

Hit Songs Through the Decades

A Machine Learning Approach to Predicting Release Periods and Analyzing Trends in Popular Music

Andrea Lovato

Elwin Freudiger

Maya Wahart

2025-05-17

Table of contents

1	Introduction	2
2	Literature review	3
2.1	Music Analytics and Machine Learning	3
2.2	Temporal Trends and Musical Evolution	3
2.3	Challenges in Modeling Popularity and Period	3
2.4	Contribution of the Present Study	4
3	Data	4
3.1	Sources	4
3.2	Description	5
3.3	Wrangling / Cleaning	6
3.4	Image Data	7
3.5	Spotting Mistakes and Missing Data	7
3.6	Listing Anomalies and Outliers	8
3.6.1	Interpretation	9
3.7	Summary statistics	10
3.8	Correlation Matrix	10
4	Exploratory Data Analysis	10
4.1	Correlation matrix	11
4.2	Variables analysis	12
4.2.1	Release year	12
4.2.2	Track numbers	13
4.2.3	Track Duration	15
4.2.4	Explicit	18

4.2.5	Popularity	22
4.3	Textual analysis	23
5	Machine Learning Methods	26
5.0.1	Supervised Learning Methods	27
5.0.2	Unsupervised Learning Methods	27
6	Results	28
6.1	Supervised Learning	28
6.1.1	Random Forest with undersampling	28
6.1.2	Random Forest on median year	30
6.1.3	Neural Networks	32
6.1.4	Conclusion	34
6.2	Unsupervised Learning	34
6.2.1	Principal Component Analysis	35
6.2.2	K-means clustering	37
7	Conclusion	41
8	Appendix	43
8.1	Hierarchical clustering	43
9	References	43

1 Introduction

The primary objective of this project is to apply machine learning techniques to a dataset of 10,000 top-charting songs from the ARIA and Billboard rankings, spanning from the 1950s to 2024. Using audio and metadata features provided by Spotify, the project aims to explore the evolution of musical trends and understand the key characteristics that define popular music across different decades.

From a machine learning perspective, the project focuses on both supervised and unsupervised learning tasks. This includes building predictive models to estimate a song's release period based on its features, and using clustering algorithms to uncover latent patterns in the data. The goal is not only to analyze historical music trends but also to assess the feasibility and accuracy of predictive modeling in the context of music analytics.

2 Literature review

2.1 Music Analytics and Machine Learning

The intersection of music and data science has increasingly attracted scholarly attention, particularly due to the rise of streaming platforms that provide large-scale, structured musical data. Several studies have explored the potential of audio and metadata features—such as those provided by Spotify—for understanding musical trends, classifying genres, and predicting popularity (Schedl et al., 2015). These features typically include danceability, energy, tempo, and valence, which encode perceptual and structural dimensions of sound and have been shown to correlate with human musical preferences (Friberg et al., 2011).

Machine learning techniques have been applied extensively to music data for tasks such as mood classification, genre detection (Kim et al., 2010), and hit song prediction (Herremans et al., 2019). Supervised models, particularly random forests and support vector machines, have demonstrated robust performance in tasks involving feature-based classification and regression. More recent work has applied deep learning architectures to raw audio input, though these approaches often demand substantial computational resources and training data (Choi et al., 2017).

Unsupervised learning, including clustering and dimensionality reduction techniques such as PCA, has also been used to explore latent patterns in musical corpora. For example, Jan Van Balen et al. (2015) used clustering to discover prototypical musical structures, while Serrà et al. (2012) applied PCA to uncover dominant stylistic trends over time.

2.2 Temporal Trends and Musical Evolution

Temporal analyses have shown that the characteristics of popular music evolve in response to technological, cultural, and economic forces (Mauch et al., 2015). For instance, changes in average song duration have been linked to shifts in radio programming, physical formats (e.g., vinyl, CDs), and digital streaming incentives. The increasing prevalence of explicit lyrics has been attributed to looser content restrictions and changing cultural norms (Pachet, 2008).

Empirical work by Interiano et al. (2018) used Spotify audio features to track the emotional content of popular music, observing a long-term trend toward increased sadness and decreased acousticness in top-charting songs. Similarly, it has been demonstrated that temporal audio descriptors could be used to model the release decade of songs with reasonable accuracy.

2.3 Challenges in Modeling Popularity and Period

While popularity is a central variable in music analytics, its measurement remains opaque. Spotify's proprietary popularity score is influenced by recent streaming activity, skips, and

playlist placements, making it a dynamic and platform-specific metric (Spotify, 2023). Consequently, its use in predictive models must be interpreted as a temporal snapshot rather than a static attribute.

Predicting a song's release year or period based on its acoustic profile poses challenges due to the high variability of musical styles within the same time frame and the enduring popularity of certain older songs. Nevertheless, several studies have demonstrated that audio features carry enough temporal signal to support classification tasks by decade or era, particularly when modeling broader stylistic shifts (Müller et al., 2010).

2.4 Contribution of the Present Study

Building on this body of literature, the present study contributes a comprehensive analysis of 10,000 top-charting songs from the ARIA and Billboard rankings, spanning more than seven decades. While prior work has focused either on genre classification or mood prediction, this study uniquely combines exploratory data analysis, unsupervised learning (e.g., PCA, clustering), and supervised modeling to assess both the temporal evolution of music and the predictive power of audio features.

Moreover, by investigating the relationships between explicitness, duration, and popularity in a temporal context, this work highlights how platform incentives and listener behaviors have shaped recent musical trends. It also addresses gaps in the literature regarding the feasibility of using readily available Spotify features for temporal classification, and evaluates the limitations of clustering methods in capturing stylistic boundaries.

Finally, this study serves as a practical application of machine learning techniques in the music domain, offering insights for musicologists, data scientists, and digital media analysts interested in the computational modeling of cultural data.

3 Data

3.1 Sources

The dataset employed in this study, titled “*Top 10,000 Spotify Songs – ARIA and Billboard Charts*”, was obtained from Kaggle, a widely used platform for sharing datasets and data science resources. It comprises a curated collection of 10,000 tracks that have achieved significant popularity, based on historical rankings from both the ARIA (Australian Recording Industry Association) and Billboard charts. This dual-source approach ensures a broad and balanced representation of commercially successful music across English-speaking markets.

The dataset spans a temporal range from the 1950s to 2024, capturing the dynamic evolution of popular music over more than seven decades. It includes metadata and audio-based

features extracted via the Spotify API, such as tempo, energy, danceability, valence, and instrumentalness, which allow for detailed computational analysis.

In addition to representing a variety of genres, artists, and time periods, the dataset reflects shifting cultural and musical preferences. As such, it provides a robust foundation for both exploratory data analysis and machine learning applications, particularly those aimed at uncovering temporal trends, predicting historical context (e.g., release period), and identifying latent patterns in music characteristics.

3.2 Description

The dataset comprises 10,000 entries and 35 variables, encompassing information related to song popularity, artist identity, release date, and various musical attributes.

Variable	Description	Category	Example
track_uri	Unique identifier for the track	character	spotify:track:123...
track_name	Name of the track	character	Bohemian Rhapsody
artist_uris	URIs of artists performing the track	character	spotify:artist:abc...
artist_names	Name(s) of the performing artist(s)	character	Queen
album_uri	Unique identifier for the album	character	spotify:album:def...
album_name	Title of the album	character	A Night at the Opera
album_artist_uris	URIs of the album's main artist(s)	character	spotify:artist:abc...
album_artist_names	Name(s) of the album's main artist(s)	character	Queen
release_date	Date the album was released	date	1975-11-21
album_image_url	Link to album cover image	character	https://i.scdn.co/image/...
disc_number	Disc number of the track in multi-disc sets	numeric	1
track_number	Track's position on the disc	numeric	11
duration_ms	Length of the track in milliseconds	numeric	354000
preview_url	URL to 30-second preview of the track	character	https://p.scdn.co/mp3-preview/...
is_explicit	Indicates if track has explicit content	logical	TRUE
popularity	Spotify popularity score (0-100)	integer	85
isrc	International Standard Recording Code	character	GBUM71029604
added_by	User who added the track to playlist	character	user_id_123

Variable	Description	Category	Example
added_at	Timestamp when track was added	datetime	2022-07-15T12:00:00Z
artist_genres	Genres associated with the artist(s)	character	rock, classic rock
danceability	How suitable a track is for dancing	numeric	0.6
energy	Intensity and activity level of the track	numeric	0.85
key	Musical key of the track (0=C, 1=G, ..., 11=A)	integer	5
loudness	Overall loudness in decibels	numeric	-5.3
mode	Modality: major (1) or minor (0)	integer	1
speechiness	Presence of spoken words in the track	numeric	0.05
acousticness	Confidence that track is acoustic	numeric	0.02
instrumentalness	Likelihood that track is instrumental	numeric	0.001
liveness	Likelihood of live audience presence	numeric	0.09
valence	Musical positiveness conveyed	numeric	0.7
tempo	Beats per minute (BPM)	numeric	120.5
time_signature	Estimated time signature	integer	4
album_genres	Genres associated with the album	character	rock, progressive rock
label	Record label	character	EMI
copyrights	Copyright info for the album or track	character	© 1975 Queen Productions Ltd.

3.3 Wrangling / Cleaning

In preparation for analysis, the original dataset was cleaned and standardized by renaming variables to follow consistent, machine-readable naming conventions. Redundant or non-essential columns—such as URIs, preview links, and metadata unrelated to audio features—were subsequently removed. The resulting dataset *spotify_vr* retains only the relevant musical, temporal, and popularity-related attributes needed for the subsequent exploratory and predictive modeling tasks.

The resulting dataset consists of 23 columns, containing only the variables relevant for the subsequent analysis after the removal of redundant and non-informative features.

We only keep the ID from the Spotify track URL by extracting the final component of the *track_uri* string.

3.4 Image Data

In order to enhance the predictions, album covers were selected to be added as input in the neural network model. The orginal dataset uses a link provided by spotify to retrieve each image.

All images are retrieved from their respective links, resized to be 64x64 pixels. The size was chosen to balance concerns of space and usefulness. Indeed, a 16x16 pixel image would provide hardly any useful information. On the other hand, a 600x600 image would provide plenty of information, but would be difficult to store. All images are then saved locally, with their Spotify Song ID as a filename. This process was performed in parallel using multi-treading. This enables a much faster processing time.

Image examples can be seen below:



Figure 1: Taylor Swift's Red Album cover resized to 64x64 pixels and represented with RGB color channels

3.5 Spotting Mistakes and Missing Data

The procedure involves identifying and counting missing values in the dataset, detecting rows containing incomplete information, and removing those with missing or empty loudness values. Then, the release year is extracted from the release date and converted to a numeric format, producing a cleaned dataset ready for further analysis.

The dataset have a total of 625 missing values before cleaning. After cleaning we still have 547 missing values, all of them concern the *artist_genres*.

Table 2: Rows with Missing or Blank Values

	x
track_uri	0
track_name	0
artist_names	0
album_name	0
release_date	0
album_image_url	0

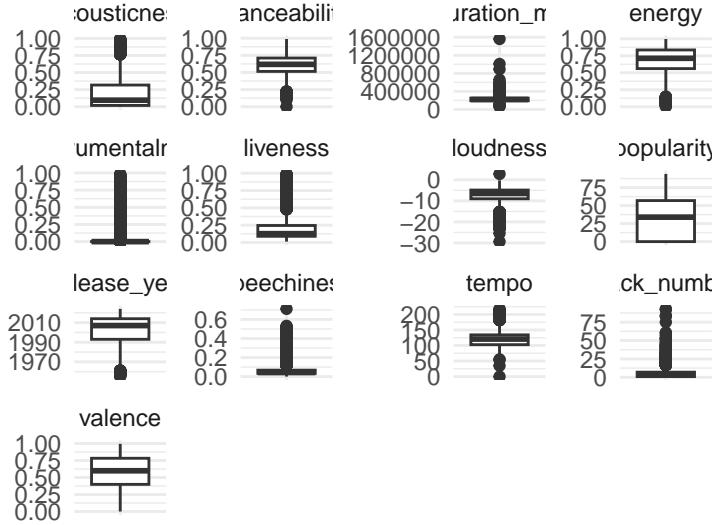
track_number	0
duration_ms	0
is_explicit	0
popularity	0
artist_genres	547
danceability	0
energy	0
loudness	0
mode	0
speechiness	0
acousticness	0
instrumentalness	0
liveness	0
valence	0

The dataset contains a total of 625 missing values, of which 551 correspond to the *artist_genres* variable. Instances with missing values in variables related to musical characteristics are excluded from the analysis to ensure data quality and consistency.

3.6 Listing Anomalies and Outliers

A subset of the dataset containing only numeric variables was selected, excluding *mode* and *time_signature*. The data was then transformed into long format to generate faceted boxplots, allowing for the visualization of the distribution and scale of each numeric variable individually.

Boxplots for Numeric Variables (Individual)



3.6.1 Interpretation

The faceted boxplots illustrate the distribution of key numeric features. Variables such as acousticness, energy, danceability, instrumentalness, liveness, valence, and speechiness are bounded between 0 and 1. Most exhibit distributions concentrated near zero, particularly instrumentalness and speechiness, which show long right tails and outliers close to 1, suggesting that while the average track lacks strong instrumental or spoken elements, some are highly characterized by them.

Duration (in milliseconds) is slightly right-skewed, with most tracks clustered around the median and a few outliers representing exceptionally long songs. Loudness is centered around negative values, consistent with its decibel scale relative to silence, and shows a compact distribution with occasional extreme lows, likely corresponding to quiet or highly dynamic tracks.

Popularity displays a broad distribution with outliers at both ends, indicating notable variability in audience reception. Tempo shows high variance and several extreme values, potentially due to anomalous entries or experimental compositions. Lastly, track number is typically low—reflecting songs positioned early in albums—though outliers suggest the presence of long compilations or inconsistent metadata.

3.7 Summary statistics

Summary statistics are generated for all numeric variables in the cleaned dataset to provide an overview of central tendencies, dispersion, and the presence of potential outliers.

track_number	duration_ms	popularity	danceability	energy	loudness
Min. : 1.000	Min. : 60093	Min. : 0.00	Min. :0.0000	Min. :2.03e-05	Min. :-29.368
1st Qu.: 1.000	1st Qu.: 192103	1st Qu.: 0.00	1st Qu.:0.5150	1st Qu.:5.61e-01	1st Qu.: -9.020
Median : 3.000	Median : 219426	Median :34.00	Median :0.6180	Median :7.13e-01	Median : -6.496
Mean : 4.938	Mean : 224245	Mean :32.56	Mean :0.6083	Mean :6.84e-01	Mean : -7.252
3rd Qu.: 7.000	3rd Qu.: 249826	3rd Qu.:57.00	3rd Qu.:0.7100	3rd Qu.:8.36e-01	3rd Qu.: -4.878
Max. :93.000	Max. :1561133	Max. :94.00	Max. :0.9880	Max. :9.97e-01	Max. : 2.769

The summary statistics table provides an overview of the central tendencies and dispersion of key variables in the dataset.

As we can see, regarding the release date, most songs have been released in more recent years. this may lead to unbalanced data. Ways to solve this unbalance will be discussed later in this paper. Regarding the track number in the song, the mean being at *4.9*, this tells us that most songs may part of an album or compilation. interestingly, the max of song number is *93* for a song named *Soul Revival* by *Johnny Diesel & The Injectors* part of a *Complete Eighties* compilation of 100 songs from the 80s.

The track duration is expressed in milliseconds, with a mean of *3 minutes and 44 seconds*. The longest song in record is a whopping 26 minutes. for *Tubular Bells - Pt. I* by *Mike Oldfield*. While this may not ring a bell (pun intended) for most readers, Amateurs of Horror may recognize this as the opening soundtrack for *The Exorcist (1973)*.

The *is_explicit* variable is highly imbalanced, with approximately 95% of the songs labeled as non-explicit. The distribution of the popularity variable is not normally distributed, with a mean value of 33, indicating that most songs fall within a lower popularity range. Variables such as danceability, energy, key, loudness, mode, speechiness, acousticness, instrumentality, liveness, valence, tempo, and time signature are audio features provided by Spotify that describe various musical characteristics of each track.

3.8 Correlation Matrix

4 Exploratory Data Analysis

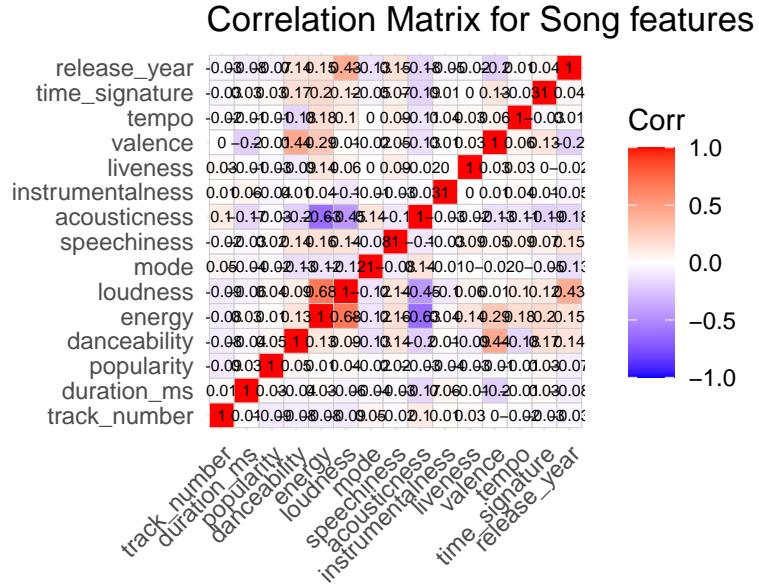
To gain a comprehensive understanding of the dataset and uncover patterns relevant to the modeling phase, an exploratory data analysis (EDA) is conducted. This phase involves the

systematic examination of the dataset's structure, distributions, and relationships among variables. We begin with summary statistics and visual inspections of key numeric features, followed by the analysis of correlations and potential multicollinearity. Subsequently, we investigate temporal trends, track characteristics, and the distribution of categorical variables such as explicit content and genre. This stepwise exploration helps identify data quality issues, potential outliers, and underlying trends that may influence or inform the subsequent application of machine learning techniques.

4.1 Correlation matrix

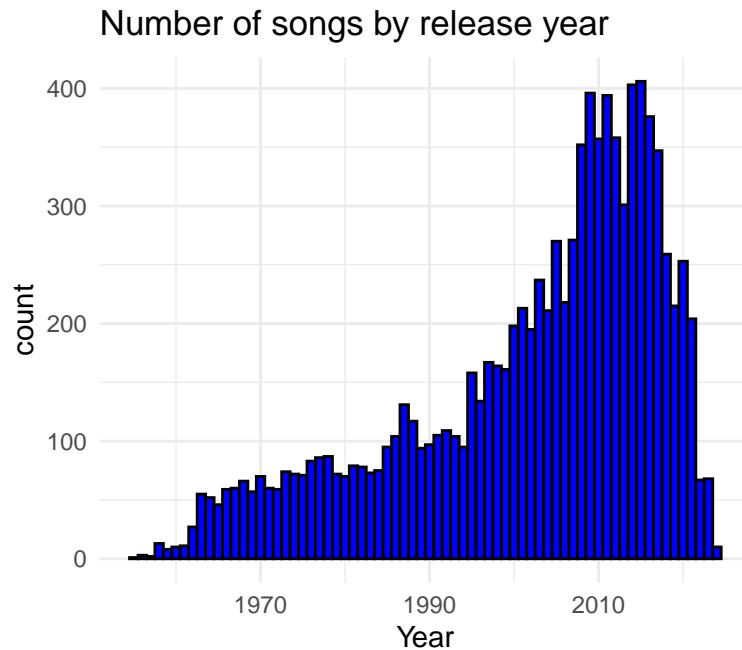
A correlation matrix is computed and visualized using a color-coded upper-triangle plot to identify linear relationships among the numeric variables in the cleaned dataset. This helps reveal patterns of association and potential multicollinearity between audio features.

The correlation matrix reveals a moderate positive association between energy and loudness, suggesting that more energetic tracks tend to be louder. A notable negative correlation is observed between acousticness and energy, indicating that acoustic songs generally exhibit lower energy levels. Overall, the absence of strong correlations among most variables suggests low multicollinearity, supporting their joint inclusion in multivariate analyses.



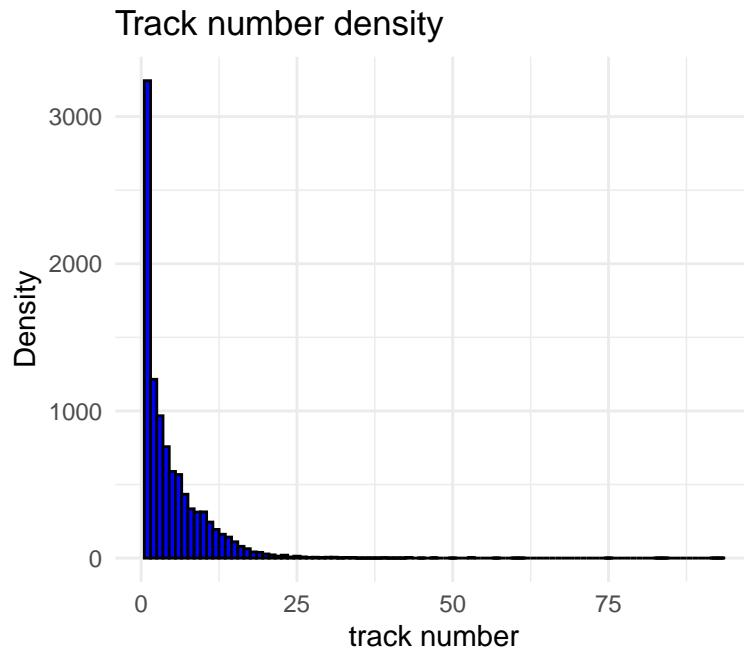
4.2 Variables analysis

4.2.1 Release year



Release year is the dependent variable in this analysis, and it is considered at the level of yearly granularity. The histogram shows a strong concentration of songs in more recent years, indicating a temporal imbalance that may affect the representativeness of earlier decades in the modeling phase.

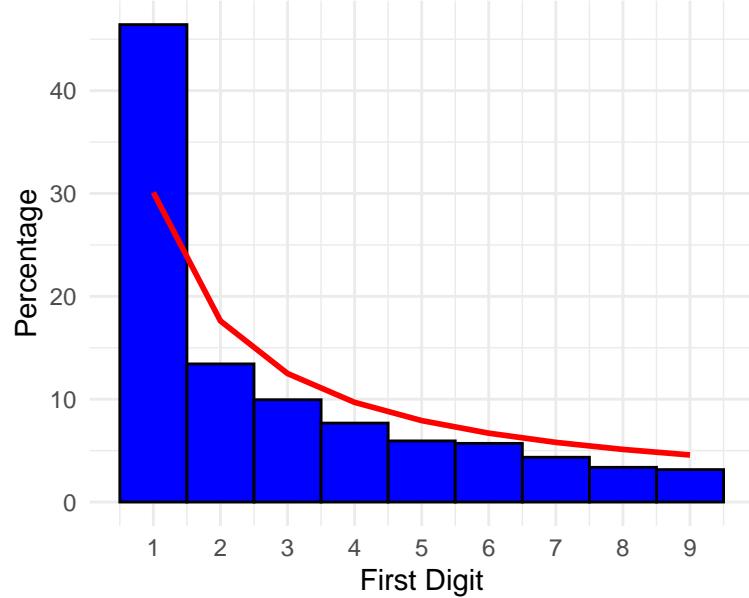
4.2.2 Track numbers



The distribution of track numbers is highly skewed toward lower values, indicating that most songs appear early in albums, while higher values likely reflect compilations or large track-lists.

It may be interesting to assess whether Benford's Law applies to the distribution of track numbers.

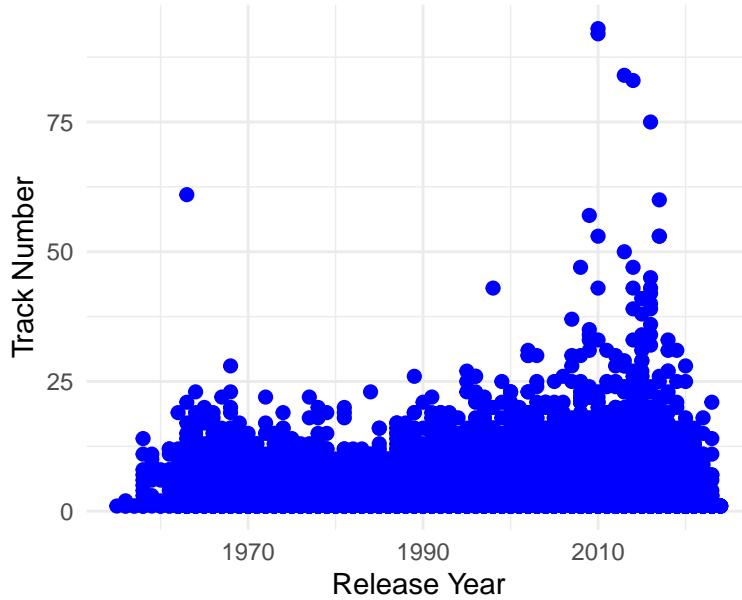
First Digit Distribution of Track Numbers vs E



The distribution of the first digit in track numbers clearly deviates from the expected pattern described by Benford’s Law. This discrepancy is likely due to structural constraints in how music albums are organized. Most releases—such as singles, EPs, and standard albums—contain a relatively small and fixed number of tracks, typically ranging from 1 to 15. As a result, lower digits, particularly 1, dominate the distribution, not because of a naturally logarithmic phenomenon but due to intentional sequencing and formatting practices in album production. This illustrates how domain-specific conventions can override general statistical laws in structured datasets.

Lastly, the evolution of track numbers over time can be examined to assess whether album structure or track positioning has changed across decades, potentially reflecting shifts in music consumption formats or production practices.

Track number depending on the release year

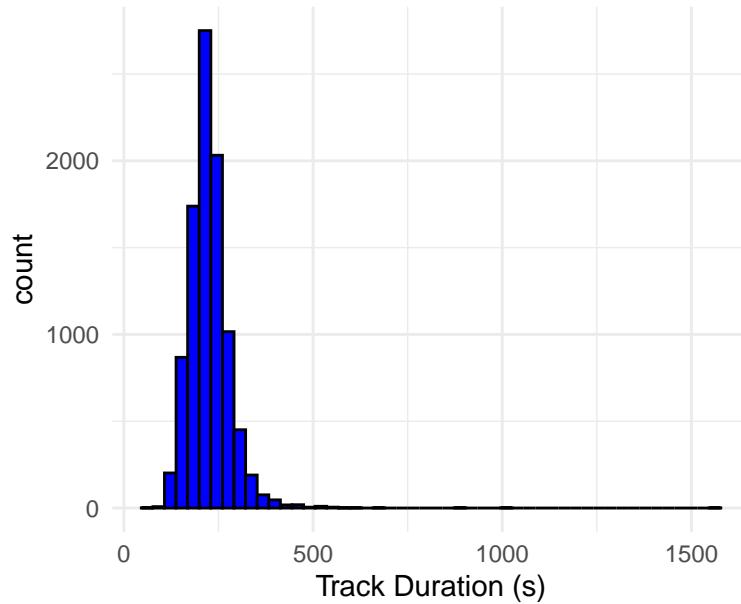


The scatter plot shows no clear trend in the evolution of track numbers over time. While the majority of songs consistently appear in early album positions across decades, some outliers—particularly in more recent years—exceed typical album lengths, likely reflecting special editions, compilations, or digital releases with extended tracklists.

4.2.3 Track Duration

The distribution of track duration, expressed in milliseconds, is examined to determine typical song lengths and to identify potential outliers, including exceptionally short or long tracks, which may affect subsequent analyses.

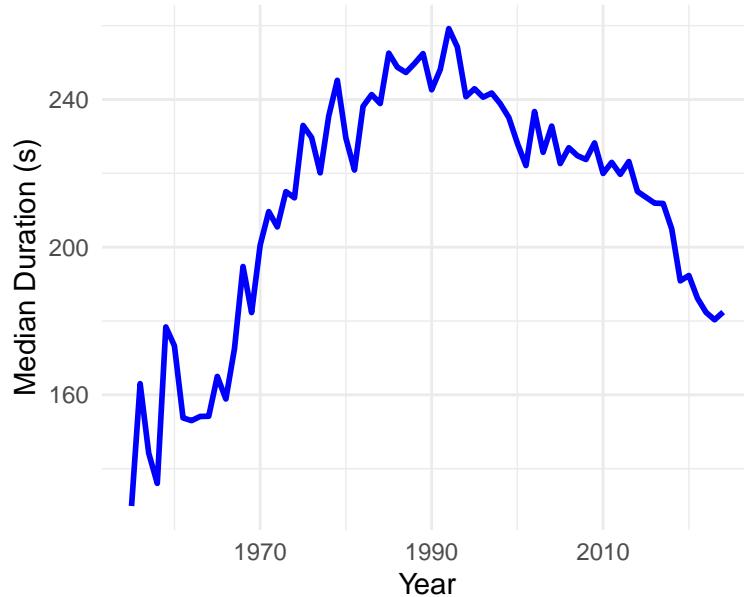
Histogram of track durations



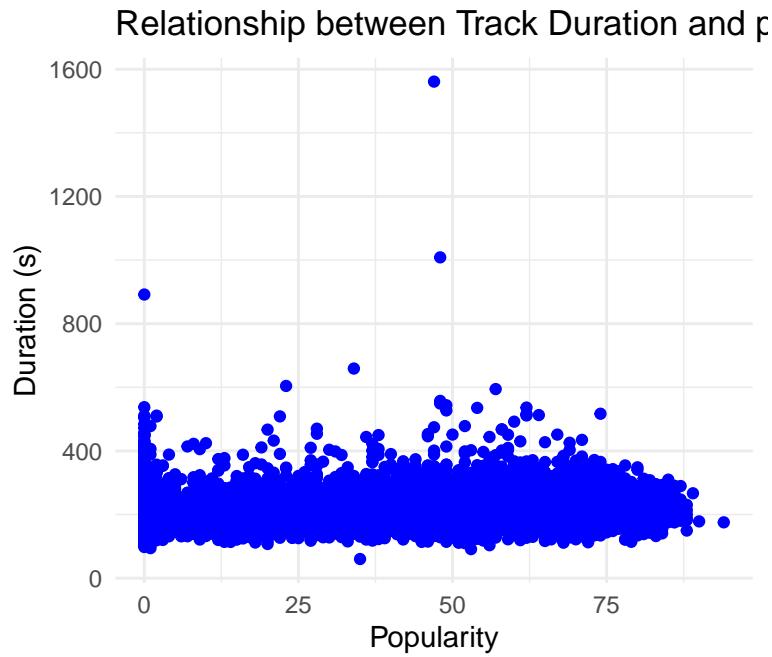
The histogram shows that the distribution of track durations is right-skewed, with most songs clustered around typical lengths and a limited number of extreme values representing unusually long tracks.

A line plot is used to visualize the evolution of median track duration over time. This allows identification of long-term trends, such as shifts in typical song lengths across decades, and highlights notable changes in production or consumption patterns.

Evolution of median track duration by year



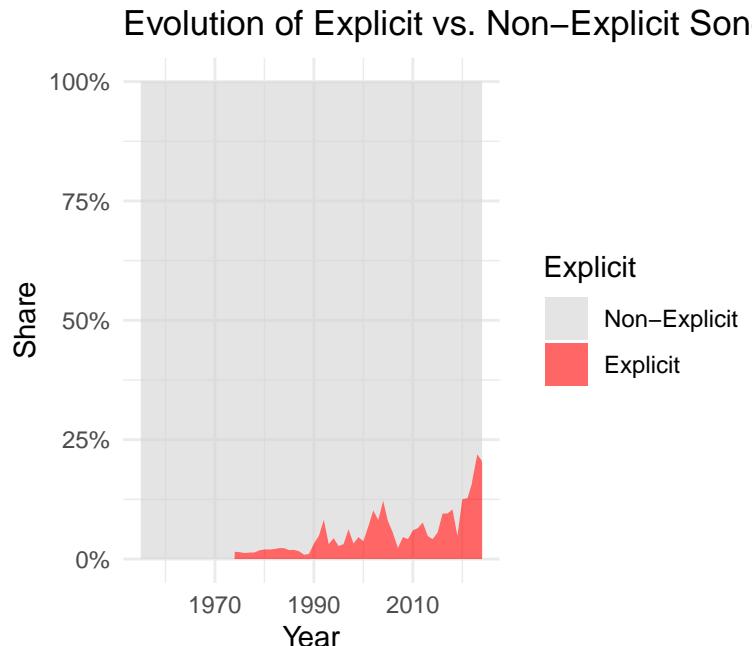
The plot shows that median track duration increased steadily throughout the second half of the 20th century, reaching a peak in 1992 at approximately 260 seconds (4 minutes and 20 seconds). A gradual decline follows, particularly from the 2010s onward. This recent downward trend is likely influenced by changes in digital consumption patterns, where streaming platforms incentivize shorter songs due to skip behavior and payout structures based on play counts.



The scatter plot shows no strong correlation between duration and popularity, though popular songs tend to have standard lengths, while very long tracks are generally less popular.

4.2.4 Explicit

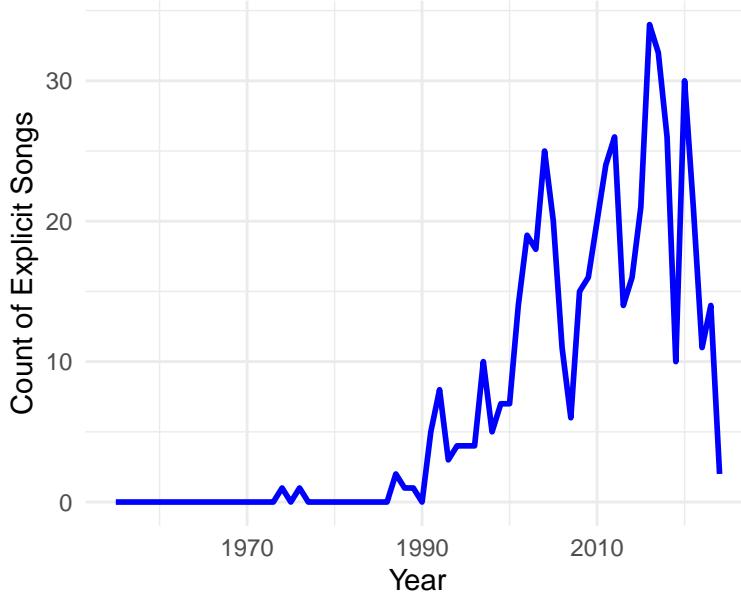
The proportion of explicit versus non-explicit songs is analyzed over time to observe how the prevalence of explicit content has evolved across different release years.



The analysis reveals a clear upward trend in the share of explicit songs over time, with a marked increase beginning in the early 2000s. This suggests a shift in lyrical content or labeling practices in the streaming era, where explicit content has become more common in popular music releases.

A line plot is now used to visualize the annual count of explicit songs over time, offering a more detailed view of their increasing presence in recent decades.

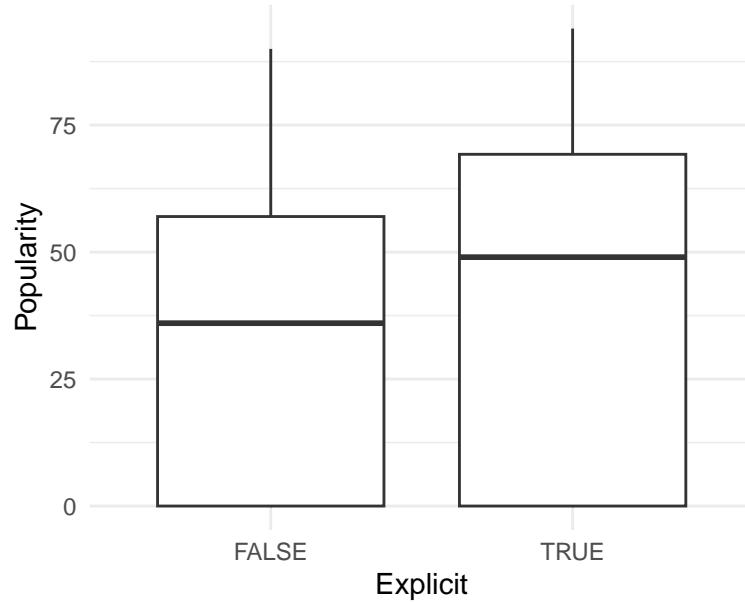
Number of Explicit Songs Over Time



The line plot confirms a rising trend in the number of explicit songs over time, particularly from the early 2000s onward. However, this pattern should be interpreted cautiously, as the explicit label is assigned by content uploaders and may not consistently reflect the presence of explicit material. In some cases, tracks with potentially explicit content are also released in censored versions to ensure broader distribution, such as radio play.

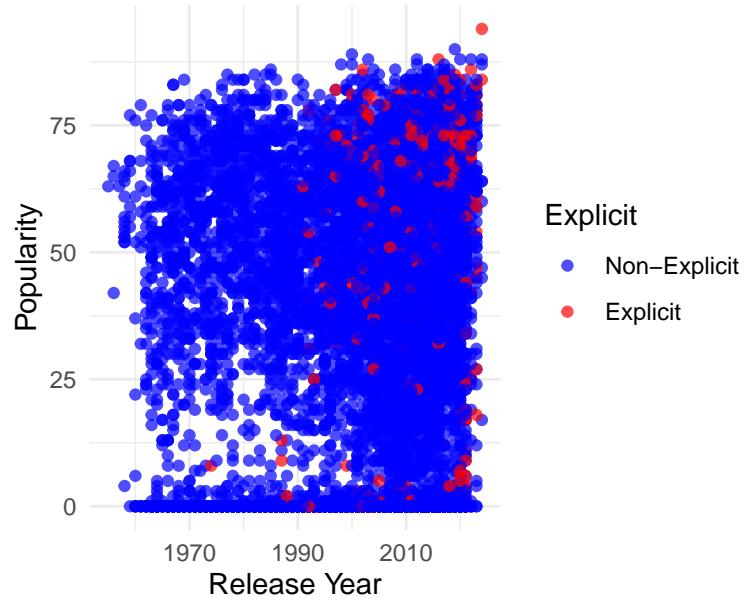
A boxplot is used to compare the distribution of song popularity between explicit and non-explicit tracks, allowing for the assessment of whether explicit content is associated with higher or lower popularity levels.

Popularity by Explicit Status



To further explore the relationship between explicit content and popularity, a scatter plot is used to visualize popularity scores over time, differentiated by explicit status. This approach helps assess whether the observed higher popularity of explicit songs is inherently tied to their content or instead influenced by temporal release patterns.

Popularity by year with explicit status



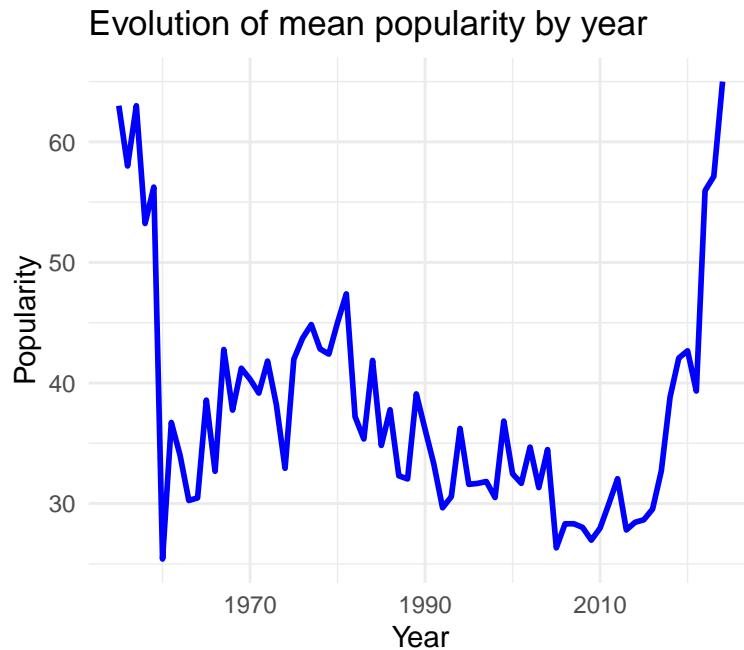
The scatter plot reveals that explicit songs (in red) are predominantly concentrated in more recent years and often exhibit high popularity scores. This supports the hypothesis that the observed popularity advantage of explicit tracks may be confounded by release date, as newer songs are both more likely to be explicit and to benefit from recency effects in popularity metrics.

4.2.5 Popularity

Popularity serves as a valuable proxy for a song's commercial success and audience reach. While the exact algorithm used to calculate Spotify's popularity score is not publicly disclosed, it is reasonable to assume that it reflects streaming frequency over a defined recent time window.

Track Name	Artist Names	Popularity
Espresso	Sabrina Carpenter	94
Cruel Summer	Taylor Swift	90
Yellow	Coldplay	89
Lose Control	Teddy Swims	89
Starboy	The Weeknd, Daft Punk	88
When I Was Your Man	Bruno Mars	88

It is important to note that the dataset was last updated in October 2024, meaning that the popularity scores reflect a specific point in time and may have since changed. However, this temporal limitation does not affect the validity of the present analysis, as the values still offer a reliable snapshot for exploring general trends and patterns in music popularity.

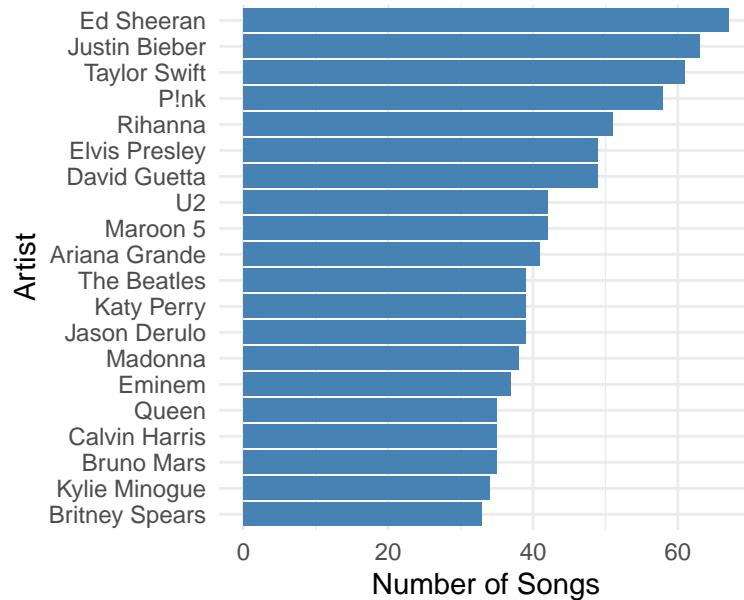


The chart shows peaks in popularity for early and recent songs. While the high popularity of recent songs is to be expected, the elevated scores for tracks from the 1960s are more surprising. A possible explanation is that artists featured in rankings and playlists remain culturally significant. Artists such as Frank Sinatra, Elvis Presley or Johnny Cash have become timeless icons that continue to be cultural relevant and popular across generations.

4.3 Textual analysis

This section performs a basic text mining task by extracting and counting individual artist names from multi-artist tracks. The goal is to identify the top 20 most frequently appearing artists in the dataset.

Top 20 Artists by Song Count



This step identifies the top 10 artists with the highest average popularity scores, offering insights into which artists consistently produce well-received songs within the dataset.

Artist Name	Average Popularity
Benson Boone	87.0
a-ha	86.0
Teddy Swims	85.5
Alphaville	85.0
Nayer	85.0
Jung Kook	84.0
Nate Ruess	84.0
Tyla	84.0
D-Block Europe	83.0
cassö	83.0

The following analysis ranks artists by their average track duration, highlighting those associated with longer musical compositions. Such patterns may reflect stylistic tendencies common in certain genres, including progressive rock, ambient, or experimental music.

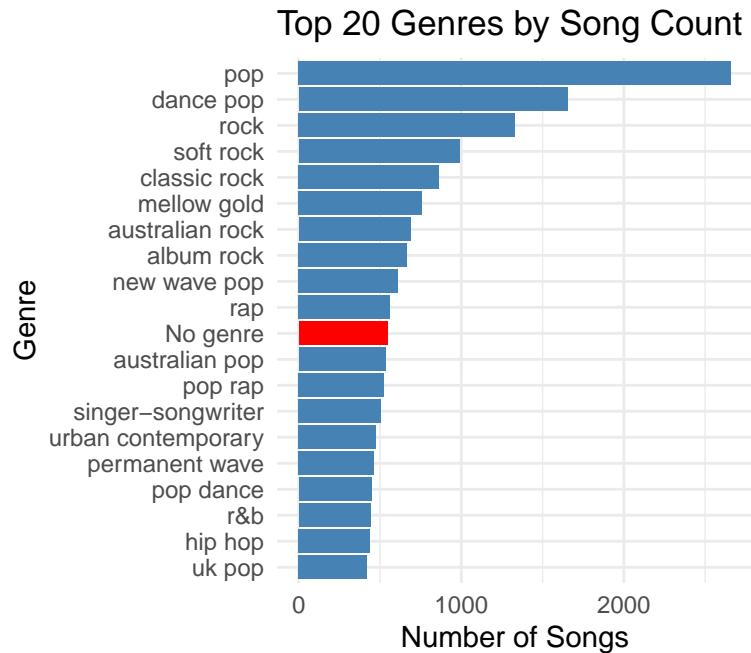
Artist Name	Avg. Song duration
The Bucketheads	891720

Mike Oldfield	633756
The Stone Roses	594453
The S.O.S Band	549000
Joe Walsh	536226
Curtis Mayfield	535333
Faithless	527240
Rollo Armstrong	527240
Sister Bliss	527240
Supercharge	523004

The following analysis focuses on the artists with the lowest average track durations, aiming to reveal trends related to brevity in musical production, which may be influenced by genre conventions or platform-driven listening behaviors.

Artist Name	Avg. Song duration
Liam Lynch	91226
Maurice Williams & The Zodiacs	97506
Fakebitcheshero	103026
The Swinging Blue Jeans	106066
Jerry Lee Lewis	111535
Hank Ballard & The Midnighters	111826
The Clovers	112200
Paul Russell	114233
Clyde McPhatter	116240
The String-A-Longs	116360

A similar type of analysis is now performed based on musical genre, in order to examine how average song characteristics vary across different stylistic categories.



Genres	Average Popularity
deep dance pop	85
melodic drill	83
indie rock italiano	82
italian pop	82
birmingham grime	80
indie r&b	79
bubblegrunge	78
sad lo-fi	78
sad rap	78
float house	77

5 Machine Learning Methods

In this study, we aim to predict the release year of a song based on its musical characteristics by applying and comparing several machine learning methods. Both supervised and unsupervised learning techniques are employed to explore the data structure and improve prediction performance.

5.0.1 Supervised Learning Methods

Supervised learning involves building models that learn from labeled data, where the outcome variable (in our case, the song's year) is known. The goal is to find a function that best predicts this outcome from a set of features.

5.0.1.1 Random Forest

Random Forest (RF) is an ensemble learning method that builds a large number of decision trees during training and outputs either the mode (for classification) or the mean (for regression) of the individual tree predictions. It combines two key ideas: bagging (training trees on bootstrapped samples of the data) and random feature selection at each split, which helps reduce variance and improve model generalization. In this context, RF will be used to predict the release year of a song as a regression task, leveraging its robustness to overfitting and capacity to handle complex interactions between features.

5.0.1.2 Neural Networks

Artificial Neural Networks (NN) are flexible, layered models composed of nodes (neurons) that compute weighted sums of inputs and apply non-linear activation functions. They are particularly effective for capturing complex, non-linear relationships in data. In our study, we use NNs with one or more hidden layers to model the relationship between musical features and the year of release, treating this as a regression problem. The model is trained using gradient descent to minimize a loss function, typically the mean squared error (MSE).

5.0.2 Unsupervised Learning Methods

Unsupervised learning does not use labeled data but instead focuses on uncovering hidden structures within the data. This can be particularly useful in exploratory phases to detect patterns and validate assumptions about the features.

5.0.2.1 K-means Clustering

K-means is a partitioning method that groups data into a predefined number of clusters (k) by minimizing the within-cluster sum of squares. The algorithm iteratively assigns each instance to the nearest cluster centroid and updates the centroids until convergence. In our analysis, K-means is used to examine whether natural groupings exist among songs based on their musical attributes, potentially revealing eras or genres that influence the release year.

5.0.2.2 Principal Component Analysis (PCA)

PCA is a dimensionality reduction technique that transforms the original features into a smaller number of uncorrelated components that retain most of the data variance. In this study, PCA is used to reduce the high-dimensional musical feature space into a more manageable form for visualization and clustering. Additionally, PCA components may be used as inputs to supervised models to mitigate multicollinearity and improve performance.

5.0.2.3 Hierarchical Clustering

Hierarchical clustering builds a tree-like structure (dendrogram) of nested clusters based on the similarity between data points. We apply the agglomerative approach, where each song starts as its own cluster, and the closest clusters are merged step-by-step. This method helps visualize relationships between songs and investigate whether certain musical profiles correspond to specific periods or styles.

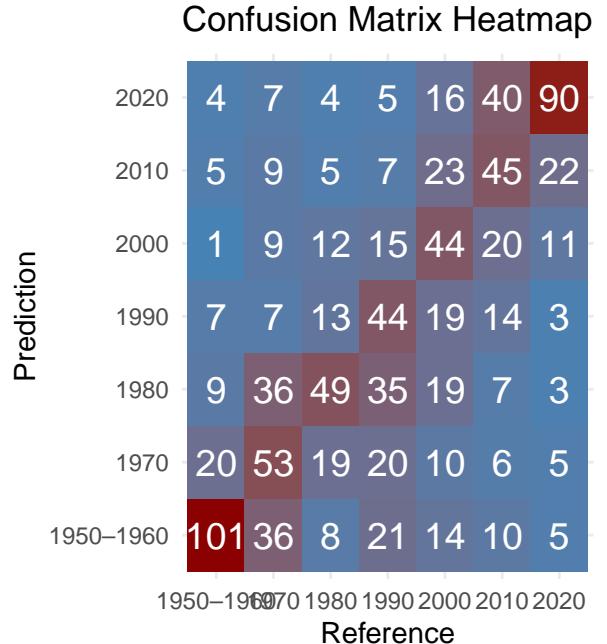
6 Results

6.1 Supervised Learning

6.1.1 Random Forest with undersampling

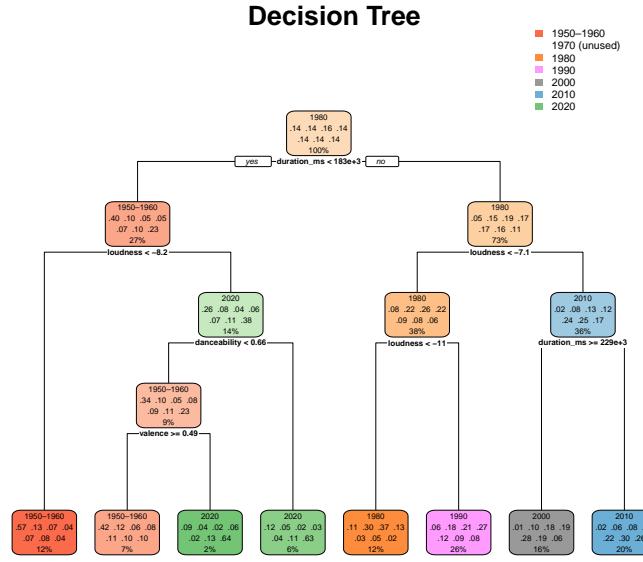
For the first supervised model, we implemented a Random Forest classifier to predict the decade of release for each song. To address the issue of unbalanced classes, we grouped songs from the 1950s and 1960s into a single category and applied the under-sampling technique, resulting in balanced groups of 470 songs per decade.

decade_group	n
1950–1960	470
1970	470
1980	470
1990	470
2000	470
2010	470
2020	470



The Random Forest model achieved an accuracy of 41.0%. The Cohen's Kappa was 0.312, reflecting only fair agreement beyond chance, and highlighting limitations in distinguishing between decades. However, the result is statistically significant ($p < 2e-16$), confirming that the model performed better than random guessing.

An analysis of the model's sensitivity and specificity revealed that it performs best in identifying songs from the 1950–1960 and 2020 periods, with sensitivities of 72.8% and 51.3% respectively. These decades likely exhibit the most distinct musical traits. In contrast, the model struggled with songs from the 1990s and 2000s, which had the lowest sensitivities (21.6% and 29.1%), possibly due to stylistic continuity across adjacent decades. Specificity, however, remained high across all decades, with values above 86% and peaking at 94.2% for the 2020 class.



The decision tree from the Random Forest revealed that features such as duration_ms, loudness, and danceability were crucial for classifying songs by decade. Older songs (1950–1960) clustered around lower loudness and shorter durations, while modern songs (post-2000) displayed higher loudness and more variability in mood (valence). The decision paths reflect a gradual stylistic evolution, supporting the hypothesis that production and mastering techniques, as well as compositional trends, shift over time.

6.1.2 Random Forest on median year

To improve predictive performance, we reformulated the problem as a binary classification task, aiming to predict whether a song was released before or after 2007, the median year in the dataset. This simplification helped the model focus on broader stylistic trends across time rather than granular decade differences.

Confusion Matrix and Statistics

Reference

Prediction after before

after	678	213
before	280	828

Accuracy : 0.7534

```
95% CI : (0.7339, 0.7721)
No Information Rate : 0.5208
P-Value [Acc > NIR] : < 2.2e-16

Kappa : 0.5045

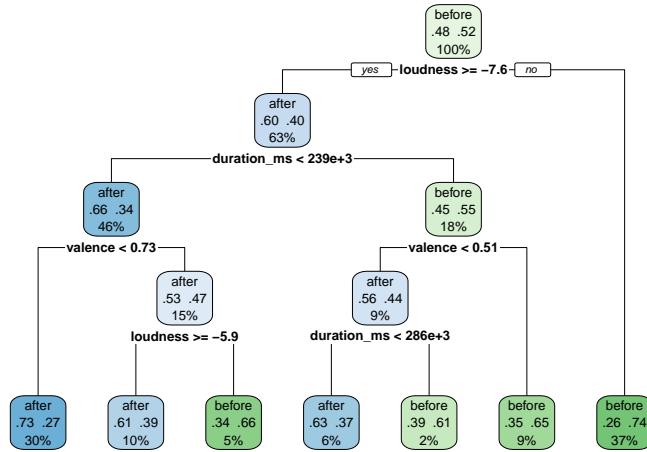
McNemar's Test P-Value : 0.002954

Sensitivity : 0.7077
Specificity : 0.7954
Pos Pred Value : 0.7609
Neg Pred Value : 0.7473
Prevalence : 0.4792
Detection Rate : 0.3392
Detection Prevalence : 0.4457
Balanced Accuracy : 0.7516

'Positive' Class : after
```

The Random Forest model achieved an accuracy of 75.3%, with a balanced accuracy of 75.2% and a Cohen's Kappa of 0.504, indicating moderate agreement beyond chance. It performed well across both classes, with a sensitivity of 70.7% (correctly identifying post-2007 songs) and a specificity of 79.5% (correctly identifying pre-2007 songs). The model's performance was statistically significant ($p < 2.2e-16$).

Decision Tree



The decision tree revealed that features such as loudness, valence, and duration_ms were crucial for distinguishing songs released before and after 2007. Songs released before 2007 tended to have lower loudness levels and lower valence, reflecting older production techniques and more subdued emotional tones. In contrast, post-2007 songs were often characterized by shorter durations, higher valence, and greater variation in energy, aligning with modern trends in pop music shaped by the digital and streaming era.

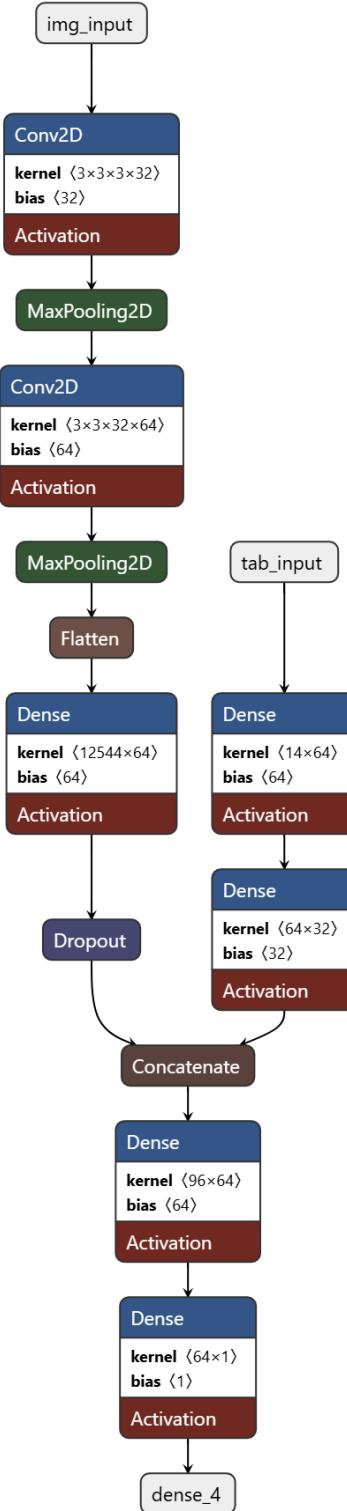
6.1.3 Neural Networks

While Random Forest models showed satisfying results when it came to classification of track decades, attempts involving regression tasks to predict the year saw issues. It was therefore decided to attempt to fit a neural network to predict the release year using tabular data as well as album covers.

6.1.3.1 Model Building

The model was built and fitted using python and the **Keras** library, allowing to easily build neural networks.

The decision was made to train both inputs, the image data in a CNN and the tabular data in a DNN. the layers were then merged. into one hidden layer thereby taking in all the information from both layers and finally making a prediction.



The model summary can be seen below:

To build the model, the release year was normalized, Therefore enabling predictions bounded between 1955 and 2024, the minimum and maximum release years. The final output layer uses a sigmoid function, making sure that forecasts are always bounded and that the model will not predict impossible years such as 2026 for example.

The model was trained with 20 epochs and a batch size of 32 returned the following metrics:

Training MSE: 0.0274 Training MAE: 0.1133

Test MSE: 0.0270 Test MAE; 0.1139

Test R² score: 0.4539

These results show no sign of overfitting, regarding the metrics themselves, one should keep in mind that these metrics are for normalized years. Therefore, a mean absolute error of 0.1139 over the 69(2024-1955) years of study means that the model is on average off by 7.8 years. The R-squared score shows that while the model may capture some of the variance, a big part is due to other factors.

6.1.4 Conclusion

To conclude on supervised learning, the random forest models show that musical features hold some predictive value for estimating a song's release period, though with limitations. In the multi-class model, predicting the exact decade proved challenging, with overlapping styles between adjacent decades affecting accuracy. Nonetheless, features like loudness, valence, and danceability captured broad production trends over time.

The binary classification model, distinguishing songs released before or after 2007, performed considerably better. This suggests that while musical features may not pinpoint exact release dates, they are effective in capturing wider historical shifts in music production and composition.

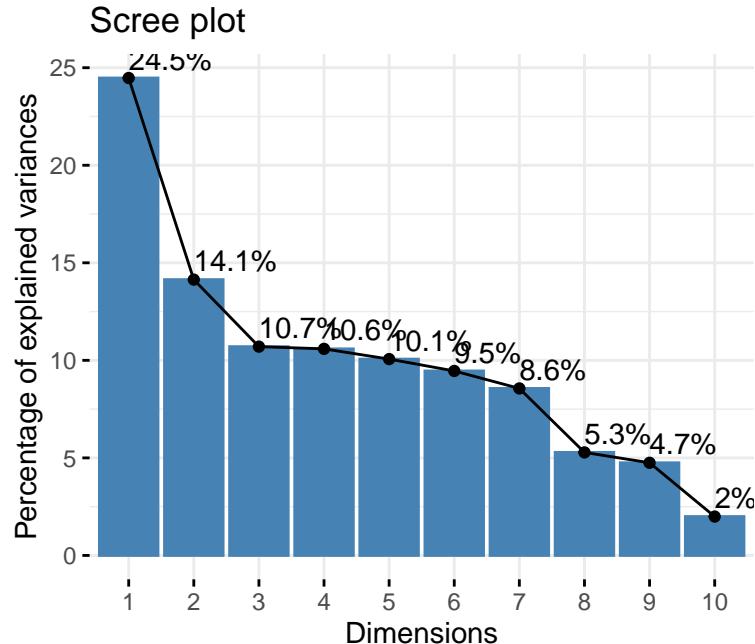
finally, the neural network model leveraging both album cover images and tabular data showed no signs of overfitting, its performance remained moderate. The average prediction error of approximately 7.8 years suggests that musical features and visual aesthetics alone are not sufficient to precisely predict a song's release year. Nonetheless, the model does capture some temporal patterns, indicating potential when combined with richer contextual data.

6.2 Unsupervised Learning

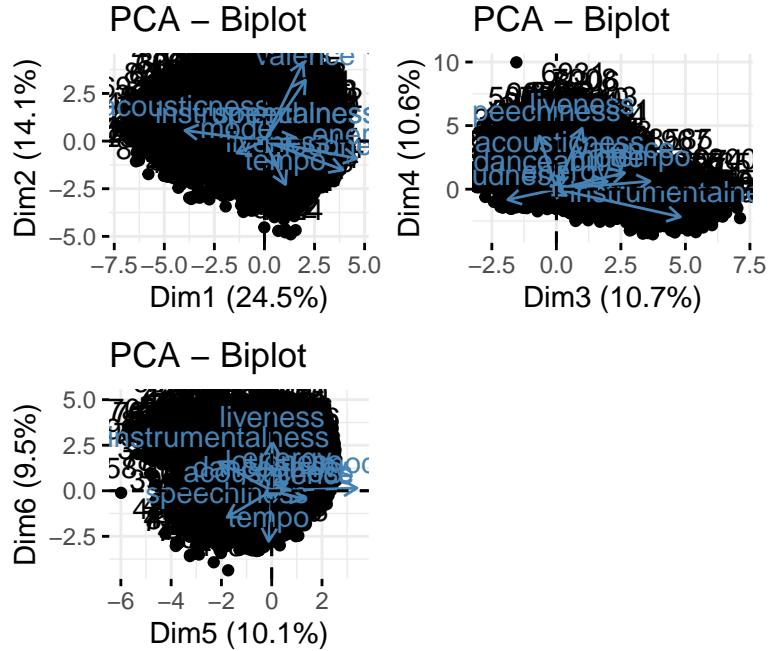
Several methods can be used for unsupervised learning, among them, K-means clustering, Hierarchical Clustering and PCA were performed, due to somewhat disappointing results for Hierarchical Clustering, only K-means clustering was retained.

6.2.1 Principal Component Analysis

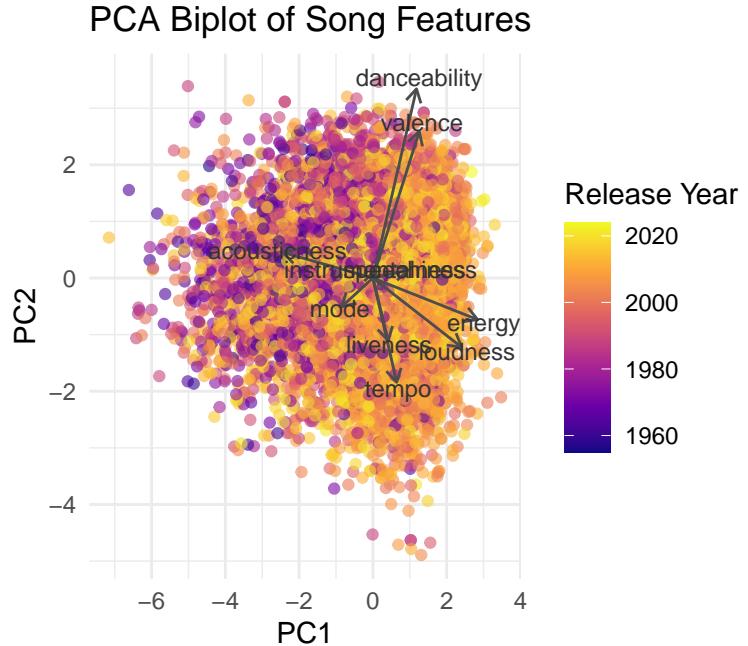
Principal Component Analysis is a dimension reduction technique that helps potentially reduce the number of variables used and that uncovers interactions between variables.



The above graph shows the percentage of explained variance by dimension. As can be observed, the variance explained by 3 dimensions is only 49.3% which is quite underwhelming.



6 dimensions do explain a bit more of the variance, but remain difficult to interpret. Some patterns do emerge. For example, Dimension 1 seems to be related to acousticness vs energy and loudness. On the other hand, Dimension 3 seems to be related with instrumentality. Dimension 5 also seems to relate with mode.



Observing the relation with dimension 1 and 2 colored by release year, Dimension 1 displays the negative relation between acousticness and the positive interaction between energy and loudness. Adding release year information tells us that older tracks tend to be more acoustic, and less loud. This conclusion seems fitting, as older tracks before the appearance of electronic instruments would clearly tend to be more acoustic. Studies have shown that music has been getting louder with time.

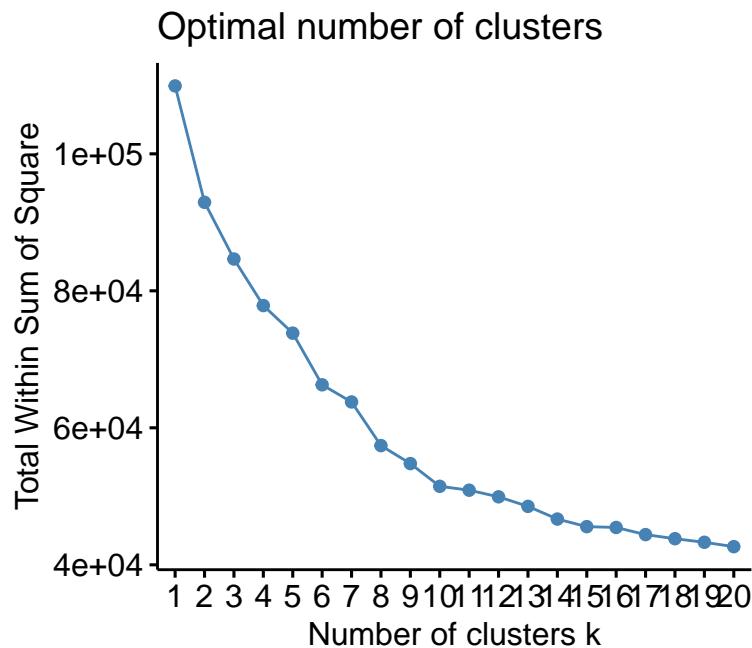
In conclusion, while this PCA analysis struggles to truly reduce the number of variables needed, some interesting interactions can be observed such as the inverse relationship between acousticness and energy/loudness which point to an evolution in music styles through the years.

6.2.2 K-means clustering

A K-means clustering analysis was conducted using the R package *factoextra*. To perform this analysis, only numeric variables related to track features were selected: **danceability, energy, loudness, mode, speechiness, acousticness, instrumentalness, liveness, valence, tempo, and time_signature**. These variables were chosen because they are metrics developed by Spotify, which does not disclose the exact methodologies behind their computation.

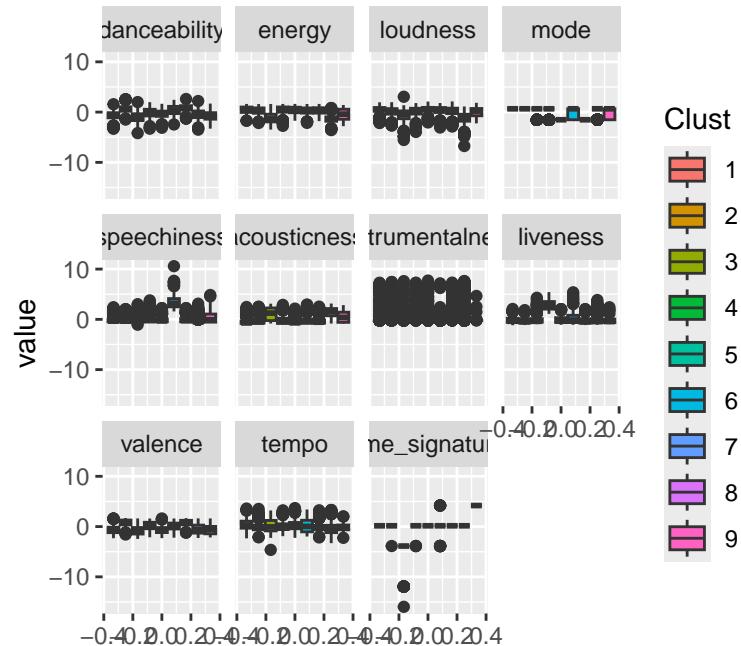
It is speculated that Spotify may use some type of clustering technique to derive these features, potentially to better categorize musical genres and preferences. This makes clustering an appropriate method for exploring patterns within the dataset.

The first step in the clustering process is to determine the optimal number of clusters. This was done using the Total Within-Cluster Sum of Squares (TWSS) as the evaluation metric. By identifying the elbow point in the TWSS curve, the appropriate number of clusters can be selected.



Based on the graph, it was decided to choose 9 clusters.

Boxplots are created to inspect these clusters.



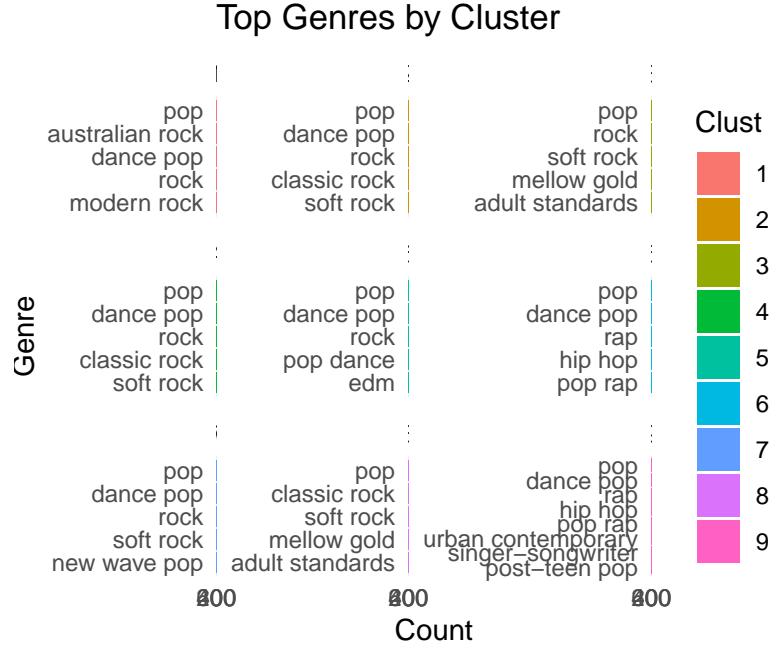
Some patterns seem to emerge from this clustering. Cluster 9 appears to have lower values

for mode compared to others. Cluster 4 has higher speechiness, while clusters 5 and 6 show slightly more acousticness. Cluster 5 also has much higher instrumentalness, and cluster 7 has a significantly lower time signature than the rest.

Cluster	count	Mean Year	Median Year
1	1918	2005.035	2009
2	2555	2000.063	2004
3	403	1999.231	2003
4	544	2000.408	2006
5	1076	2007.619	2011
6	510	2008.267	2010
7	1326	2003.388	2007
8	1610	1995.105	1997
9	51	2007.608	2012

Looking at the release year, some patterns emerge. In terms of count, clusters **1** and **9** contain the most observations. Cluster 9 stands out for having the lowest median release year overall. It is also important to note that differences between the mean and median values, as well as the overall shift toward recent years, are likely due to an imbalance in the dataset, with more recent tracks represented.

The same analysis can be applied to identify genres within each cluster.



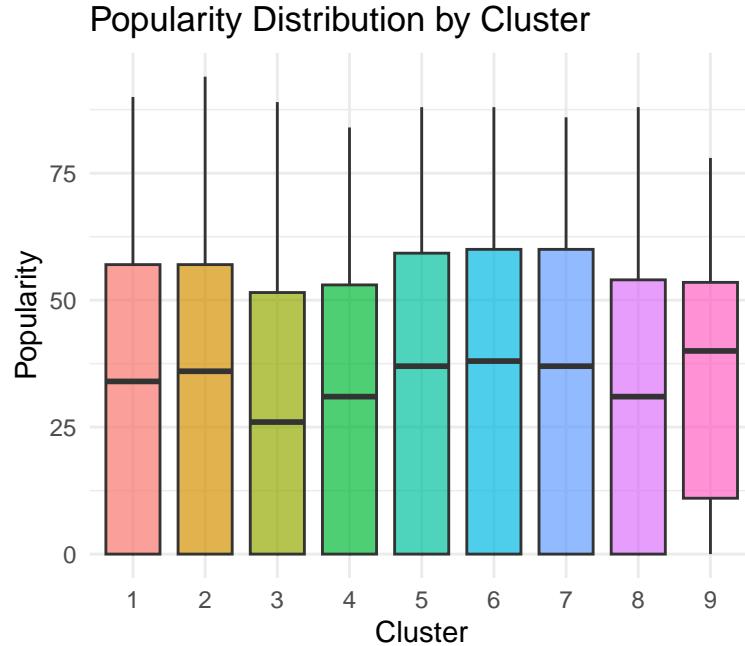
While this graph is interesting, it does not account for the overall distribution of genres. For instance, *Pop* and *Dance Pop* are by far the most present genres in the dataset with respectively 2660 and 1655 instances. To address this, we compute the relative proportion of each genre $rel_i n_c cluster = cluster_n / total_n$, for example, *pop* appears 599 times in Cluster 1 out of 2660 total appearances. Only genres with at least 500 appearances are considered to focus only on main genres instead of niche ones.

Clust	artist_genres	rel_in_cluster	total_n
1	australian rock	0.3284884	688
1	rock	0.2806622	1329
2	new wave pop	0.3944354	611
2	album rock	0.3288288	666
3	No genre	0.0603291	547
3	singer-songwriter	0.0554455	505
4	singer-songwriter	0.0871287	505
4	album rock	0.0795796	666
5	pop	0.1736842	2660
5	dance pop	0.1480363	1655
6	rap	0.3155080	561
6	pop rap	0.2982792	523
7	dance pop	0.2217523	1655
7	new wave pop	0.1882160	611
8	mellow gold	0.3157199	757
8	singer-songwriter	0.3108911	505
9	rap	0.0124777	561
9	singer-songwriter	0.0118812	505

We identify the top two genres for each cluster.

- **Cluster 1:** Dominated by Pop songs.
- **Cluster 2:** Rock is the most frequent genre.
- **Cluster 3:** Features Singer-Songwriter and Mellow Gold, often from artists like Eric Clapton, George Harrison, or Bob Dylan—suggesting influences from folk, country, and soft rock.
- **Cluster 4:** Primarily contains rap music.
- **Cluster 5:** Largely composed of tracks labeled “No genre” and New Wave, likely reflecting more niche or lower-popularity songs.
- **Cluster 6:** Again features Soft Rock and Mellow Gold, indicating some similarity with Cluster 3.
- **Cluster 7:** Includes many songs with no genre and some Singer-Songwriter.
- **Cluster 8:** Appears to be focused on Australian music.
- **Cluster 9:** Dominated by Pop and Dance Pop.

This clustering also reveals trends related to popularity. Some clusters may represent more popular songs than others. This can be visualized easily using boxplots.



Surprisingly, the differences in popularity between clusters are not as pronounced as one might expect. Nevertheless, Cluster 4—dominated by rap—appears to be the most popular. Clusters 5 and 7, which include many songs with “No genre” labels, seem to be the least popular.

To conclude, one key takeaway from this clustering analysis is that the large volume of tracks likely affects the clarity of the results. Still, meaningful insights can be drawn—especially regarding genres. The analysis suggests that genres are, at least partially, distinguishable using Spotify’s audio features.

Release year and popularity offer additional context. The main insight is that our dataset contains significantly more recent songs, and that these tend to be more popular. Currently, genres such as *Pop*, *Dance Pop*, *Rap*, and *Hip-Hop* dominate popular music.

7 Conclusion

This study set out to explore whether machine learning techniques could effectively predict a song’s release period and uncover stylistic patterns in popular music using audio and metadata features. By combining supervised and unsupervised learning methods, we evaluated the extent to which musical characteristics reflect temporal trends across time.

The results demonstrate that while musical features alone do not allow for precise year-by-year prediction, they do contain significant temporal signals that support broader classifications. The Random Forest model, when applied to a binary classification task distinguishing songs released before and after 2007, performed particularly well. This suggests that machine learning is capable of capturing the evolution of popular music at a macro level, even if it struggles with more granular distinctions between adjacent periods.

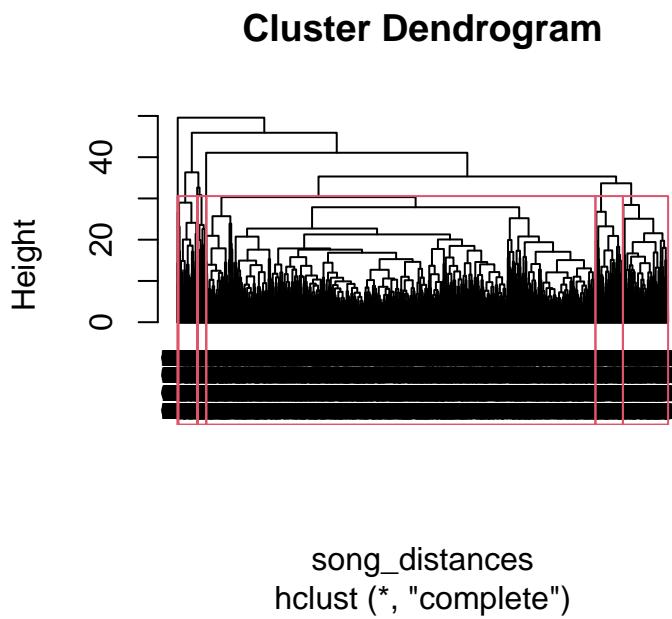
The neural network model combining tabular audio features and album cover images showed no sign of overfitting and demonstrated moderate predictive capacity. However, the average error of approximately 8 years highlights the difficulty of predicting the exact release year of a song, even with enriched inputs. While the model was able to capture certain temporal patterns, a substantial portion of the variance remains unexplained, suggesting that additional factors may play a critical role in determining a song's place in time.

Unsupervised clustering provided complementary insights, revealing that genres and stylistic groupings are partially recoverable using Spotify's audio features. However, the clarity of these results was affected by the large volume of data and imbalance in the temporal distribution of songs.

In summary, this study contributes to the growing field of music analytics by demonstrating how machine learning can uncover meaningful patterns in musical data. While current models are limited in their ability to deliver highly accurate release year predictions, they offer valuable tools for tracing stylistic evolution and classifying music into broader historical categories. Future research could enhance these findings by integrating lyrical analysis, social metadata, and regional diversity to better capture the cultural complexity behind popular music's evolution.

8 Appendix

8.1 Hierarchical clustering



Due to poor performance from the hierarchical clustering, the decision was made to only focus on k-means clustering. In subsequent studies, a deeper focus on this clustering technique could be conducted.

9 References

- Beesa, P., Naregavi, V., Imandar, J., & Thatte, S (2023). *Songs popularity analysis using Spotify data: An exploratory study*. In Vidhyayana – An International Multidisciplinary Peer-Reviewed E-Journal, Volume 8, Special Issue 7, 4th National Student Research Conference on “Innovative Ideas and Invention in Computer Science & IT with its Sustainability” (pp. 211–223). <https://www.researchgate.net/publication/384286217>
- Choi, K., Fazekas, G., Sandler, M., & Cho, K. (2017). *Convolutional recurrent neural networks for music classification*. 2017 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2392–2396. <https://doi.org/10.1109/ICASSP.2017.7952585>
- Friberg, A., Erkut, C., & Bresin, R. (2011). *Perceptual ratings of musical parameters*. Proceedings of the International Conference on Music Perception and Cognition (ICMPC), 193–198.

Gómez-Cañón, J. S., Cano, E., Eerola, T., Herrera, P., Hu, X., Yang, Y.-H., & Gómez, E. (2021). Music emotion recognition: Toward new, robust standards in personalized and context-sensitive applications. *EEE Signal Processing Magazine*, 38(6), 106–114.

Herremans, D., Martens, D., & Sørensen, K. (2019). *Dance hit song prediction*. University of Antwerp Operations Research Group Applied Data Mining Research Group, University of Antwerp.

Interiano, M., Kazemi, K., Wang, L., Yang, J., & Komarova, N. L. (2018). *Musical trends and predictability of success in contemporary songs in and out of the top charts*. Royal Society Open Science, 5(5), 171274. <https://doi.org/10.1098/rsos.171274>

Jan Van Balen, L., Burgoyne, J. A., Bountouridis, D., Müllensiefen, D., & Honingh, A. (2015). *Corpus analysis tools for computational hook discovery*. Proceedings of the 17th ISMIR Conference, 227–233.

Kim, Y. E., Schmidt, E. M., Migneco, R., Morton, B. G., Richardson, P., Scott, J., ... & Turnbull, D. (2010). *Music emotion recognition: A state of the art review*. Proceedings of the 11th International Society for Music Information Retrieval Conference (ISMIR), 255–266.

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), 150081. <https://doi.org/10.1098/rsos.150081>

Müller, M., Ewert, S., & Grosche, P. (2010). *Musical structure analysis using audio features*. In Müller, M. (Ed.), *Fundamentals of Music Processing* (pp. 135–154). Springer. https://doi.org/10.1007/978-3-319-21945-5_6

Pachet, F. (2008). *Hit song science is not yet a science*. Proceedings of the 9th International Conference on Music Information Retrieval (ISMIR), 355–360.

Schedl, M., Gómez, E., & Urbano, J. (2015). *Music information retrieval: Recent developments and applications*. Foundations and Trends in Information Retrieval, 8(2-3), 127–261. <https://doi.org/10.1561/1500000042>

Serrà, J., Corral, Á., Boguñá, M., Haro, M., & Arcos, J. L. (2012). *Measuring the evolution of contemporary Western popular music*. Scientific Reports, 2(1), 521. <https://doi.org/10.1038/srep00521>

Spotify. (2023). *How Spotify's popularity score works*. Spotify Developer Documentation. Retrieved from <https://developer.spotify.com/>