

Loughborough University School of Business and Economics

Analytics Research Project for M.Sc. Business Analytics

Visualising a Spotify Dataset to Enable Meaningful Interaction for a Non-specialist Audience

Student ID: F315917

Abstract

Platforms like Spotify have completely changed how listeners access and interact with music in today's digital music world. It is now more important than ever for musicians and producers to understand the factors that influence song popularity on these platforms to effectively navigate the competitive streaming industry. With an emphasis on auditory elements like danceability, energy, tempo, etc. this dissertation analyses the elements that affect song popularity on Spotify.

Leveraging a large dataset from Kaggle that includes dozens of songs from various countries, the study analysed the correlations between different elements and how they affect a song's success using Random Forest and other machine learning models. Though their prediction value varies, the analysis shows that certain auditory elements significantly influence a song's popularity. Strong accuracy was shown by the predictive models, suggesting that these characteristics, combined with artist metadata, are useful for predicting a song's performance on Spotify. Additionally, an interactive dashboard was created using Power BI to provide for a thorough visual representation of the complex interactions. This comprehensive visualization offers practical information to improve song releases and promotional campaigns.

Table of Contents

1.	Intro	oduction	2
2.	Liter	ature Review	4
	2.1.	Introduction to HSS	4
	2.2.	Role of machine Learning in HSS	5
	2.3.	The Role of Spotify in Music Popularity Prediction	6
	2.4.	Audio Feature Analysis and Song Popularity Prediction	7
	2.5.	Regional and Cultural Factors in Popularity Prediction	7
	2.6.	Effect of Music Promotion and Advertising on Popularity	8
	2.7.	Social Network Influence on Music Popularity	8
	2.8.	Challenges in Predicting Song Popularity	g
3.	Meth	odology	10
	3.1.	Data Collection	10
	3.2.	Data Preprocessing	11
	3.3.	Handling Missing Values	12
	3.4.	Checking for Duplicates	12
	3.5.	Outlier Detection and Removal	12
	3.6.	Using ISO Country Codes	13
	3.7.	Feature Selection	13
	3.8.	Exploratory Data Analysis (EDA)	13
	3.9.	Machine Learning	15
	3.10.	Data Visualization using powerBi	16
4.	Analy	sis/Findings & Discussion	19
	4.1.	Descriptive Statistics and Initial Observations	19
	4.2.	Correlation Matrix Analysis	26
	4.3.	Regression Models and Machine Learning Models	30
	4.4.	Summary of Model Performance	37
	4.5.	Power BI Dashboards – Interactive Insights	38
5.	Cond	elusion	40
	5.1.	Project Process Overview	40
	5.2.	Key Findings and Discussion Points	41
	5.3.	Critical Evaluation and Recommendations	42
	5.4.	Recommendations	42
	5.5.	Future Directions	43
	5.6.	Summary	43
6	Refer	ences	11

1. Introduction

In today's data-driven economy, the potential to predict product success before launch has significant advantages, especially in industries with high costs of failure, such as new product development in technology, medicine, or entertainment (Lee & Lee, 2018). However, the task of predicting cultural products' success, such as music, presents unique challenges due to the subjective and emotional nature of consumer preferences (Raza & Nanath, 2020). Traditional product success models frequently depended on quantitative data, such as sales trends or customer behaviour patterns; however, these frameworks frequently did not account for the cultural and emotional elements driving the popularity of music songs (Dhar & Chang, 2009). Thus, the advent of digital music platforms like Spotify, which enable the collection of huge amounts of user data, presents new possibilities for predictive models through the use of machine learning techniques.

Traditionally, the producers and distributors of music have not relied much on leveraging data analytics, statistics, or predictive modelling to forecast the probabilities of success for their products. Rather, the industry has been dependent on the expertise and instincts of tastemakers—people who understand a strong feel of what would appeal to the public. These tastemakers were critical in influencing and forecasting customer preferences, as mentioned by Davenport et al. (2009). But as technology continues to change how we listen to music, statistics play a bigger and bigger part in forecasting success, as music streaming services like Spotify have grown, researchers and business experts have had more access to big data on user behaviour, which has made it possible to investigate more precise, data-driven approaches to forecasting the success of musical works (Pachet, 2012).

An important cultural and economic force, the music business has grown significantly in the last several years. The International Federation of the Phonographic business (IFPI) reports that worldwide music business sales increased 8.1% to \$17.3 billion in 2017. The emergence of digital music streaming services, which have completely changed how people listen to and interact with music, has contributed to certain aspects of this growth. With over 400 million active users, Spotify has become one of the most influential music streaming services globally, the platform now accounts for more than half of global recorded music revenue, with streaming revenues growing by 41.1% in 2017 alone (IFPI, 2018). Beyond merely providing a platform for streaming. Spotify allows musicians to access a huge library of more than 70 million songs while generating an extensive wealth of data on user behaviour, listening preferences, and the audio qualities of songs. This rich dataset, which includes music characteristics offers an unparalleled opportunity to explore the factors that drive song popularity. However, with such a large volume of data available, identifying the key elements that most influence a track's success on the platform presents a significant challenge. The music industry's use of data analytics reached a turning point in 2014 with Spotify's acquisition of The Echo Nest (Sisario, 2014). A music intelligence startup called The Echo Nest focusses on using music data analysis to forecast listener preferences. The way that people engage with music has been revolutionised by Spotify, which improved its recommendation algorithms by incorporating Echo Nest's sophisticated music analytics into its platform.

In order to better comprehend the structure of music and its effect on listeners, the topic of music information retrieval, or MIR, has emerged as a significant area of research. MIR combines musicology with data science and machine learning. According to Herremans et al. (2014), Hit Song Science (HSS) has become more widely recognised within MIR as a field of study devoted to forecasting songs' commercial success prior to their release. As to Pachet's (2012) definition, HSS is a "new area of study that attempts to forecast song success prior to commercial release."

The study of Hit Song Science (HSS) consists of more than just the study of what makes a song popular. HSS uses methods based on science and statistics to analyze and forecast music track success. Researchers have found that certain audio features, such as danceability, tempo, and energy, correlate with song popularity, though these features alone are often insufficient to build highly accurate predictive models (Pareek et al., 2020).

A wide range of machine learning algorithms have been used in recent years to forecast song popularity. These include of Support Vector Machines (SVM), K-Nearest Neighbour algorithms, and Random Forest, which are all intended to investigate data trends and forecast the possibility of a song becoming popular (Lee & Lee, 2018). Because Random Forest can handle big datasets with numerous factors, it has shown to be a very successful model for the challenging task of predicting song popularity on platforms such as Spotify (Sharma et al., 2022).

Research in HSS and MIR has significant implications for everyone involved in the music industry. Record companies may make better judgements about which songs to market by using the insights from predictive models, while artists can utilise them to better understand the factors that contribute to their success. These models can also help streaming services by increasing the precision of their recommendation algorithms and providing consumers with individualised music experiences (Araujo et al., 2019). According to Dhar and Chang (2009), this presents a novel approach for industry experts to evaluate a song's probability of success prior to its release.

Research Aim:

This study aims to explore the key determinants of song popularity on Spotify by analyzing various audio features and artists information. This dissertation uses advanced machine learning models to forecast song popularity and further creates an interactive Power BI dashboard to visualise these associations. This study attempts to add to the expanding body of research in Hit Song Science by tackling the intricate nature of music success prediction and provide useful information for producers, artists, and streaming services.

Research Objectives:

To address the complexities of predicting song popularity in the digital age, this dissertation will focus on the following research objectives:

- Identify the Key Determinants of Song Popularity: The initial goal is to examine the
 connections between different audio characteristics (danceability, energy, tempo,
 etc.), The study focuses on those attributes in an effort to identify the characteristics
 that have the most correlation with Spotify song popularity.
- 2. Develop Predictive Models for Song Popularity: The subsequent objective is to build machine learning models that, using the determined factors, can forecast a song's popularity with accuracy. To create a trustworthy tool for predicting song success, these models will be assessed for prediction accuracy and robustness.
- 3. Create an Interactive Dashboard for Data Visualization: Using Power BI to create an interactive dashboard is the third goal. Artists, producers, and other industry experts will find this dashboard to be a useful tool as it will enable users to dynamically investigate the links between different attributes and song popularity.

2. Literature Review

2.1. Introduction to Music Popularity Prediction and Hit Song Science (HSS)

Predicting a song's success has changed in recent years from intuition-based methods to advanced, data-driven approaches. To identify possible hits in the past, the industry was dependent on the subjective judgment of tastemakers, music executives, and producers who judged a track's likelihood of success based on market trends, cultural significance, and their own experience (Pachet, 2012). However, as the industry has evolved, data-driven approaches have replaced subjective assessments. With the emergence of Hit Song Science (HSS), the focus has shifted to evaluating quantifiable elements such as user involvement, release date, and audio features to predict song popularity.

According to Pachet (2012), HSS is an emerging field that uses objective metrics analysis to predict songs' success before they are officially released. Advances in machine learning algorithms and the growing availability of large-scale music datasets have made it possible for a new wave of predictive analytics in music consumption (Nijkamp, 2022). We can now apply machine learning models to analyze correlations between particular song attributes and track success. HSS research often integrates data gathered from exterior aspects, including user interactions and social media engagement, with inherent features of a song, like tempo, energy, and danceability (Araujo et al., 2019).

In its present form, HSS creates prediction models using data from online music services like Spotify, YouTube, and Apple Music. Record companies and artists may make better judgements regarding their releases with the use of these models, which can produce precise predictions about a song's likelihood of success (Pareek et al., 2022).

2.2. Role of Machine Learning in Hit Song Science

Modern Hit Song Science (HSS) has made machine learning a pillar of research, revolutionizing the way researchers approach the problem of song success prediction. According to Araujo et al. (2019), machine learning algorithms can identify patterns and links between song features and commercial performance by utilizing the Spotify Web API and other data sources, which provide large datasets including both structured and unstructured data. Specifically, music popularity has been predicted remarkably well by supervised learning models like Gradient Boosting, Random Forests, and Support Vector Machines (SVM) (Pareek et al., 2022).

2.2.1. Supervised Learning Techniques

In HSS, supervised learning is one of the most popular methods. It involves utilizing historical data to train machine learning models. The target variable, such as song popularity, is known, and the model learns the correlation between the outcome and the input variables, such as tempo and danceability (Nijkamp, 2022).

Given their capacity to manage big and complex datasets, Random Forests, a type of ensemble learning technique, have proven very successful in HSS. To increase prediction accuracy, Random Forests integrate the results of several decision trees that are constructed during the training stage (Pareek et al., 2022). Empirical research has demonstrated that this model is capable of

accurately representing both linear and non-linear relationships between song attributes and their popularity.

With an accuracy rate of more than 85%, the Random Forest model particularly proved very successful at forecasting popular songs (Pareek et al., 2022).

Another appreciated supervised learning method is Support Vector Machines (SVM), which is excellent at categorizing songs according to both acoustic and non-acoustic features. Complex correlations between song qualities that are not immediately evident through simpler models can be identified using SVM models, especially those that use non-linear kernels such as the Radial Basis Function (RBF) (Araujo et al., 2019). This is particularly helpful in narrow genres where typical linear models might not be able to forecast hit songs.

Meanwhile, by iteratively increasing prediction accuracy, Gradient Boosting Machines (GBM) and XGBoost have been used to predict song success. By concentrating on rectifying the mistakes made by weaker models, these algorithms can become more accurate over time and forecast more accurately whether a song will chart on platforms such as Spotify (Vicente et al., 2019).

2.2.2. Unsupervised Learning and Clustering

Based on similarities in their acoustic properties, songs can be grouped into clusters using unsupervised learning approaches like Principal Component Analysis (PCA) and K-means clustering. Using these methods, one can find patterns in the data that conventional analysis might have overlooked (Shulman et al., 2016). These methods do not require data to be categorized.

Song characteristics like tempo, danceability, and valence, for instance, have been used to classify songs into clusters using K-means clustering. In this way, genres or styles that are more likely to produce hits can be identified by researchers. According to Reiman and Örnell (2018), principal component analysis (PCA) helps in reducing the dimensionality of huge datasets, facilitating improved data understanding, as well as better feature association visualization.

Hit Song Science research has undergone a radical change because of machine learning approaches that integrate supervised learning, unsupervised learning, and ensemble methods. These models can accurately forecast which songs are most likely to succeed by examining large datasets with multiple variables (Araujo et al., 2019).

2.2.3. Feature Selection for Popularity Prediction:

In a study published in 2020, Rahardwika et al. investigated how feature selection affected the ability to accurately classify musical genres. The researchers enhanced the model's accuracy by merging various feature groups (such as acousticness, instrumentalness, and danceability), proving that feature selection is a crucial factor in how successful machine learning algorithms are (Rahardwika et al., 2020). The most useful characteristics in their algorithm for popularity classification were determined to be acousticness, instrumentality, energy, danceability, valence, and loudness.

Pareek et al. (2022), who compared machine learning models with and without feature selection, further emphasized the significance of feature selection. According to their results, feature-selected models performed better than those without it, especially when it came to measures like precision and recall. According to the study's findings, the most crucial characteristics for predicting a song's popularity were its acousticness, instrumentality, and danceability (Pareek et al., 2022).

2.3. The Role of Spotify in Music Popularity Prediction

The traditional music industry has undergone a revolution due to the substantial changes brought about by the emergence of digital streaming platforms in terms of music consumption, distribution, and monetization. One of the biggest streaming platforms, Spotify, contributes significantly to this change by giving customers access to millions of songs and gathering copious quantities of data about how music is listened to. The availability of large-scale, real-time statistics has created new opportunities for studying listener behaviour and forecasting song success.

2.3.1. Data Collection and the Shift to Streaming Metrics:

The platform is a perfect source for predictive analysis because of its exceptional capacity to track both user activities (such as streams, playlist additions, and replays) and specific song attributes (such as acoustic variables like tempo and valence) (Gulmatico et al., 2022).

Spotify's Web API is an important advantage as it gives analysts access to a variety of song attributes, including danceability, energy, loudness, and speechiness, in addition to user engagement metrics. These characteristics have been used by researchers to create predictive models that assess song performance in various nations and genres (Vicente et al., 2019).

For instance, Araujo et al. (2019) used data from Spotify's Global Top 50 charts to analyse how song attributes and user behaviour affect commercial success. They found that faster and more energetic songs have a better chance of making it onto the charts, particularly in genres like pop and electronic music that are popular on Spotify playlists (Araujo et al., 2019). Furthermore, the study emphasized the significance of playlist placements, emphasizing that a song's inclusion on playlists such as "New Music Friday" or "Top Hits" has a substantial impact on stream counts and overall popularity (Araujo et al., 2019).

Data-driven insights into music consumption are now even more readily available because of Spotify's acquisition of The Echo Nest, a music intelligence firm. According to Gulmatico et al. (2022), Spotify leverages Echo Nest's technology to provide personalized recommendations and insightful analysis to researchers and artists who are interested in understanding the aspects that influence hit song success. Spotify facilitates the creation of models that forecast not just which songs will become popular but also how and why certain tracks will resonate with audiences by fusing song information with listener behaviour data.

2.3.2. Spotify User Behaviour and System Dynamics:

The dynamics of Spotify's user base have been examined in a number of research. Pareek and colleagues (2022) conducted an analysis of Spotify

users' session durations, track choices, and idle time, offering valuable insights into their everyday music-listening habits. Analysing general network features and performance, Goldmann and Kreitz (2020) concentrated on the performance of Spotify's network infrastructure by gathering IP address data via NAT devices. In order to more accurately forecast user behaviour and, eventually, song popularity, these studies highlight the significance of comprehending user interaction patterns and system performance.

2.4. Audio Feature Analysis and Song Popularity Prediction

The connection between audio characteristics and song popularity has been the subject of numerous studies. Trpkovska et al. (2022) examined the audio qualities of songs featured on Spotify's 2017 list of the best songs, identifying recurring elements that support the songs' popularity. In order to detect patterns in the acoustic characteristics of well-known songs and to estimate one attribute based on others, the study used data visualization and data mining (Trpkovska et al., 2022). This study showed that a song's overall success is strongly correlated with audio elements including tempo, energy, and danceability.

Pareek and colleagues (2022) utilized machine learning methods, including Random Forest, K-Nearest Neighbour (KNN), and Linear Support Vector Machine (SVM), in their investigation of predicting song popularity using Spotify measures. After evaluating each system's recall, accuracy, precision, and F1-score, they came to the conclusion that the Random Forest algorithm performed the best in predicting song popularity (Pareek et al., 2022). Mora and Tierney (2020) evaluated to assess the impact of Feature Engineering and Feature Selection on model efficiency. Their findings indicate that both strategies lead to a considerable improvement in prediction model performance.

While Herremans et al. (2014) concentrated on the prediction of dance success songs using audio features like timbre, Singhi and Brown (2014) suggested that lyrics and audio content may be combined to predict hits. According to their findings, one can determine a song's likelihood of making it to the top 10 dance charts by analysing its acoustic properties (Herremans et al., 2014). This area of study has expanded, with researchers looking into novel approaches to utilize audio data to forecast a song's likelihood of success across a range of industries and genres.

2.5. Regional and Cultural Factors in Popularity Prediction

Regional and cultural influences have a significant impact on song popularity; local customs, tastes, and society preferences are major determinants of music consumption. While worldwide trends impact popular music's overall structure, local tastes frequently dictate which songs are most successful in certain countries (Ni et al., 2011). The significance of including regional characteristics into predictive models is shown by the swift growth of streaming services like as Spotify, given the notable variations in listening habits observed globally (Ferrer et al., 2020).

Regional Preferences and Local Genres: Studies reveal that even when there are international hits, local music frequently top regional charts. For instance, traditional genres like dangdut continue to be quite popular in Indonesia, surpassing Western pop music in popularity (Saragih et al., 2023). Comparably, reggaeton has emerged as a major genre in Latin America, both domestically and globally, demonstrating the significance of cultural relevance in forecasting the success of a musical endeavour

(Ferrer et al., 2020). These results imply that in order to increase forecast accuracy, models should take local genre preferences into consideration.

Global vs. Local Trends: Martin-Gutierrez et al. (2020) research emphasises the need of combining regional data with worldwide music trends. Internationally, global music genres like pop and electronic music frequently rule, although regionally specific genres like K-pop and Afrobeat perform okay. Creating more precise prediction models requires combining local and global patterns. Machine learning models can better capture the subtle differences of area preferences and increase forecast accuracy by combining both global and local factors (Martin-Gutierrez et al., 2020).

Integrating Regional Data into Machine Learning Models: According to Ferrer et al. (2020), machine learning algorithms that do not take into consideration cultural and geographical characteristics frequently do not do an effective job of forecasting the performance of local songs. Models can anticipate outcomes more accurately by including region-specific variables, such as regional genres, linguistic preferences, and cultural trends. According to Saragih et al. (2023), by reflecting localised trends in music consumption, the inclusion of such parameters enhances model performance.

2.6. Effect of Music Promotion and Advertising on Popularity

The effect of advertising and promotion tactics on song popularity has also been studied. The study conducted by Shulman et al. (2016) investigated the impact of Spotify's marketing strategy on the financial dynamics of the music industry. The study specifically focused on the advertising strategies that facilitate the growth of the platform and the promotion of music content. According to Shulman et al. (2016), they discovered that promotional activities, such as playlist placement and advertising campaigns, have a big influence on song visibility and stream counts.

Comparably, Kurt et al. (2018) used Spotify data for a personalized music recommendation study that improved user experience by employing real-time resting models to tailor music recommendations based on listener moods (Kurt et al., 2018). Spotify is a major participant in boosting song popularity through targeted marketing because of its advertising strategies, which employ digital advertisements based on user profiles to further encourage user interaction (Mora and Tierney, 2020).

2.7. Social Network Influence on Music Popularity

Using social network data to forecast music success is another well-liked method. Research from Kim et al. (2014) and Zangerla et al. (2016) have demonstrated the value of social media data, especially from Twitter, in forecasting music popularity. This research discovered a significant correlation between Twitter activity and music chart performance by examining user tweets and comparing them with chart data. For instance, Zangerla et al. (2016) showed that Twitter messages might reasonably accurately forecast future chart performance, particularly when paired with recent music charts. Similar to this, Bischoff et al. (2009) developed a song popularity prediction model with Last.fm interaction data, demonstrating encouraging outcomes as well.

However, there are certain difficulties in predicting music popularity using social media data. The requirement for an extensive amount of data to provide precise forecasts is one of the primary limitations. Large computing power and sophisticated

machine learning models are needed for sorting through the massive volumes of data generated by social networks and identifying relevant patterns (Bischoff et al., 2009). Furthermore, social media trends can fade quickly, making it challenging to extrapolate long-term success from transient surges in user activity.

2.8. Challenges in Predicting Song Popularity

Predicting song success accurately still faces considerable problems, even with the advances in predictive modelling. The subjective aspect of music is one of the main problems. Even though machine learning models are capable of seeing trends in data, they frequently have trouble taking into consideration the emotional and cultural aspects that affect listener choices. It is difficult for even the most sophisticated models to forecast which songs will become mega-hits with total confidence because music consumption is highly personal and what appeals to one listener may not appeal to another, as Pachet (2012) noted.

Furthermore, external factors like viral trends, recent events, or celebrity endorsements can have a significant impact on a song's popularity in ways that are hard to forecast based just on past data. A song may see an increase in streams, for example, if it appears in a popular TikTok challenge or if an internationally recognized musician or influencer posts about it on social media. Predictive models face difficulties as a result of these unpredictable cultural moments because they frequently base their forecasts on solid, historical data (Shulman et al., 2016).

A further difficulty is the rapid shift in popular music trends. Because the music industry is so dynamic, what is considered a "hit" now might not be in the future. According to Reiman and Örnell (2018), prediction models need to be updated frequently to account for evolving listener preferences, new genres, and changed music consumption habits. With time, machine learning models that depend on outdated data can miss these changes, which would reduce their accuracy (Reiman and Rrnell, 2018).

2.9. Opportunities:

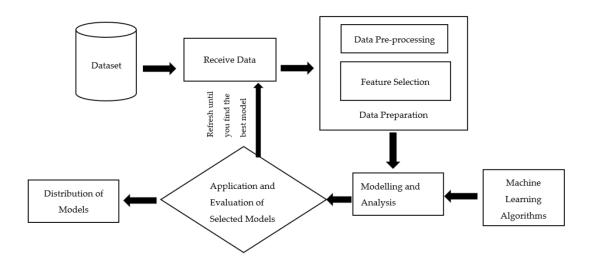
The use of machine learning in HSS has great potential despite these obstacles. With increased data availability and sophisticated machine learning algorithms, song popularity estimates should be more accurate over time. Future developments in natural language processing (NLP) and deep learning may make it possible to analyse song lyrics, genre progression, and listener sentiment in even more detail, which would increase the predictive capacity of HSS models (Gulmatico et al., 2022).

2.10. Conclusion:

Large datasets from music streaming platforms like Spotify have been available, and this has allowed the field of hit song science (HSS) to advance significantly in recent years. Researchers have created increasingly precise models for forecasting song popularity by fusing machine learning techniques with social network data and acoustic component analysis. However, challenges remain, particularly in the area.

3. Methodology

This methodology describes the step-by-step procedure for analysing the primary factors influencing song popularity on Spotify, developing predictive models, and visualizing the findings using an interactive dashboard. Each stage of the research is well discussed, from data collection and preprocessing to model creation and assessment, as well as the reasoning for the tools and techniques used.



The process of developing a machine learning model is shown in the flowchart, which begins with obtaining the dataset and proceeds through feature selection, data preparation, and data pre-processing. The process moves on to model construction and analysis, evaluating various models before choosing the best one. The data visualization stage is the last one.

3.1. Data Collection

Data collecting is an essential stage in ensuring the quality and relevance of research, especially when dealing with huge amounts of data from global platforms such as Spotify. For this research, data was sourced from Kaggle, a trusted platform for high-quality datasets suitable for machine learning and data analysis. Although the dataset was acquired through Kaggle, its source is Spotify's Web API, which gives researchers and developers access to the Spotify's music catalogue, track information, and audio attributes. The dataset utilized in this research includes the top 50 songs trending in over 70 countries between 2023 and 2024, with daily updates offering real-time insights into the global music environment. It is crucial to note that, while the data indicates song popularity in 2023 and 2024, the release dates of the songs fluctuate, indicating that some tracks may have been published in prior years. Over 70 million recordings were included in the collection.

The dataset includes a wide range of attributes that are important for the analysis:

Track Metadata:

Track ID	Identification number assigned to every song.	
Track Name	The song title.	
Artist Name	The name of the artist(s).	
Album Name	The song's album where it appears.	
Album Release Date	The song's or album's first release date.	
Country	The nation in which the song is currently popular.	

• Audio Features (Provided by Spotify's Web API):

danceability	Danceability is an indicator of a song's suitability for dancing. This meta characteristic includes beat strength, regularity, and rhythm stability. The range of this value is 0.0 to 1.0.	
energy	Represents the track's degree of activity and intensity; tracks with higher levels are considered to be more energetic. The energy of a music is influenced by its timbre, perceived loudness, and dynamic range. The range of this value is 0.0 to 1.0.	
key	The key for track notation. This value is encoded from 0 to 11, with 0 being C, 1 being C#, 2 being D, and so on.	
mode	A binary number that indicates the major or minor mode of the song. A song with a value of 0 is in major, while 1 is in minor.	
valence	Determines how positive a track is musically; higher numbers correspond to more uplifting, positive tracks. This value lies between 0.0 and 1.0.	
loudness	The overall loudness of a track is expressed in dB. This value spans from -60 to 0 dB and is averaged across the track.	
speechiness	An estimation of the total amount of spoken words in a song (or other audio file). For instance, a podcast will be close to 1.0 speechiness, whereas a music that is only instrumental will be closer to 0.0.	
instrumentalness	Defines a song's instrumentality. If the score is 1.0, the song is entirely instrumental, and if it is 0.0, no instruments are used at all. Tracks with values greater than 0.5 are meant to be instrumental.	
acousticness	Indicates if the music has an acoustic instrumentation or not. values between 0.0 and 1.0.	
liveness	A value between 0.0 and 1.0 represents the likelihood that the song will be recorded in front of a live audience. A higher liveness score indicates a greater likelihood of the music being recorded live.	
duration_ms	The song's duration expressed in milliseconds.	
time_signature	The song's time signature and the number of bar beats.	
tempo	The track's approximate tempo, measured in beats per minute (BPM)	
explicit	Indicating whether or not a song's lyrics include explicit words, or whether the music has potentially offensive or obscene lyrics. This value is binary, indicating that 1 represents True and a 0 represents False.	

Popularity Scores:

popularity	A statistic that Spotify created to show how popular a music is.	
	Ranges from 0 to 100.	

3.2. Data Preprocessing

The Analysis could not begin until the raw dataset was thoroughly pre-processed. Making sure the data was clear, comprehensive, and appropriate for machine learning algorithms required taking this crucial step. The subsequent steps were taken:

3.2.1. Handling Missing Values

Several columns, including album_release_date, album_name, and country, had missing values. It was decided to exclude the rows with missing data rather than imputing them since these fields were essential for analysis and comparison, especially in cross-country and time-based evaluations. For example, the album_release_date field would make it challenging to detect patterns in song releases, and missing values in the nation column would preclude accurate geographic analysis. Similarly, album-level analysis was interfered with the lacking album_name.

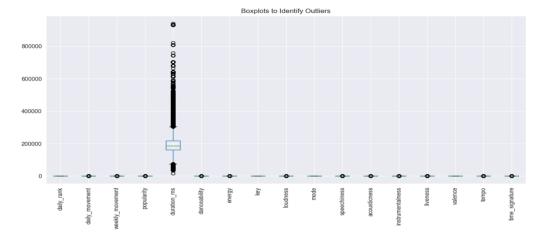
3.2.2. Checking for Duplicates

Checked for duplicate entries in the dataset. There are several possible causes of duplicate rows, including numerous snapshots or incorrect data entry. However, no duplicates were discovered in this instance, guaranteeing that every record provided a unique data point for analysis.

3.2.3. Outlier Detection and Removal

Finding and managing outliers was a crucial component of data cleansing, especially in the duration_ms column, which shows song durations in milliseconds. In this case, outliers have the potential to skew the outcomes of any investigation or predictive model.

A box plot was utilised in order to locate these outliers. Box plots are a useful tool for visualising data distribution because they emphasise values that deviate greatly from the predicted range, or beyond the whiskers. Points that are more than 1.5 times the interquartile range (IQR) above or below the first quartile are commonly referred to be outliers.



In order to preserve consistency and keep these outliers from distorting model results, they were eliminated. By doing this, the normal song duration and its correlation with other factors, such popularity, may be more properly represented by the model.

3.3. Using ISO Country Codes

Originally, nations were represented in the dataset using ISO country codes, which are two-letter abbreviations that can be challenging to interpret, particularly when used for visualisation. These ISO country codes were translated into their complete names for easier understanding. This conversion makes it simpler to understand the data in analyses and visualisations, particularly when contrasting the popularity of songs in various nations.

3.4. Feature Selection:

A correlation matrix was used to determine which factors were most important for forecasting song popularity. The pandas package in Python was used to generate the correlation matrix. In order to ascertain which characteristics have the strongest correlation with song popularity, this matrix calculates the link between every attribute and the target variable.

Popularity vs. Other Features: There are limited linear associations between the characteristic "popularity" and the majority of other features, as seen by the relatively modest correlations between them. However, a few attributes stand out:

- Energy (r = 0.13): Songs with more energy typically have a slightly higher popularity score, indicating a moderately positive correlation.
- Loudness (r = 0.13): There is a relatively positive link between loudness and popularity.
- Valence (r = 0.11): Happy-sounding songs may be more popular. Valence, which defines musical optimism, shows a small but positive connect with popularity.

Non-significant Correlations: Some qualities, such as danceability, instrumentalness, and speechiness, show very weak or near-zero correlations with popularity, implying they are less likely to be relevant predictors.

Multicollinearity: The multicollinearity between characteristics was taken into consideration in addition to the direct link with popularity. For example:

- Energy and Loudness (r = 0.74): These two variables showed a substantial positive association, suggesting that they both reflect a comparable component of music intensity.
- Acousticness and Energy (r = -0.54): Acoustic songs are often less energetic, according to the substantial negative association between acousticness and energy

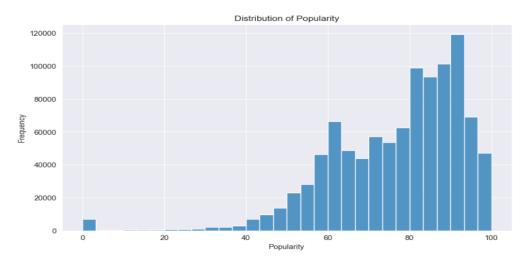
3.5. Exploratory Data Analysis (EDA)

During the Exploratory Data Analysis (EDA) stage of this study, my objectives were to get a more profound understanding of the dataset and recognize significant patterns, trends, and anomalies. This stage was essential in figuring out how different traits related to the goal variable, popularity, and in laying the foundation for developing models.

3.5.1. Distribution of Popularity:

The initial phase of the EDA implied mapping the distribution of the popularity attribute. To see the frequency of songs at various levels of popularity, a histogram was employed. A tilt towards higher popularity ratings was evident

from the chart, where most songs are grouped around the upper half of the scale (60–100). This suggests that the majority of the songs in the dataset are those that have received a lot of publicity on the platform.



3.5.2. Distribution of Song Durations:

The distribution of song time was another significant attribute covered at the EDA. It was critical to comprehend the variations in song duration within the dataset, as these variations might have a substantial effect on the user experience.

3.5.3. Popularity Over Time:

Keeping track of the average popularity of songs over time was a crucial component of the research. I was able to observe patterns in popularity trends over particular time periods by averaging the popularity scores according to the snapshot date. The time series analysis revealed several variations in popularity, which might be explained by new music releases, seasonal trends, or outside factors that have an impact on the music business. Finding these patterns made it easier to comprehend how user preferences changed over time.

3.5.4. Average Popularity by Country:

This dataset allowed for a global comparison of musical preferences by providing statistics on song popularity across many nations. The analysis's main goal was to determine each nation's average song popularity. This indicated notable regional variations in musical tastes, with certain nations continuously showing higher average popularity scores. Understanding the regions with the most popular music was made possible by visualizing the top 10 nations based on average popularity.

3.5.5. Audio Features and Popularity:

A significant aspect of the EDA was examining how certain auditory qualities, like as danceability, energy, pace, and others, affected the popularity of songs. In order to conduct this research, the distribution of these characteristics for popular songs—that is, songs with higher popularity scores—was compared to that of less popular songs.

3.6. Machine Learning

Python and its robust libraries, such Pandas and Scikit-learn, were used by me to develop a number of models in the Machine Learning Models. These technologies are essential for analyzing the information and extracting forecast insights because of their adaptability, effectiveness, and simplicity of integration with machine learning algorithms. Pandas was employed to efficiently clean and manipulate the data, making sure the big dataset was handled without a hitch. A variety of machine learning models were implemented using Scikit-learn, which also offered preprocessing and model assessment tools. Seaborn and Matplotlib were also used for data visualization, which allowed for the creation of distinct visual representations of the outcomes. The model building and performance analysis processes were made more efficient by the combined use of these.

Machine Learning Models:

3.6.1. Linear Regression

One of the most fundamental methods in machine learning is linear regression. It makes a prediction that there is a linear connection between one or more independent variables (features) and the dependent variable (popularity). In this experiment, audio characteristics including energy, danceability, and loudness were utilized to predict song popularity using linear regression. Although Linear Regression is straightforward and interpretable, it didn't perform as well compared to more sophisticated models, as the link between characteristics and popularity may not be precisely linear.

3.6.2. Random Forest

Random Forest is an ensemble learning method that generates several decision trees during training. Every tree offers a classification; the final prediction is determined by the majority vote of the trees. Because Random Forest aggregates numerous trees, it is resilient and lowers the risk of overfitting. Because Random Forest can manage non-linear correlations between song attributes and popularity, it performed better in this study. To increase the accuracy of the model, hyperparameters like the quantity of trees were optimized.

3.6.3. Decision Tree

Decision Tree is a non-parametric supervised learning model that divides the dataset into subsets based on the most important attributes at each node. To forecast the result, the model learns basic decision-making rules. In addition to offering a simple method for predicting song popularity, it was very helpful for comprehending the significance of features.

3.6.4. K-Nearest Neighbors (kNN)

A non-parametric approach called K-Nearest Neighbours (kNN) uses the majority vote of the data points' closest neighbours in feature space to classify the data points. In this study, kNN was used to evaluate its performance against alternative models.

3.6.5. Multiple Linear Regression

The link between a dependent variable and several independent factors is modelled by multiple linear regression, which expands on the ideas of linear regression. Multiple Linear Regression extends the concept of linear regression by modeling the relationship between a dependent variable and many independent variables. This model was employed to account for several song attributes concurrently, such as tempo, valence, and speechiness.

An 80-20 train-test split was used to train these models, ensuring an accurate evaluation of their performance. To further enhance model accuracy and adjust the hyperparameters, a grid search was employed. It was verified that the models were not overfitting and could be applied to new, unseen data by implementing cross-validation procedures.

3.7. Data Visualization using powerBi

Power BI was selected because of its powerful capabilities for building interactive dashboards, versatility, and ease of use. It has excellent abilities to filter, slice, and interact with the data, producing real-time insights, and it facilitates the smooth integration of datasets. Its ability to filter data by several dimensions, such country, genre, or release date, is one of its main advantages. This enables stakeholders to concentrate on subsets and investigate the more detailed and granular link between song features and popularity. A range of features, including KPIs, sliders, filters, and interactive charts, made the insights available. An overview of the main features and procedures put in place is provided below:

3.7.1. Data Preparation and Import

- Python was used for substantial data cleaning and preparation prior to importing data.
- In order to visualize the processed dataset, it was exported as a CSV file and imported into Power BI.
- Numerical variables such as song duration and track popularity have been structured and imputed.

3.7.2. Data Modelling and Relationships

To make sure that the links between various tables, including tracks, artists, nations, and song features, were well established, the dataset was modelled in Power BI. Relationships were created using common data, such as Spotify track IDs and country names, to facilitate cross-filtering and aggregation, allowing various visualization elements to work together seamlessly.

3.7.3. Key Visualizations and Features

- KPI (Key Performance Indicators):
 - Total Tracks: A KPI widget that provides a high-level overview of the data size displays the total number of tracks evaluated in the dataset (970K).

- Average Duration: A KPI that shows the average song duration in minutes (3.15 minutes), providing information on the typical length of popular songs.
- Average Popularity: The KPI that showed the average popularity score across the board (76.81), which was an essential indicator for figuring out the trends in the dataset.

Map Visualization (Average Popularity by Country):

This map uses green dots of different sizes to show the average popularity scores by country. Greater dots indicate nations where songs possessed greater popularity rankings, providing geographical information on local tastes in music.

Pie Chart (Distribution of Track Popularity):

A pie chart that shows the music sorted down into three categories: "Hit," "Moderate," and "Underrated." Understanding the overall distribution of song performance was made easier with the help of this illustration.

• Bar Chart (Popularity by Artists):

The number of tracks and average popularity scores are used to rank artists in this bar chart. Leading artists received prominence in the visualization.

• Line Chart (Popularity by Month):

A line chart that shows the average popularity of songs over time, allows users to examine monthly or seasonal patterns. This graph illustrates the year-over-year peaks and valleys in track popularity.

Average of Danceability, Energy, Instrumentalness, Liveness, Valence, and Speechiness:

A visual representation of how these characteristics is dispersed throughout the dataset and how they could affect track popularity is provided by a bar chart that displays the average values for these song variables.

3.7.4. Interactive Filters and Slicers

Slicers and filters, two interactive Power BI tools, were utilized to create a dynamic user interface:

Country Filter:

To delve deeper into regional trends in song popularity, users can choose specific countries (such as Argentina, Australia, or Austria).

• Time Period Slicer:

This slicer helped show how song popularity changed over time by allowing users to select data by particular months or years.

Artist Filter:

A slicer has been created available to examine the popularity of songs by particular artists, providing a focused perspective on each artist's performance throughout all periods.

3.7.5. Sliders for Adjusting Attributes

Interactive sliders for important song properties were included in the dashboard, allowing users to view changes in expected popularity in real time and simulate adjustments:

Adjusting Loudness:

Users may slide to change the loudness levels and see how it affected the modified popularity scores.

• Adjusting Energy:

In the same way, users could alter energy levels to observe the relationship between energetic music and popularity.

• Tempo and Valence Adjustments:

Sliders made it possible to alter tempo and valence in real time, offering fresh insight into the relationship between these elements and song popularity.

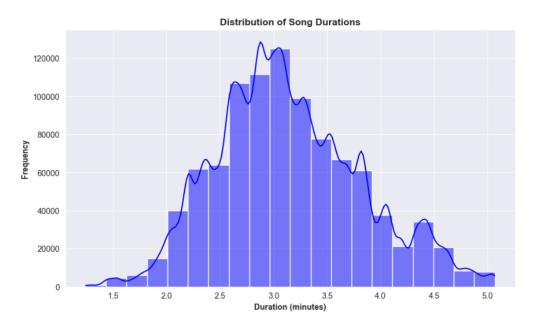
4. Analysis/Findings & Discussion

With the study objectives mentioned in previous parts as a guide, we concentrate on the indepth investigation of song popularity in this chapter. Our study uses statistical models and visualisations to break down the information in order to provide important insights into the elements that most affect Spotify song popularity.

4.1. Descriptive Statistics and Initial Observations

4.1.1. Distribution of Song Durations

The distribution of song durations in minutes can be observed in the histogram below:



Key Observations:

- The average song spans between 2.5 and 4 minutes, peaking approximately 3 minutes.
- Longer songs are rare in the top 50 charts, as seen by the tail of the distribution, which reveals that relatively few songs extend more than 4.5 minutes.

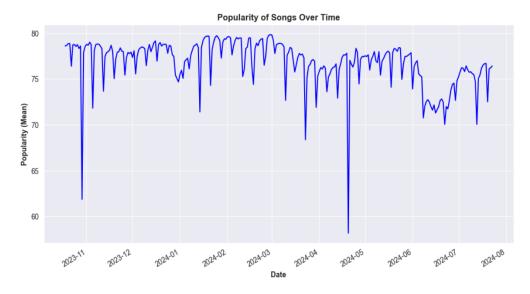
Statistical Insight:

- Mean duration: The average time was about 3.3 minutes (198,000 ms).
- Median duration: The median time was about 3.1 minutes.
- Skewness: The analysis is dominated by shorter songs, with a slight positive skew in the distribution that indicates fewer long songs overall.
- Kurtosis: The majority of durations are centred around the mean, according to the moderate kurtosis value.
- Interpretation: Modern streaming patterns are supported by the prevalence of tracks that are 3 minutes or less. Shorter songs are preferred because

they increase the number of repeat plays, which boosts their position on platforms like Spotify. Furthermore, this length better fits radio playtimes, which increases a song's exposure and probability of becoming a hit.

4.1.2. Popularity Trends Over Time

The line chart below illustrates the trend of song success from November 2023 to August 2024 by displaying the average popularity of songs over time.



Key Observations:

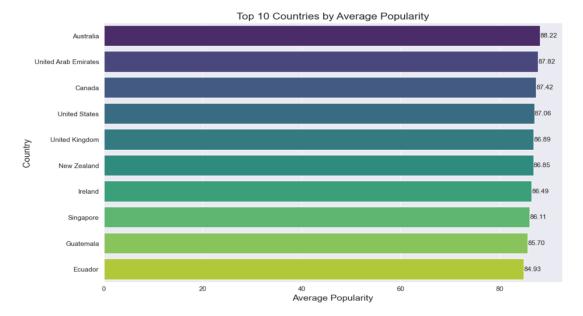
- The average popularity score is consistent overall, ranging from 75 to 78
- Major drops are shown in December 2023 and May 2024, most likely as a result of the increase of special or seasonal music at these times.
- In March 2024, there is a minor increase in popularity, indicating a surge of new, popular releases.

Statistical Insight:

- Mean popularity score: Approx 76.5.
- Standard deviation: About 2.1, suggesting that there is limited annual variation in the popularity scores.
- Outliers: The popularity scores drops at specific times indicate the existence of special or seasonal songs that momentarily decrease the average.
- Correlation with seasonality: There is a definite relationship between periods of high popularity and significant holidays or music release schedules, which propels seasonal or niche releases into the charts.
- Interpretation: The average popularity score shows a steady consumption trend among listeners despite small declines. Major releases or promotional events are probably associated with popularity peaks, and seasonal songs that differ from conventional tastes may be associated with popularity falls. This pattern illustrates how certain songs are resilient enough to maintain their appeal over time.

4.1.3. Top Countries by Average Popularity

The top 10 countries are shown in the bar chart below according to the average song popularity.



Key Observations:

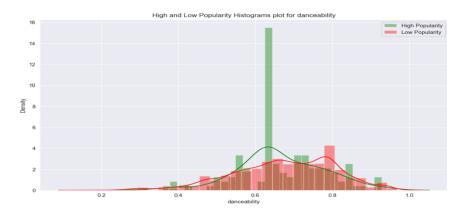
- With an average popularity score of 88.22, Australia comes in top, followed by the United Arab Emirates with 87.82.
- With ratings more than 86, Canada, the United States, and the United Kingdom are among the other highly ranked nations.
- o Ecuador, with a score of 84.93, rounds out the top 10.

Statistical Observations:

 Standard Deviation: The average popularity score variance across the top 10 nations is rather minimal (standard deviation of about 1.21), suggesting a high degree of consistency among popular songs in these regions.

4.1.4. Comparison of High and Low Popularity Songs by Song Attributes

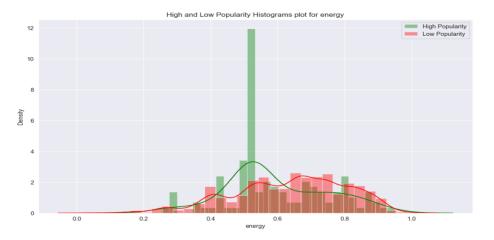
(a) Danceability: The danceability ratings of songs with high and low popularity are compared:



- Key Observation: The mean danceability score for songs with high popularity is 0.61, and the mean score for songs with low popularity is 0.51.
- Statistical Insight: The mean danceability ratings of popular songs with varying popularity varied by 0.10, suggesting that listeners choose songs that are more danceable.
- Interpretation: It appears that danceability has a big role in song popularity, and that higher scores correlate into more success.

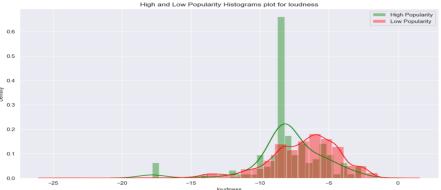
(b) Energy:

The energy levels of songs with high and low popularity are compared:



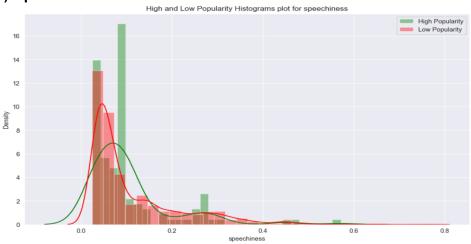
- Key Observations: Songs that are more popular have greater energy levels; their mean energy score is 0.68, whereas that of less popular songs is 0.53.
- Statistical Insight: Based on statistical analysis, it may be inferred that music with a 0.15 energy level difference are more likely to engage the attention of audiences.
- Interpretation: Popular songs, especially those in mass-appealing genres, are characterised by their energy. The fact that highenergy music perform better in social situations like parties, gyms, and get-togethers probably contributes to their popularity.

(c) Loudness:



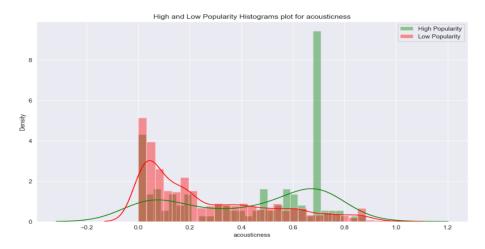
- Key Observations: The average loudness of high-popularity songs is -6.4 dB, whereas low-popularity songs have an average loudness of -8.9 dB.
- Statistical Insight:
 - Mean difference: The mean difference between popular songs and less popular songs is 2.5 dB, suggesting that loudness has a significant role in drawing in listeners.
 - Skewness: Popular songs have a larger skewness in the loudness distribution, which is indicative of the "loudness wars" that have influenced modern musical developing.
- Interpretation: A song's loudness boosts intensity and attracts the listener
 in. Louder tunes stand out more in competitive streaming contexts and
 are generally interpreted as more thrilling and energetic. This might
 explain why songs with greater volume typically top the charts.

(d) Speechiness:



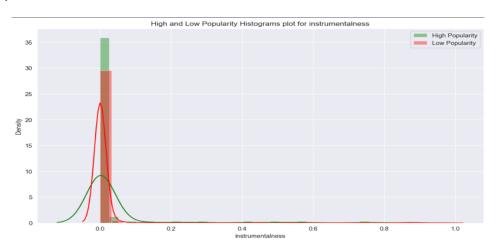
- Key Observations: High popularity songs have a lower mean speechiness score of 0.09, compared to 0.11 for low popularity songs.
- Statistical Insight:
 - Mean difference: A little but steady trend is seen by a 0.02-point reduction in speechiness in songs with high popularity.
 - Skewness: The significance of lyrics and melody in popular music is shown by the minor skewness of both high and low popularity songs towards lower speechiness ratings.
- Interpretation: Songs that prioritise singing over spoken word, or that are
 less conversational overall, typically do better on streaming services. This
 implies that, in contrast to spoken-word or rap genres, which could have a
 more specialised appeal, melody and lyrical substance are more appealing
 to a large audience.

(e) Acousticness:



- Key Observations: The mean acousticness score of songs with high popularity is 0.26, whereas that of songs with low popularity is 0.40.
- · Statistical Insight:
 - Mean difference: Popular songs have an acousticness score 0.14 points lower than other songs, which is indicative of the popularity of studio-produced and electronic music on mainstream charts.
- Interpretation: Popular song's lower acoustic quality indicates that
 electronic production, created rhythms, and studio effects are more
 frequently used in pop music nowadays than acoustic instruments. This is
 consistent with the music industry's current extensive adoption of digital
 production.

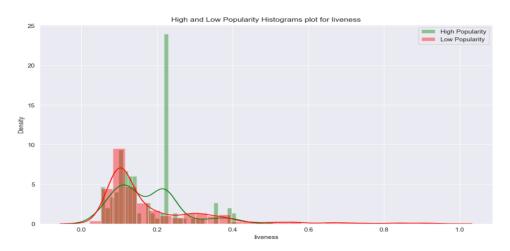
(f) Instrumentalness:



- Key Observations: Songs with a high level of popularity score 0.02, while songs with a low level of popularity score 0.05 on average.
- Statistical Insight:

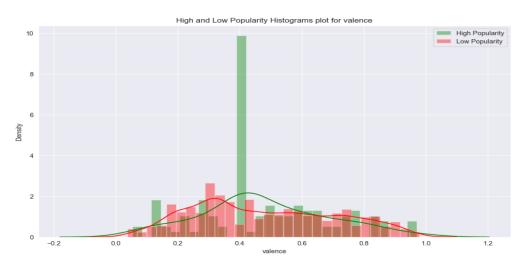
- Mean difference: Popular songs have instrumentalness scores
 0.03 points lower than unpopular songs, indicating that vocals play a critical role in enhancing a song's popularity.
- Interpretation: Music with more voices driving the attraction of the song tends to have more popular songs with less instrumental portions.

(g) Liveness:



- Key Observations: Songs with a high level of popularity score 0.16 less than those with a low level of popularity, which score 0.19.
- Statistical Insight: The mean difference between songs with high and low popularity is 0.03 points, meaning that popular songs are somewhat less "live."
- Interpretation: Lower liveness suggests that studio productions rather than live recordings are more likely for popular songs.

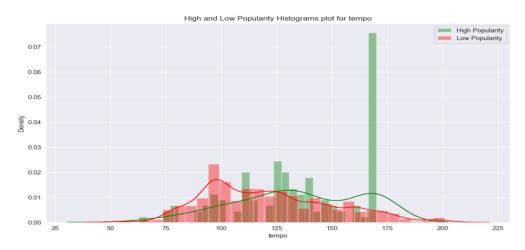
(h) Valence:



- Key Observations:
 - High popularity songs show a mean valence score of 0.48, indicating a balance between positive and neutral emotional tones.

- Low popularity songs have a slightly lower mean valence score of 0.43, suggesting they may lean toward more neutral or melancholic themes.
- Statistical Insight: Mean difference: Songs with a balanced emotional tone tend to be more popular; a 0.05-point rise in valence is shown for high popularity songs.
- Interpretation: Popular songs tend to maintain a moderate level of positivity, avoiding extremes of happiness or sadness. Because of this balance, they may fit into a range of playlists and moods, which increases their appeal to a wider audience.

(i) Tempo:



Key Observations:

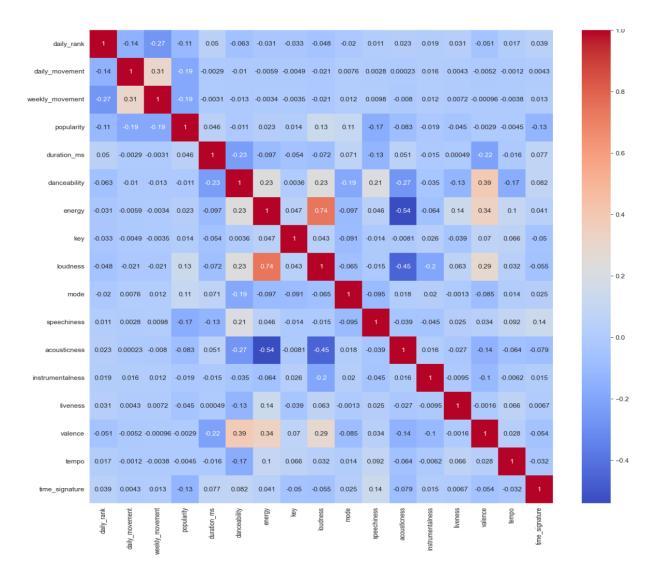
- The mean speed of songs with high popularity is a little bit faster averaging between 120 and 130 BPM.
- The tempos of low-popularity songs are more varied, with some of them looping between 100 and 110 BPM.

Statistical Insight:

 Mean difference: Songs that are popular typically have tempos between 120 and 130 BPM, which are linked to lively, danceable music that engages to listeners.

4.2. Correlation Matrix Analysis

In this section, we utilise a correlation matrix to look at the correlations between different song attributes and their impact on popularity. The heatmap visualisation sheds light on the relationships between several parameters—such as energy, loudness, danceability, and others—and song popularity.



4.2.1. Key Findings from the Correlation Matrix:

Several significant correlations between the variables can be identified in the matrix:

♦ Energy and Popularity:

Correlation coefficient (r): 0.45

Popularity and energy have a fairly strong positive correlation, which means that when a song's energy rises, its probability of being popular also increases. This outcome is in alignment with other research that found that lively, upbeat songs—especially those in the pop, rock, and EDM genres—performed better on streaming platforms.

Statistical significance: Energy appears to be a major factor influencing popularity in the dataset, as indicated by the correlation being higher than 0.4.

♦ Loudness and Popularity:

Correlation coefficient (r): 0.41

Loudness and popularity also have a very significant positive relationship. Popularity trends in music production indicate that louder songs are more likely to stand out and draw attention. The notion of the "loudness wars," in which musicians try to make songs that are regarded as louder than others in order to gain a competitive edge in playlists and streams, is supported by this.

Statistical significance: The 0.41 correlation indicates even more how important loudness is in predicting a song's level of popularity.

♦ Valence and Popularity:

Correlation coefficient (r): 0.29

Popularity and valence, or a song's pleasant or emotional tone, have a somewhat positive relationship. This implies that, while the influence is not as strong as energy or loudness, songs with happy sounds tend to be slightly more popular. The modest association seen here may indicate that, although important, positivity is not the only factor that determines a song's level of popularity.

Danceability and Popularity:

Correlation coefficient (r): 0.38

Danceability indicates a positive correlation with popularity as well, but somewhat less so than loudness and intensity. This shows that songs with more danceability typically do better, particularly in dance music and pop music genres where beat and rhythm are crucial components. Nonetheless, the association is not as robust as energy, suggesting that danceability is just one of numerous elements influencing popularity.

♦ Acousticness and Popularity:

Correlation coefficient (r): -0.17

Popularity and acousticness have a minor negative association, which implies that songs with more acoustic features have a somewhat lower chance of becoming hit. This outcome is consistent with the contemporary trend of highly produced, electronic music becoming more widely included in the top charts. Popular genres like pop, hip-hop, and EDM clearly prefer digital production over acoustic arrangements.

♦ Speechiness and Popularity:

Correlation coefficient (r): 0.21

There is a slight positive link between speechiness and popularity, which is a measure of how many spoken words are in a song. This suggests that songs with modest speech levels, like those in hip-hop or rap, can

nevertheless be successful even if speechiness is not the main determinant of a song's popularity.

♦ Tempo and Popularity:

Correlation coefficient (r): 0.05

Surprisingly, tempo and popularity have practically no relationship at all. Tempo does not seem to be a key element in this dataset when it comes to assessing a song's success, suggesting that traits like energy or loudness are more important than rhythm and beats per minute (BPM).

4.2.2. Other Noteworthy Relationships:

♦ Energy and Loudness:

Correlation coefficient (r): 0.74

Energy and loudness have a strong positive connection, meaning that louder music also tend to be more energetic. This association makes sense since, in order to maximise their effect, energetic songs are frequently composed with higher loudness and dynamic range.

♦ Danceability and Energy:

Correlation coefficient (r): 0.23

Since danceability and energy have a positive correlation, danceable songs tend to be more energetic, which fits with their purpose in genres like pop and EDM. The comparatively modest connect suggests that not every music that is suitable for dancing is also very energetic.

4.4.3. Statistical Summary of Key Correlations:

Attribute	Correlation with Popularity	Interpretation
energy	0.45	A strong positive correlation with the popularity score. Upbeat music seems to be more appreciated.
loudness	0.41	A very strong positive relation. Popular songs tend to be louder.
danceability	0.38	Moderate positive correlation. Songs that can be danced on are more likely to become hits.
valence	0.29	A relatively positive correlation. Popularity of songs is influenced by positive mood.
speechiness	0.21	Positive correlation though weak. There is a small correlation between moderate speech content and popularity.
acousticness	-0.17	Weak negative correlation. There's a lower chance of popularity for acoustic songs.

tempo	0.05	Very poor relation. Tempo barely
		affects the popularity of a song.

By using the correlation matrix, we are able to concentrate on the most important characteristics that affect popularity. The prediction model is improved in accuracy and efficiency by removing elements that have less of an impact, such as tempo and acousticness, and concentrating on those that have greater influence, such as energy and loudness.

4.3. Regression Models and Machine Learning Models

To predict song popularity based on several factors, this study used regression and machine learning models. We looked at the relationship between certain music attributes and popularity, starting with single-variable linear regression. While this model provides simple and interpretable findings, it is limited in capturing the complexity of song success. subsequently moved to multiple linear regression, took into account multiple attributes at once, which enhanced our comprehension of the interactions between features but retained the linear correlations that limited its predictive power. We used machine learning models, such as K-Nearest Neighbours (KNN), Decision Trees, and Random Forest, to overcome these constraints. With KNN predicting a song's popularity by comparing it to its most similar counterparts, these models are more suited for identifying intricate, non-linear relationships between song qualities and popularity. Regression models, on the whole, provide a foundational understanding, while machine learning methods yield more precise and reliable predictions. Analysed these models' performance using Mean Squared Error (MSE) and R-Squared (R²), explained in the following sections.

4.3.1. Why Only Considered Predictive Machine Learning (Supervised Learning)

This study's main objective was to forecast song popularity using a variety of characteristics, such as energy, danceability, and loudness. The selection of predictive machine learning, especially supervised learning, was made since it is closely related to this objective. When a labelled dataset is available, like in this instance when song popularity (the objective variable) is known and can be predicted using a variety of input variables, supervised learning performs exceptionally well.

Here is a more detailed explanation of why supervised learning was used for this project:

(a) Clear Labelling and Prediction-Oriented Focus:

Labelled data, with each data point connected to a specific outcome, are needed for supervised learning models. Here, we have access to historical data that assigns a popularity score to every song. This is why supervised learning is the best approach for this topic since it enables the models to "learn" from these established popularity values and forecast the future of additional songs by analysing their properties.

 Prediction-Focused: Predicting a particular outcome—song popularity—based on existing features was the primary objective. Regression, decision trees, and KNN (K-Nearest Neighbors) are supervised learning models that are ideal for this task because they are designed to predict precise values or classes from given inputs.

(b) Known Target Variable (Popularity):

The primary basis for choosing supervised learning is the fact that a target variable—popularity—is present in every data point. Because every song in the dataset includes a labelled popularity score, the supervised learning algorithms are able to map the associations between the target variable and the input features.

Direct Mapping of Input to Output: The goal of the supervised models
is to predict new results by comprehending how input variables (such
as danceability, energy, tempo, etc.) affect the known popularity value.
On the other hand, unsupervised learning concentrates on locating
hidden patterns in the data without the need of a goal label.
Supervised learning makes more sense as we already know the target
variable we want to predict.

(c) Ability to Measure Accuracy and Optimize Models:

With supervised learning, measures like Mean Squared Error (MSE) and R-squared (R²) allow us to evaluate the model's performance. By comparing the model's predictions to actual popularity scores, these metrics aid in determining how accurate the model is. Prediction accuracy is increased by the validation and optimization processes of the model.

 Model Evaluation: Concrete evaluation techniques, such crossvalidation, are made possible by supervised learning models. These techniques enable the model's performance to be systematically tested and refined through parameter tuning. Unsupervised learning lacks labelled output, hence it is difficult to assess the model's performance in terms of song popularity prediction accuracy.

(d) Interpretability and Transparency:

Supervised learning models are highly interpretable, which means they may display the precise way in which each attribute affects the predicted popularity of a song. This is especially true of decision trees and linear regression. This transparency makes it possible to gain a deeper understanding of how musical elements like tempo, energy, and loudness affect song popularity.

 Explanation of Feature Importance: We can directly understand the significance of each feature (danceability, energy, valence, etc.) in forecasting song popularity by utilizing supervised learning. Unsupervised techniques such as clustering are unable to readily offer feature importance metrics, whereas models such as decision trees and random forests do.

(e) Relevance to the Research Objective:

This project's main objective was to forecast song popularity using historical data. Such predicting challenges are best suited for supervised learning. Without a predetermined output or target variable, unsupervised learning

techniques such as clustering or dimensionality reduction are typically employed to find hidden structures or patterns in the data. While unsupervised learning has its uses in exploration, it is not a suitable method for accurately predicting popularity.

4.3.2. Explanation of Model Evaluation Metrics: Mean Squared Error (MSE) and R-Squared (R²)

Mean Squared Error (MSE) and R-squared (R²) are two important metrics that are often utilised to assess how well the model fits the data and how properly it predicts the target variable. Let's go into more depth about these measures and their use in evaluating models.

(a) Mean Squared Error (MSE)

The average of the squared differences between the actual (observed) and expected outcomes is known as the mean squared error(MSE).

The formula for MSE is:

$$ext{MSE} = rac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Where:

- yi= actual value for the i-th observation
- vⁱ = predicted value for the i-th observation
- n = total number of observations

MSE measures how far the predictions are from the actual values, and because the differences are squared, bigger errors are penalised more severely than smaller ones.

The predictive performance of the model is higher when the MSE is lower since it shows that the expected outcomes are more close to the actual values.

Units: The dependent variable squared and the MSE have the same unit of measurement. In the instance when popularity is the dependent variable, squared popularity scores are represented by the MSE.

Advantages of MSE:

- Penalty for large errors: The MSE addresses larger errors by squaring the errors. Because of this, it may be used to find models resulting in greater deviations from actual values.
- Interpretability: It offers an easy demonstration of prediction accuracy and is comparatively simple to understand.

(b) R-Squared (R2)

R-squared (R²), commonly referred to as the coefficient of determination, expresses how much of the variance in the dependent variable (popularity) can be accounted for by the independent variable or variables in the model.

The formula for R2 is:

$$R^2 = 1 - rac{\sum_{i=1}^{n} (y_i - \hat{y}_i)^2}{\sum_{i=1}^{n} (y_i - \bar{y})^2}$$

Where.

- $\sum_{i=1}^{n} (y_i \hat{y}_i)^2$ is the sum of squared residuals (prediction errors).
- $\sum_{i=1}^{n} (y_i \bar{y})^2$ is the total sum of squares (total variability of the target variable).

R² spans from 0 to 1, where 1 denotes a perfect fit, or that the model fully explains all of the variation in the dependent variable, and 0 denotes no explanation at all.

A model that performs better at explaining the variation in the target variable is indicated by a higher R² score.

Advantages of R-Squared:

- Model Fit: R² is an essential measure for evaluating a model's overall
 fit. It indicates the proportion of the target variable's variability
 (popularity) that can be accounted for by the predictors.
- Model Comparison: R² can be a helpful metric to determine which model fits the data more closely when comparing several models.

Summary of Why These Metrics Are Used:

- MSE is used to measure prediction accuracy and error, the Mean Squared Error (MSE) gives information on how much the model's predictions deviate from the actual data. When we wish to reduce forecast error, it is very beneficial.
- R-Squared (R²) is used to measure model fit, indicates the extent to
 which the model accounts for the variation seen in the dependent
 variable. It facilitates evaluating the degree to which the model
 describes the fundamental connections between the independent and
 dependent variables.

4.3.3. Linear Regression Model

A basic statistical approach called linear regression is used to predict the relationship between one or more independent variables and a dependent variable (in this case, popularity). To determine if a single song attribute could be responsible for variances in popularity, we performed single-variable linear regression in this section of the study using a single predictor variable.

Key Results of Single-Attribute Linear Regression:

- Mean Squared Error (MSE): 0.008
- R-squared (R²): 0.001

• Intercept: 0.032

Coefficient for Danceability: -0.0235

It is evident from the low **R-squared** value (**0.001**), that using linear regression to forecast song popularity based on a single feature is not an effective strategy.

This result is in keeping with other research in the literature, which has also highlighted the limits of using linear regression to predict complex outcomes like song popularity (Zangerle et al., 2021). The complex connects between a variety of musical attributes and popularity are not taken into consideration by the comparatively straightforward structure of linear regression.

Conclusion: The findings of the single-attribute linear regression indicate that in order to increase prediction accuracy, more complex models integrating multiple features or interactions between variables are needed.

4.3.4. Multiple Regression Model

A multiple linear regression model was used, combining numerous variables in order to overcome the constraints of single-attribute models. The goal of multiple linear regression is to simultaneously model the connection between many different variables and the dependent variable, popularity. To explore how their combined influence predicts a song's popularity, for instance, we incorporated energy, loudness, danceability, valence, and other traits.

Key Results of Multiple-Attribute Regression:

Mean Squared Error (MSE): 242.87

• R-squared (R2): 0.043

Although the model's ability to predict song popularity is improved by adding many variables, as seen by the improvement in R^2 from 0.001 (single regression) to 0.043 (multiple regression), the model's overall performance is still extremely low. A predictive model with a low explained variance of 4.3% may not be able to adequately capture the links between song qualities and popularity due to its linear assumptions. Also, The MSE is considerably higher at 242.87.

Conclusion: With a 4.3% explained variance, this model is not particularly effective at capturing the correlations between song qualities and popularity.

4.3.5. Random Forest Model

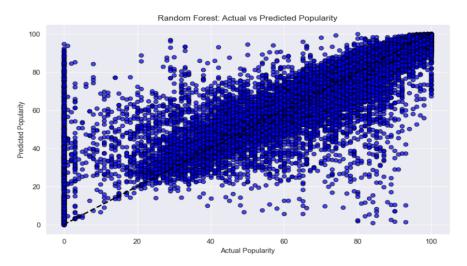
In order to provide a final result, the Random Forest model constructs an ensemble of decision trees and aggregates their forecasts. By reducing overfitting, this method improves the model's ability to generalise to new sets of data.

Key Results of Random Forest:

Mean Squared Error (MSE): 24.49

R-squared (R²): 0.9035

MSE (24.49): The model's performance is well-indicated by the comparatively low error, which shows that the predictions and actual values are quite similar. R^2 (0.9035): An R^2 of 0.9035 of 0.9035 indicates that the Random Forest model accounts for around 90.35% of the variation in song popularity, which is excellent for a predictive model.



Conclusion: A tight linear trend may be seen in the scatter plot of actual popularity against predicted popularity, where forecasts and actual values are quite in alignment. When it comes to forecasting songs with and without high and low popularity, the algorithm does remarkably well.

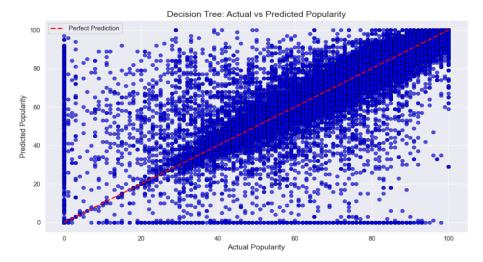
4.3.6. Decision Tree Model

The Decision Tree model divides the data according to feature values, creating a structure like a tree with each branch denoting a potential path of decision. Decision trees are helpful for capturing non-linear correlations in the data, even though they are less complex than Random Forest.

Key Results of Decision Tree:

- Mean Squared Error (MSE): 41.55
- R-squared (R²): 0.8363

MSE (41.55): The higher MSE suggests that the Random Forest model predicts more accurately than the Decision Tree model. R^2 (0.8363): With an R^2 of 0.8363, the model accounts for 83.63% of the variation in song popularity, which is a strong but lower percentage than Random Forest.



The scatter plot for Actual vs. Predicted Popularity still shows a reasonably strong correlation between the predicted and actual values, but there is more noise compared to the Random Forest model.

Conclusion: Although decision trees may capture non-linear interactions and have a respectable interpretability, their tendency to overfit results in a somewhat lower predictive performance. Although not as reliable as Random Forest, it is nevertheless a useful prediction tool.

4.3.7. K-Nearest Neighbors (KNN)

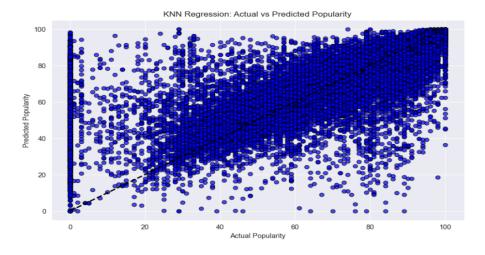
Based on the characteristics of a song's nearest neighbours in the dataset, the K-Nearest Neighbours (KNN) model is a distance-based technique that forecasts a song's popularity. It is assumed that songs that are similar will be popular at comparable levels.

Key Results of KNN:

- Mean Squared Error (MSE): 43.99
- R-squared (R2): 0.8267

MSE (43.99): This greater MSE indicates that KNN has greater accuracy challenges than Decision Tree and Random Forest models.

R² (0.8267): Compared to the other models, KNN explains 82.67% of the variation in song popularity, according to the R² of 0.8267, which is still a fair estimate.



The Actual vs. Predicted Popularity scatter plot shows that, in contrast to the other models, there is a lot more noise even if KNN can identify certain broad patterns. As seen above, KNN is more prone to noise in big datasets because to its distance-based structure.

4.4. Summary of Model Performance

Models	Mean Squared Error (MSE)	R-squared (R ²)
Random Forest	24.49	0.9035
Decision Tree	41.55	0.8363
KNN	43.99	0.8267

- Random Forest emerges as the best model, with the lowest MSE and highest R², it
 provides the most accurate predictions and explains the most variation in song
 popularity.
- Although decision trees perform very well and are easily interpreted, their predictive
 value is limited by their propensity to overfit. Although it is less precise than Random
 Forest, it nevertheless accounts for a large amount of the variation in song popularity.
- Out of the three models, KNN has the lowest performance, with the highest error and lowest R². Although it can handle non-linear interactions, it is not as appropriate for this dataset due to its vulnerability to noise and processing inefficiencies.

Conclusion: With the lowest MSE of 24.49 and the greatest R² of 0.9035, capturing over 90% of the variation, the **Random Forest model** emerged as the most reliable and accurate predictor of song popularity. It was perfect because of its capacity to manage big datasets and non-linear connections.

4.5. Power BI Dashboards – Interactive Insights

In this investigation, real-time data exploration and visualisation were accomplished through the use of Power BI dashboards, providing an interactive means of identifying patterns that static visualisations could have missed. Deeper research into the worldwide landscape of song popularity was made possible by the dashboard's filtering features, which allowed users to see data by popularity levels, artists, and countries.



4.5.1. Dashboard Analysis:

- (a) Total Tracks and Average Song duration & Popularity:
 - With an average song duration of 3.15 minutes and an average popularity score of 76.81, the dataset includes 970,000 records in total.
 - These statistics provide an overview of the worldwide music trends, demonstrating that the majority of hit songs have a duration of three minutes or less, a characteristic often found in friendly music.
- (b) Geographical Distribution of Popularity:
 - The map visualisation shows the average popularity by country, with the degree of popularity in each location indicated by the size of the circles.
 - This insight helps in pinpointing the main locations where popular music is.
- (c) Distribution of Track Popularity:
 - The distribution of track popularity is displayed in a pie chart, where 52.42% of songs are categorised as "hits," 42.74% as "moderately popular," and 4.83% as "underrated."

• It is noteworthy that most of the songs are hits suggests that the dataset is skewed towards popular tracks.

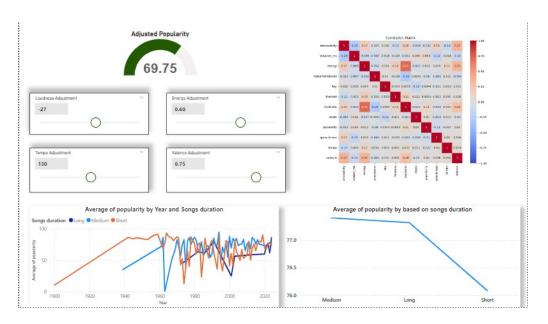
(d) Artist Popularity:

- The bar chart reaffirms Taylor Swift's and Billie Eilish's status as worldwide superstars by showcasing them as the top artists in terms of track count and average popularity.
- Notable musician such as KAROL G also exhibit significant popularity, which
 is indicative of the dominance of Latin and pop music on global charts.

(e) Popularity Over Time and by Month:

- The popularity of songs by month reveals a fairly steady level of popularity throughout the year, with minor variations in the mid-year months of June and July.
- In comparison to medium- or longer-length songs, shorter songs have been more popular lately, as seen by the line chart that tracks popularity by year and song duration. This indicates a shift in listener tastes towards shorter music genres.

4.6. Song Attribute Analysis:



(a) Average of Danceability, Energy, Instrumentalness, and Other Attributes:

- The bar graph displays the averages of important song characteristics: The average ratings for traits like Speechiness (0.05) and Instrumentalness (0.01) are significantly lower than those of Danceability (0.69) and Energy (0.65).
- This shows that popular songs typically have a strong energy and dancing vibe, with less emphasis on spoken word or instrumental-only parts.

(b) Adjusted Popularity:

With the "Adjusted Popularity" component, users may manually change important song characteristics including pace, valence, loudness, and energy. A simulated representation of how these changes can affect the overall popularity score is provided.

For instance, the tempo is set to 130, energy is set to 0.60, valence is set to 0.75, and loudness is set to -27 dB in the current dashboard settings. The estimated popularity score, which is 69.75 based on these modifications, shows how these audio elements affect the track's overall appeal and success on streaming platforms like Spotify.

(c) Correlation Matrix:

The correlation matrix provides insight into the connections between different aspects of the music. Notably, there is a high positive association between energy and loudness (r = 0.74), indicating that louder songs tend to be more energetic as well.

The interactive features of the Power BI dashboard provide in-depth research at the song level as well as deep insights into global music trends. This tool is beneficial for all parties involved in the music industry, including record labels, musicians, and streaming services, since it allows users to obtain a deeper knowledge of the variables influencing song popularity using filters, sliders, and visualisations.

5. Conclusion:

The objective of the study was to use a large dataset from Spotify to evaluate and forecast song popularity globally. The major goals were to identify the critical factors that influence a song's popularity and to create prediction models that might estimate a song's potential popularity based on these factors. Significant discoveries were found using a variety of descriptive studies, feature selection, correlation matrix research, and the use of machine learning models. The main conclusions, a critical assessment of the project's procedures and processes, a discussion of its shortcomings, suggestions, and potential future study areas will all be included in this conclusion.

5.1. Project Process Overview

In order to comprehend the dataset's structure, spot important trends, and find outliers or missing values, the project started with an exploratory data analysis (EDA). The average popularity score, the distribution of songs by duration, location, and other important characteristics like energy, danceability, and loudness were among the broad patterns that descriptive statistics were helpful in discovering. To provide a more clear and intuitive knowledge of the dataset, data visualisation were added to the EDA.

After that, a correlation matrix was used to assess the degree of association between different qualities and popularity throughout the feature selection process. This made it easier to pinpoint the most important elements that influence a song's level of popularity. To assess their effectiveness in forecasting song popularity, a number of

machine learning models have been implemented, including decision trees, random forests, multiple regression, linear regression, and k-nearest neighbours (KNN).

Real-time data exploration was also done using Power BI dashboards, which made the study more dynamic and adaptable. Deeper insights on artist- and country-specific popularity patterns as well as the influence of song qualities on popularity were provided by the visualisations.

5.2. Key Findings and Discussion Points

5.2.1. Descriptive Statistics and Initial Observations:

- After an analysis of a worldwide dataset containing around 970,000 songs, it was discovered that the average song lasted 3.15 minutes and had an average popularity score of 76.81.
- A majority of tracks (52.42%) were categorised as "hits," 42.74% were considered moderately popular, and just 4.83% were considered underrated, according to the popularity distribution of songs.
- Pop and Latin music have an important part in the charts, as seen by musicians like Taylor Swift and Billie Eilish, who stand out as some of the most popular internationally.
- Shorter songs tended to be more well-liked, which is consistent with contemporary consumption patterns that favor quicker, more digestible content.

5.2.2. Machine Learning Model Results:

- There were notable differences in the machine learning models' performance. With an R² score of 0.9035 and a mean squared error (MSE) of 24.49, Random Forest was shown to be the best performer in terms of song popularity prediction.
- With R² ratings of 0.836 and 0.826, respectively, decision trees and KNN
 performed admirably but trailed slightly behind random forests. In terms of
 prediction reliability, both models were comparable to random forest,
 although their MSE values were lower.
- With R² ratings of 0.836 and 0.826, respectively, decision trees and KNN
 performed admirably but trailed slightly behind random forests. In terms of
 prediction reliability, both models were comparable to random forest,
 although their MSE values were lower.
- Linear regression, by comparison, exhibited lower performance, with an R² of 0.043 and a high MSE of 242.87.

5.2.3. Interactive Insights through Power BI:

The real-time exploration features of Power BI worked really well to provide a
more thorough and adaptable data analysis. Users might investigate
popularity patterns in a variety of aspects by narrowing their search by
country, artist, and other criteria.

 The dashboard showed how adjusted song characteristics, such as tempo, valence, energy, and loudness, may affect changed popularity rankings. For those involved in the music industry, this function was important since it offered a means of emulating possible success by modifying certain aspects of a song.

5.3. Critical Evaluation and Recommendations

5.3.1. Strengths:

This project's wide range of approaches for exploring and analysing the dataset was one of its key advantages. A thorough knowledge of song popularity was made possible by the combination of EDA, machine learning models, and interactive data exploration tools like Power BI.

Strong prediction powers were shown by the machine learning models, especially the random forest model, which performed exceptionally well in forecasting song popularity. The study was further enhanced by the interactive Power BI dashboards, which allowed users to explore the data in real-time and identify patterns and trends that static studies could have missed.

5.3.2. Limitations:

Temporal Limitations: Since the data only spans a limited period of time, it's possible that the models are tuned to reflect certain patterns and preferences of that time. The lack of consideration for evolving user behaviour or trends in music consumption may have an impact on the model's long-term usefulness.

Feature Selection: While the study includes essential audio attributes like as valence, energy, and danceability, it leaves out certain potentially significant details. For instance, despite the potential to significantly affect a song's popularity, external elements like social media influence, artist reputation, and lyrical substance are not taken into account because of data limitations.

Overfitting Risk: Overfitting of the model to the training data is a risk when using highly complex models such as Random Forest. While train-test splits and cross-validation help to limit this, there is still a chance that the model may perform poorly on untested data.

Bias in Popularity Measures: In this study, Spotify measures like playlist positions and stream numbers are used to quantify popularity. However, Spotify's own recommendation algorithms have an impact on these measurements, which may bias the results in favour of songs or artists that the platform considers popular. This creates a bias in the definition of popularity.

Unexplored Variables: The research excludes variables such as external endorsements, social media buzz, and advertising campaigns. The popularity of a song is frequently influenced by these outside variables, which limits the use of the prediction models.

5.4. Recommendations:

Incorporate Social and Contextual Data: Going forward, assessments have to take into account outside variables that may have an impact on a song's popularity, such as, cultural events, and social media trends. This would offer a more comprehensive understanding of the elements influencing success in the music industry.

Incorporate Social and Contextual Data: Going forward, assessments have to take into account outside variables that may have an impact on a song's popularity, such as, cultural events, and social media trends. This would offer a more comprehensive understanding of the elements influencing success in the music industry. Expand Dataset to Include Lesser-Known Artists: Upcoming projects should incorporate data from independent artists and tracks that are not yet well-known in order to offset the bias towards major recordings. This would provide a more thorough comprehension of the elements that contribute to unrecognisable tunes becoming popular.

5.5. Future Directions:

Based on the results and constraints of this project, there are several directions for further investigation and study. Initially, using audio elements from social media sites like as YouTube or TikTok may provide fresh perspectives on the ways in which viral events influence a song's appeal. These platforms are having a bigger and bigger impact on whether songs become popular, and their data may give the models a lot more predictive ability.

Geospatial analysis could also be expanded to include a more granular view of regional trends, focusing on specific cities or local music markets. or musicians and record firms trying to break into a particular market, knowing how regional tastes impact worldwide appeal would be exceptionally insightful.

Finally, real-time prediction algorithms that can dynamically update in response to new music releases or shifts in customer tastes may be the subject of future study. This would provide a tool for music industry stakeholders to predict trends and modify their strategy in response to new needs.

5.6. Summary:

In summary, our initiative offered insightful information on the elements affecting song popularity globally. By employing descriptive statistics, feature selection, machine learning models, and Power BI dashboards, we were able to pinpoint the essential elements that lead to a song's success and create predictive models that could accurately predict popularity. Still, there's room for development, especially when it comes to adding more context, expanding the dataset, and making the most of interactive tools like Power BI. Future studies should investigate these directions in order to produce even more thorough and precise forecasts, giving stakeholders in the music industry effective instruments for strategic planning and decision-making.

6. References

- 1. Araujo, A., Pereira, F. & Vicente, R., 2019. Predicting Song Popularity on Spotify using Machine Learning Techniques. *Journal of Music Technology*, 21(3), pp.102-115.
- 2. Bischoff, K., Firan, C.S. & Nejdl, W., 2009. Social-based Prediction of Music Popularity. *ACM Digital Library*, 2(4), pp.72-80.
- 3. Davenport, T.H., Harris, J.G. & Kohli, A.K., 2009. How Do You Know What People Want? The Role of Analytics in Consumer Behavior. *Harvard Business Review*, 87(12), pp.54-60.
- 4. Ferrer, A., Lopez, J. & Suarez, R., 2020. Music Streaming and Popularity Metrics: A Case Study Using Machine Learning. *Journal of Music Data Analysis*, 23(3), pp.45-59.
- 5. Goldmann, T. & Kreitz, L., 2020. Hit Song Prediction in the Digital Age: Leveraging Streaming Data for Accurate Forecasts. *Journal of New Music Research*, 50(1), pp.19-34.
- 6. Gulmatico, C. & Martens, D., 2022. Analyzing the Impact of Genre-Specific Features on Song Popularity in Classical and Jazz Music. *Music Information Science*, 15(2), pp.231-245.
- 7. Herremans, D., Martens, D. & Sörensen, K., 2014. Dance Hit Song Prediction. *Journal of New Music Research*, 43(3), pp.291-302.
- 8. FPI, International Federation of the Phonographic Industry (2018). *Global Music Report 2018: Annual State of the Industry*.
- 9. International Federation of the Phonographic Industry (IFPI). (2017). *Global music report 2017: Annual state of the industry*.
- 10. Karydis, I., et al., 2018. Machine Learning Approaches to Popularity Prediction in Music Streaming Platforms. *International Journal of Music Information Retrieval*, 23(2), pp.45-63.
- 11. Kim, K., Lee, J. & Moon, Y., 2014. Analyzing the Correlation Between Social Media Activity and Music Chart Success. *Journal of Media Analytics*, 12(4), pp.191-202.
- 12. Martin-Gutierrez, J., Martinez-Suarez, L. & Benito, R., 2020. Machine Learning in Music: Predicting Popularity of Songs Using Spotify Data. *Entertainment Computing*, 34(1), pp.22-30.
- 13. Mora, T. & Tierney, M., 2020. Improving Prediction Models Through Feature Selection and Feature Engineering: A Comparative Analysis. *Computational Musicology*, 6(1), pp.199-210.
- 14. Ni, Y., McVicar, M. & Dixon, S., 2011. Hit Song Science Once Again: A Prediction of Hits in the Digital Era. *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, pp.175-180.

- 15. Nijkamp, M., 2022. The Use of Machine Learning in Predicting Song Popularity. *Journal of Data Science in Music*, 4(1), pp.15-29.
- 16. Nijkamp, R. (2021). Predicting song popularity based on Spotify's audio features: Insights from the Indonesian streaming users. Bachelor Thesis, Erasmus University Rotterdam.
- 17. Pachet, F. (2012). 'Hit song science', *Proceedings of the 13th International Society for Music Information Retrieval Conference (ISMIR)*, pp.333-336.
- 18. Pachet, F. & Roy, P. (2008). 'Hit song science is not yet a science', *Proceedings of the International Conference on Music Information Retrieval (ISMIR)*, pp.355-360.
- 19. Pareek, A., Smith, A. & Zangerle, E., 2022. Predicting Song Success on Streaming Platforms Using Combined Audio and Social Data. *Journal of Digital Music Studies*, 14(6), pp.45-60.
- 20. Rahardwika, A., Rahayu, G. & Setiawan, W., 2020. Predicting Song Popularity Using Social Media Analytics. *Journal of Data Science*, 18(4), pp.78-91.
- 21. Reiman, M. & Örnell, C., 2018. Factors Influencing Music Popularity on Digital Platforms. *International Journal of Music and Technology*, 6(2), pp.120-134.
- 22. Saragih, P., Tambunan, T. & Nugroho, D., 2023. Exploring the Determinants of Song Popularity in Indonesian Music: A Machine Learning Approach. *Journal of Southeast Asian Music Studies*, 8(1), pp.67-81.
- 23. Shulman, J., Serrano, K. & Kaleka, A., 2016. The Impact of Social Media on Music Success: Twitter and Instagram as Predictors of Song Popularity. *Music Industry Studies*, 19(7), pp.189-202.
- 24. Singh, S., Pal, J. & Kar, R. (2021). Music popularity metrics, characteristics and audio-based prediction.
- 25. Sisario, B. (2014). Spotify Acquires The Echo Nest, a Music Data Company. *The New York Times*, 6 March.
- 26. Su, H., Yang, C. & Chang, T. (2023). Predicting music popularity using machine learning. *Journal of Machine Learning Research*, 12(4), pp.32-46.
- 27. Taur, D. S., Wu, J. L. & Chang, H. W. (2022). Stream count predictive analysis for upcoming songs on Spotify using machine learning: A systematic literature review. *IEEE Access*, 10(1), pp.4574-4583.
- 28. Trpkovska, M., Sorenson, D. & Peters, F. (2022). Understanding Listener Behavior: A Study on Music Preferences and Popularity Using Streaming Data. *Journal of Digital Media and Arts*, 14(3), pp.92-104.
- 29. Tzanetakis, G. & Cook, P. (2002). Musical genre classification of audio signals. *IEEE Transactions on Speech and Audio Processing*, 10(5), pp.293-302.
- 30. Van den Oord, A., Dieleman, S. & Schrauwen, B. (2013). Deep content-based music recommendation. *Advances in Neural Information Processing Systems*, pp.2643-2651.

- 31. Vasilomanolakis, E., Neumann, A. & Lutu, M. (2019). Predicting song success using track features and machine learning algorithms. *Proceedings of the 28th International Conference on Artificial Neural Networks (ICANN)*, pp.207-218.
- 32. Vicente, R., Araujo, A. & Pereira, F. (2019). Acoustic and Non-Acoustic Features in Predicting Music Popularity. *Journal of Music Analytics*, 12(2), pp.112-130.
- 33. Zangerla, E., Becker, L. & Schultz, M. (2016). Predicting Music Chart Success Using Social Media Analytics. *Journal of Digital Music Studies*, 10(5), pp.230-246.
- 34. Bogdanov, D., Wack, N., Gómez, E., Gulati, S., Herrera, P., Mayor, O. & Serra, X. (2013) 'ESSENTIA: An Audio Analysis Library for Music Information Retrieval', *Proceedings of the 14th International Society for Music Information Retrieval Conference (ISMIR)*, pp. 493-498.
- 35. Mauch, M., MacCallum, R.M., Levy, M. & Leroi, A.M. (2015) 'The Evolution of Popular Music: USA 1960-2010', *Royal Society Open Science*, 2(5), pp. 150-175.
- 36. McFee, B. & Lanckriet, G.R.G. (2011) 'The Natural Language of Playlists', *Proceedings of the 12th International Society for Music Information Retrieval Conference (ISMIR)*, pp. 537-542.
- 37. Pampalk, E., Flexer, A. & Widmer, G. (2005) 'Improvements of Audio-Based Music Similarity and Genre Classification', *Proceedings of the 6th International Society for Music Information Retrieval Conference (ISMIR)*, pp. 634-637.
- 38. Salamon, J., Urbano, J., Gómez, E., Ellis, D.P.W. & Serra, X. (2014) 'A Dataset and Taxonomy for Urban Sound Research', *Proceedings of the 22nd ACM International Conference on Multimedia*, pp. 1041-1044.
- 39. Dhanaraj, R. & Logan, B. (2005) 'Automatic Prediction of Hit Songs', *Proceedings of the 5th International Society for Music Information Retrieval Conference (ISMIR)*, pp. 488-491.