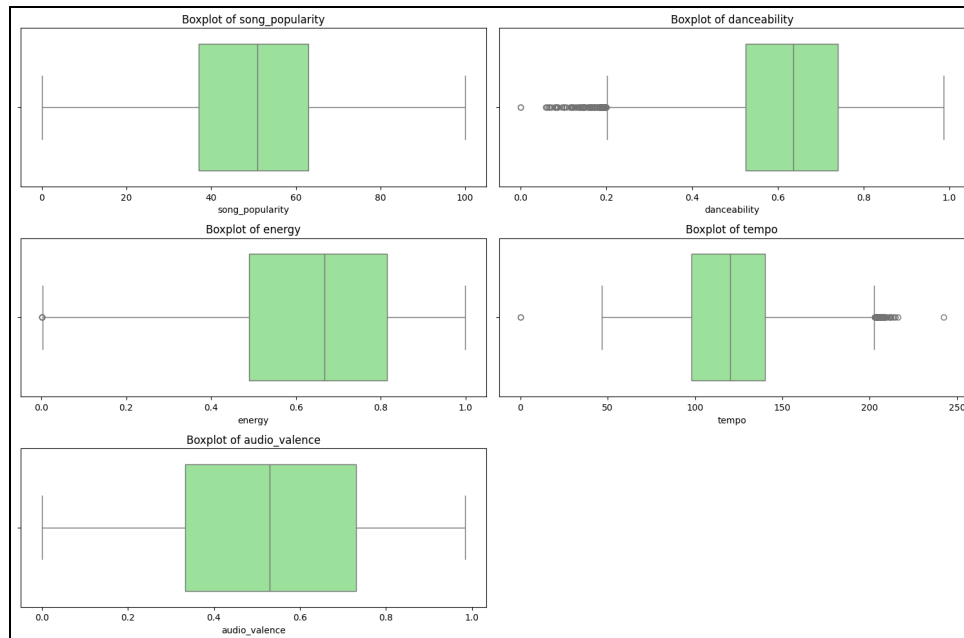
## Correlation Matrix:



**Observation:** **Energy** and **loudness** show a strong positive correlation (+0.77), indicating that louder songs tend to have higher energy. Negative correlations are observed between **acousticness** and **energy** (-0.68), suggesting that acoustic songs are less energetic.

**Insight:** Understanding these correlations helps identify interdependencies between features and informs feature selection for predictive models.

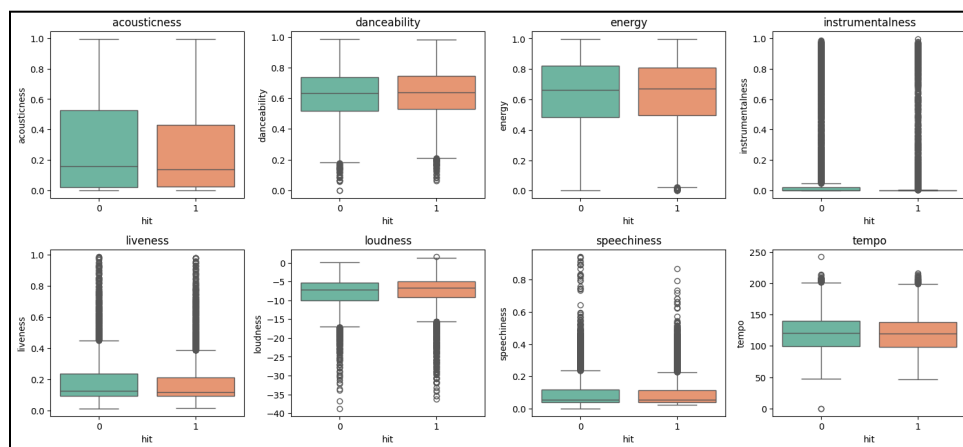
### Box Plots of Numeric Features:



**Observation:** Features like **danceability**, **energy**, and **valence** show consistent median values with minimal outliers, indicating stable distributions. Tempo displays a broader range, with some high outliers.

**Insight:** Key features contributing to song popularity are evenly distributed, with variability in tempo suggesting diverse song pacing.

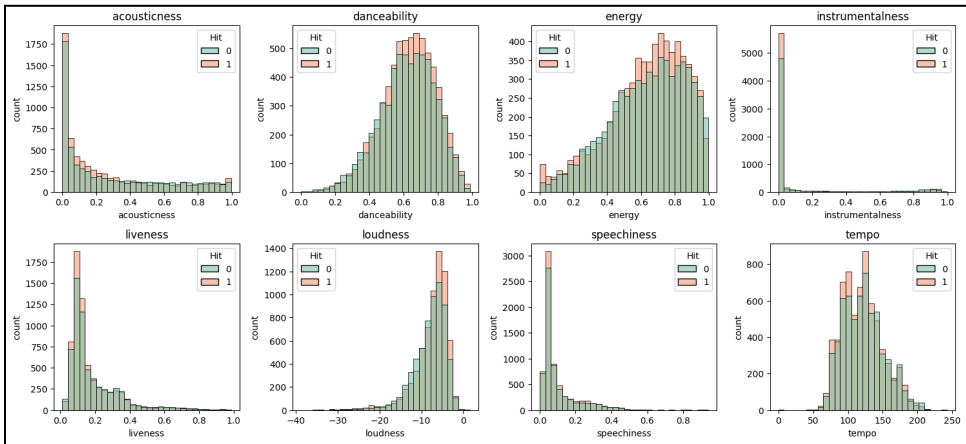
### Box Plots of Hit vs. Non-Hit Songs:



**Observation:** Hit songs have higher median values for **danceability**, **energy**, and **loudness**, while non-hit songs tend to have higher **acousticness**. Hits also show slightly lower **speechiness** values.

**Insight:** Popular songs exhibit characteristics like being highly energetic, danceable, and loud, while acoustic tracks are less likely to be hit.

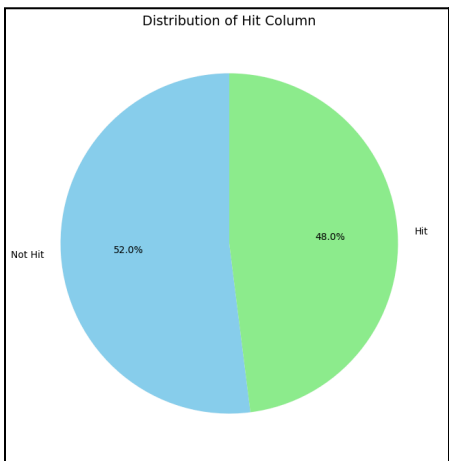
**Histograms for Hit vs. Non-Hit Songs:**



**Observation:** Hit songs are more concentrated in mid to high ranges for **danceability**, **energy**, and **loudness**, whereas non-hit songs dominate the lower ranges for these features. Both categories show similar tempo distributions.

**Insight:** Audio features like energy and loudness play a crucial role in distinguishing hit songs from non-hits, while tempo has less impact.

**Pie Chart of Hit Distribution:**



**Observation:** The dataset is relatively balanced, with 48% hit songs and 52% non-hit songs.

**Insight:** The balanced dataset ensures that machine learning models trained on this data will not be overly biased toward one category.

# Methodology

The methodology employed in this project focuses on analyzing audio features to identify patterns that distinguish hit songs from non-hit songs and to predict song popularity using regression models. The process included data preprocessing, exploratory data analysis (EDA), feature selection, and the application of machine learning models. Here is a detailed summary of the methodology:

## Data Preprocessing

- **Dataset Source:**
  - The dataset, obtained from Kaggle, consists of 15 audio features for over 12,000 songs.
- **Handling Missing Data:**
  - The dataset was inspected for missing or null values, which were addressed by either imputation or removal, ensuring data integrity.
- **Hit Classification:**
  - A new column, **Hit**, was created to classify songs based on their popularity. Songs with a **popularity** score above 50 were labeled as hits (**1**), and those below 50 were labeled as non-hits (**0**).
- **Feature Scaling:**
  - Features with different ranges (e.g., **duration\_ms**, **loudness**) were standardized using **StandardScaler** to normalize their distributions and prevent model bias toward features with larger scales.

## Model Implementation

Two regression models, **Multilinear Regression** and **Polynomial Regression**, were implemented to analyze and predict song popularity.

## A. Multilinear Regression

- **Objective:**
  - Capture linear relationships between features and song popularity.
- **Implementation:**
  - The model was trained using normalized features to ensure fair contributions.
  - Coefficients were examined to understand the influence of individual features on popularity.
- **Evaluation:**
  - Training and testing performance was evaluated using Mean Squared Error (MSE) and R-squared ( $R^2$ ) metrics.

## B. Polynomial Regression

- **Objective:**
  - Capture nonlinear relationships and interactions between features.
- **Feature Transformation:**
  - Polynomial features up to degree 2 were generated using `PolynomialFeatures` from Scikit-learn.
  - Interaction terms (e.g., `danceability energy`, `energy tempo`) were included to capture the combined effects of features.
- **Model Fitting:**
  - The transformed features were scaled and used to train a linear regression model.
- **Evaluation:**
  - Polynomial regression was evaluated similarly to multilinear regression, with MSE and  $R^2$  as performance metrics.

## Model Validation

- **Train-Test Split:**



- The dataset was split into 80% training and 20% testing to evaluate generalization performance.
- **Cross-Validation:**
  - K-fold cross-validation was used to ensure robustness and reduce the risk of overfitting.
- **Comparison:**
  - The performance of both models was compared:
    - Multilinear regression served as the baseline.
    - Polynomial regression showed improved performance by capturing nonlinear trends.

## Metrics for Evaluation

- **Mean Squared Error (MSE):**
  - Quantified the average error between predicted and actual popularity scores.
- **R-squared ( $R^2$ ):**
  - Measured the proportion of variance explained by the model, indicating its effectiveness.

The methodology combined rigorous preprocessing, insightful EDA, and effective model implementation to analyze song popularity. By leveraging both linear and nonlinear regression models, the approach provided interpretable insights and improved predictive power, forming a robust framework for identifying and understanding hit songs.

## Results

	Multilinear Regression:	Polynomial Regression (Degree 2):
<b>Training Performance:</b>	<ul style="list-style-type: none"><li>• Mean Squared Error (MSE): 0.0938</li><li>• R-squared (<math>R^2</math>): 62.39%</li></ul>	<ul style="list-style-type: none"><li>• MSE: 0.0838</li><li>• <math>R^2</math>: 66.43%</li></ul>
<b>Testing Performance:</b>	<ul style="list-style-type: none"><li>• MSE: 0.0943</li><li>• <math>R^2</math>: 62.22%</li></ul>	<ul style="list-style-type: none"><li>• MSE: 0.0860</li><li>• <math>R^2</math>: 65.57%</li></ul>
<b>Observations:</b>	<p>Multilinear regression explained approximately 62% of the variance in song popularity.</p> <p>Features like <b>speechiness</b>, <b>audio valence</b>, and <b>instrumentalness</b> were identified as key contributors.</p>	<p>Polynomial regression improved model performance by capturing nonlinear relationships and interactions between features.</p> <p>Key nonlinear features included <b>audio_mode<sup>2</sup></b>, <b>song_popularity<sup>2</sup></b>, and interactions like <b>danceability energy</b> and <b>energy tempo</b>.</p>

### Feature Importance:

- **Top Features:**
  - Speechiness, audio valence, and instrumentalness were the most influential features in the multilinear model.
  - In the polynomial model, squared terms and feature interactions like **audio\_mode<sup>2</sup>** and **danceability energy** had a significant impact.

## **Observations and Insights from Results**

Hit songs are characterized by higher danceability, energy, and loudness, while lower acoustiness suggests a preference for digital over acoustic elements. Polynomial Regression revealed the importance of nonlinear effects and feature interactions, such as the combined influence of energy and tempo, in shaping song popularity. Both models performed reliably, with  $R^2$  values above 60%, though Polynomial Regression outperformed Multilinear Regression, demonstrating its ability to capture complex relationships. Key features like speechiness, audio valence, and instrumentality were found to significantly impact song popularity, while consistent performance across datasets highlighted model robustness and generalizability.

## **What Can Be Learned from the Results**

The results highlight that hit songs tend to be highly energetic, danceable, and less acoustic, aligning with listener preferences for engaging tracks. Nonlinear analysis, as captured by Polynomial Regression, provides deeper insights into the complex dynamics of song popularity. These findings offer actionable guidance for the music industry, where artists can optimize features like danceability and energy, and marketers can focus on promoting tracks with hit-like characteristics. Additionally, the study lays a foundation for future research by integrating external factors like social media trends and listener demographics to enhance predictions.

# Conclusion

The analysis of music popularity and trends provided valuable insights into the factors that influence the success of songs and demonstrated the predictive power of machine learning models in identifying potential hits. By leveraging a dataset of audio features and applying exploratory data analysis, regression modeling, and feature importance evaluation, the project achieved its goal of understanding the characteristics of popular music and predicting song popularity.

## Key Takeaways

### 1. Understanding Popular Songs:

- Hit songs are characterized by higher **danceability**, **energy**, and **loudness**, suggesting that engaging and rhythmic tracks resonate most with audiences.
- Lower **acousticness** in hits reflects a preference for modern, synthetic sounds over acoustic instrumentation.
- Emotional attributes like **valence** (happiness) and moderate **speechiness** levels contribute to the success of songs, balancing musical and spoken-word content.

### 2. Importance of Nonlinear Relationships:

- The inclusion of polynomial terms and interaction effects in Polynomial Regression revealed that relationships between audio features and popularity are not strictly linear. For example, the interaction of **danceability** and **energy** enhances a song's appeal, and the quadratic effect of **audio mode** demonstrates the importance of feature transformations in capturing complexity.

### 3. Feature Contributions:

- **Speechiness**, **audio valence**, and **instrumentalness** emerged as the most influential features, emphasizing the role of vocal and emotional characteristics in determining a song's popularity.
- Interaction terms and higher-order features, such as **audio\_mode^2** and **danceability energy**, highlighted the combined impact of multiple features on listener preferences.

#### 4. Model Performance:

- Both Multilinear Regression and Polynomial Regression successfully modeled song popularity, with Polynomial Regression showing superior performance ( $R^2 = 66.43\%$  for training and  $65.57\%$  for testing). The consistent results across training and testing datasets indicate robust models with minimal overfitting.

#### 5. Data and Analysis Insights:

- The balanced dataset (48% hits, 52% non-hits) enabled unbiased training and evaluation of models.
- The correlation analysis revealed strong interdependencies between features (e.g., **energy** and **loudness**) and guided feature selection.

### Challenges Encountered

#### 1. Multicollinearity:

- High correlations between features (e.g., energy and loudness) posed challenges in interpretation and necessitated careful feature selection and scaling.

#### 2. Data Imbalance in Certain Features:

- Features like **instrumentalness** and **liveness** were highly skewed, limiting their predictive contribution and requiring careful handling to ensure model reliability.

#### 3. Complexity of Nonlinear Patterns:

- Polynomial Regression, while improving model performance, introduced additional complexity and computational costs, emphasizing the need for careful regularization and validation.

### Practical Implications

#### 1. For Artists and Producers:

- Insights into key audio features can guide music production, helping artists create tracks that align with listener preferences for hit songs.

- Optimizing features like **danceability**, **energy**, and **valence** can improve the likelihood of success in competitive markets.

## **2. For Marketers and Record Labels:**

- Predictive models can help identify promising songs early, optimizing marketing strategies and promotional efforts.
- Understanding the characteristics of hits allows targeted promotion to specific audiences.

This study demonstrates the power of data-driven approaches in understanding and predicting music popularity. By analyzing audio features and leveraging machine learning models, the project provides actionable insights for artists, producers, and marketers to create and promote successful tracks. While the models effectively capture key patterns, the results highlight the need for continuous refinement through the integration of external factors and advanced techniques. Ultimately, this project lays a strong foundation for future research in the intersection of music, data science, and audience engagement.