# Success in music Production: Analysis of Hot 100 Songs Over Time

By: Sayedmahdi Raghib

Sayedmahdi.raghib@rwth-aachen.de

Advisor : Prof. Jürgen Lerner
Date: 04.07.2023

**RWTH**AACHEN
UNIVERSITY

# Contents

Objective

Approach

Results

Reflection

Limitations

Future Work

# Objective

Research Questions

## Objective

The goal of this analysis is to provide valuable insights to music industry players on what makes a song successful:

- Analyze the **successful** songs that made it to the Weekly Hot 100 singles chart by Billboard with other **unsuccessful** songs that are not in the chart to examine their **characteristics** such as genre, tempo and danceability.

- **Sentiment analysis** on both successful and unsuccessful songs

- **Top keyword analysis** on lyrics of the songs

# Billboard Hot 100

- Billboard Hot 100 is the music industry standard record chart in the United States for songs, published weekly by Billboard magazine.

- Chart rankings are based on sales (physical and digital), radio play, and online streaming in the United States.
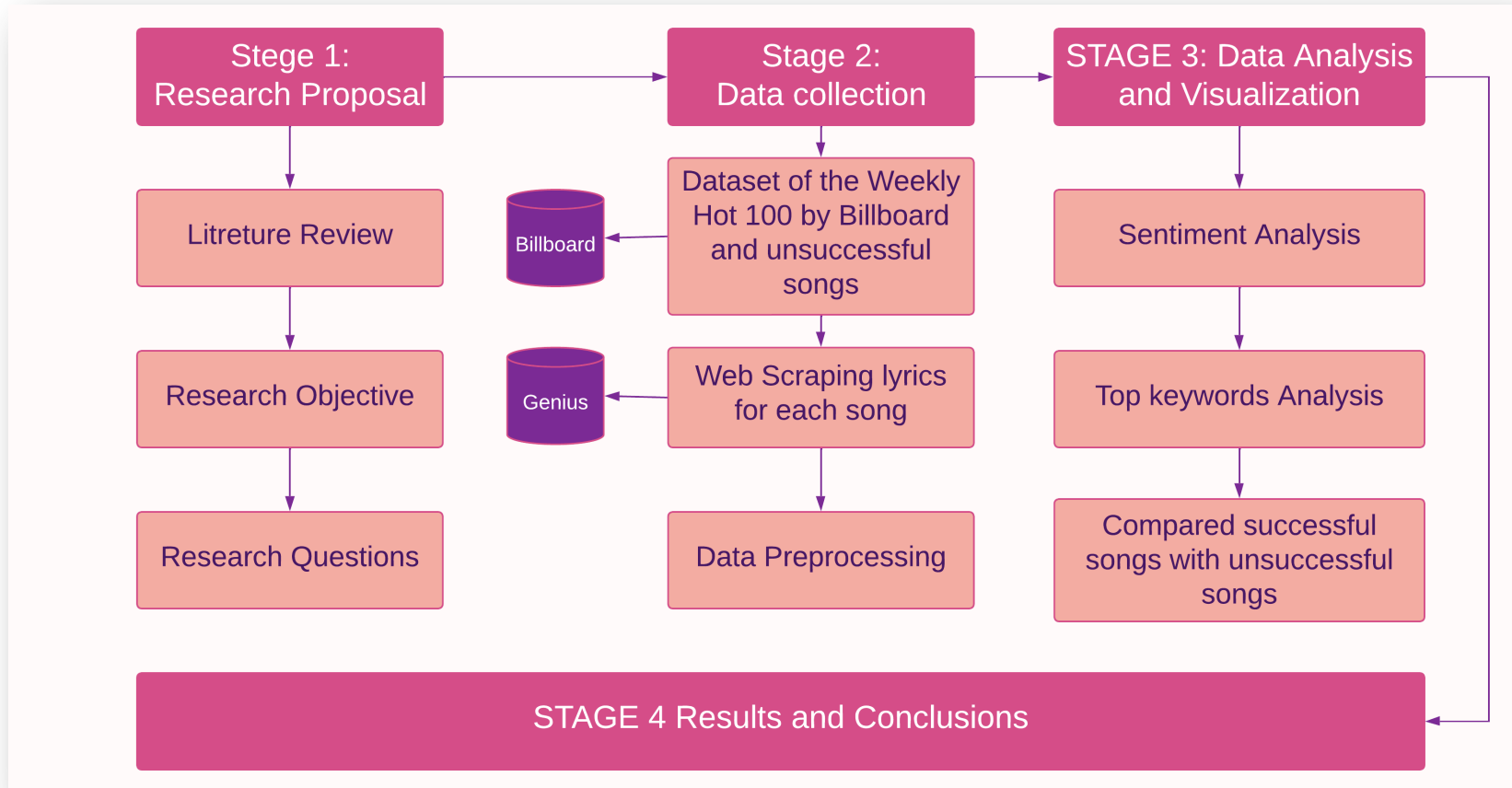
# Research Questions

1. What is a pattern of the longevity of hot songs? (Some hit songs appear in the list once and some more)

2. What patterns of repeated success among the artists? (Unique hit songs released by artists.)

3. What is the temporal trend of the audio features?

   - Do tracks become more danceable or more relaxed over the past years based on the audio attribute levels?

   - Is there a seasonal mood change?

## Research Questions

5. How has the sentiment of lyrics in popular music changed over time?

6. What are the top keywords that have been consistently present in popular song lyrics over time?

7. Is there any correlation between the sentiment of the song's lyrics and its chart performance? Do More positive or negative songs tend to be more successful?

8. How has the vocabulary used in popular music changed over time? Are certain words more common in certain eras?

# Approach

How I have tackled the task

# Research Framework

## Data Description

**Successful songs dataset** provided us with information every weekly Hot 100 singles chart from Billboard.com from 1958 to 2019. (data.world):

- Total number of successful songs is 21741.

**Unsuccessful songs dataset** provided us with any songs that have never appeared in the weekly Hot 100 singles chart from Billboard.com. (Kaggle.com)

- Total number of unsuccessful songs is 23671.

## Data Description

**Audio Features**

- **acousticness** - A confidence measure from 0.0 to 1.0 of whether the track is acoustic. 1.0 represents high confidence the track is acoustic.

- **danceability** - Danceability describes how suitable a track is for dancing based on a combination of musical elements including tempo, rhythm stability, beat strength, and overall regularity. A value of 0.0 is least danceable and 1.0 is most danceable.

- **energy -** Energy is a measure from 0.0 to 1.0 and represents a perceptual measure of intensity and activity. Typically, energetic tracks feel fast, loud, and noisy. For example, death metal has high energy, while a Bach prelude scores low on the scale.

RWTH AACHEN UNIVERSITY

## Data Description

**Audio Features**

- **instrumentalness** - Predicts whether a track contains no vocals. "Ooh" and "aah" sounds are treated as instrumental in this context. Rap or spoken word tracks are clearly "vocal". The closer the instrumentalness value is to 1.0, the greater likelihood the track contains no vocal content.

- **loudness** - The overall loudness of a track in decibels (dB). Values typical range between -60 and 0 db.

- **valence** - A measure from 0.0 to 1.0 describing the musical positiveness conveyed by a track. Tracks with high valence sound more positive (e.g. happy, cheerful, euphoric).

# Approach

**Lyrics Scraping**

- I scraped the lyrics for the songs from the website genius.com using a python API client called LyricsGenius.

# Approach

## Sentiment Analysis

- Performed sentiment analysis on the lyrics used [TextBlob](#) library.

# Approach

## Keyword Extraction

- Extraction the keywords from lyrics of each song by combining the Rapid Automatic Keyword Extraction (RAKE) algorithm with the NLTK toolkit, known as [rake-nltk](#).
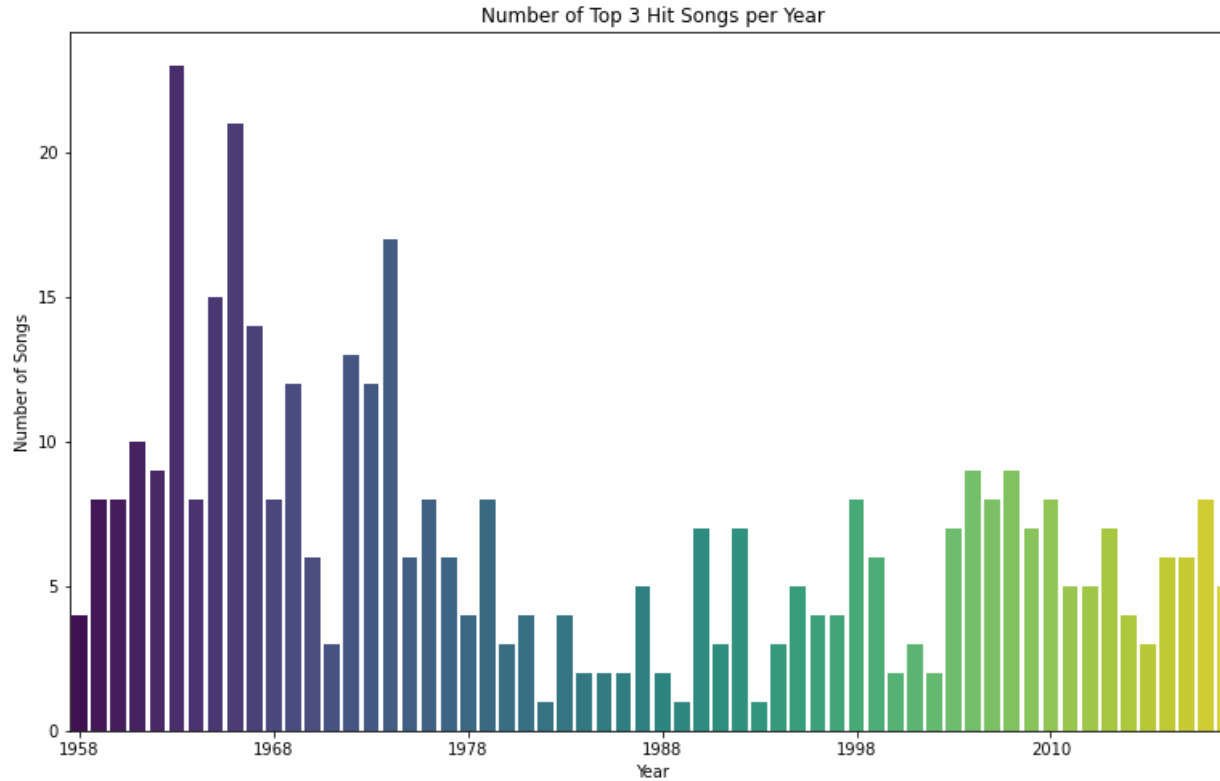
# Results

# Results

- What is a pattern of the longevity of hot songs?

# Results

- What is the number of top hits on Billboard over time ?



Number of Top 3 Hit Songs per Year
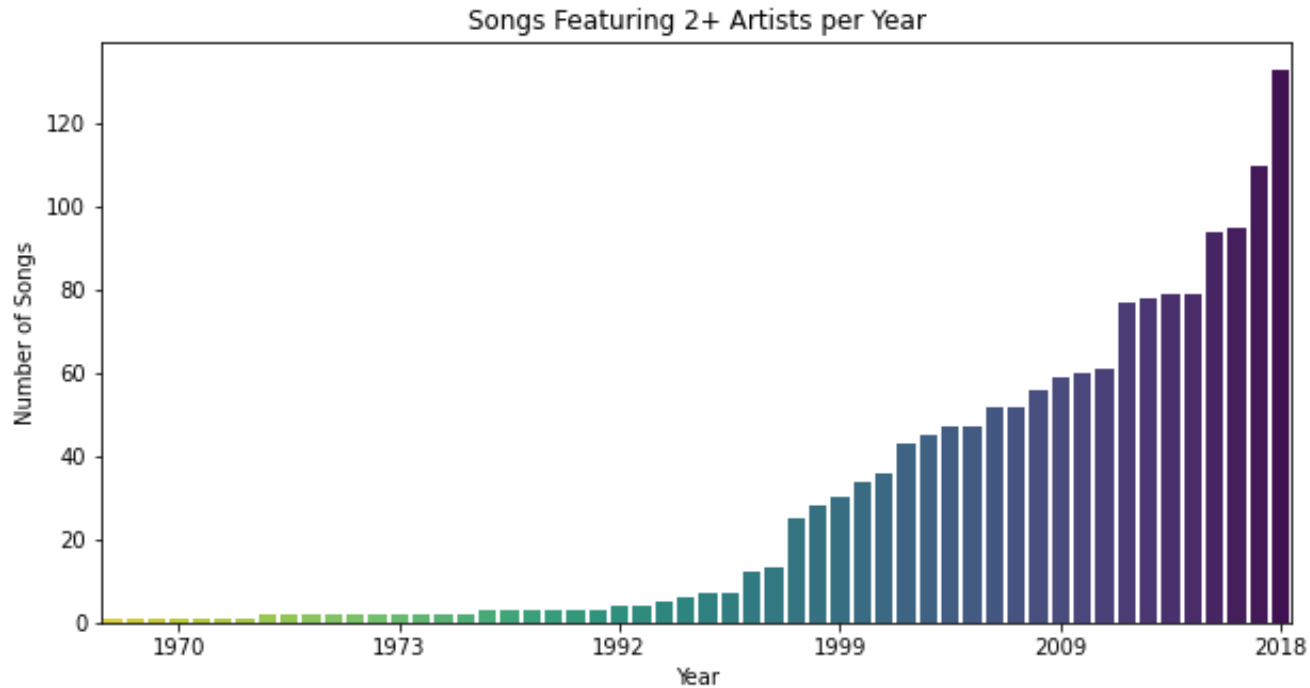
# Results

- What patterns of repeated success among the artists?



Number of Songs per Artist



Number of Songs per Artist per Year

# Results

- What patterns of repeated success among the artists?



Songs Featuring 2+ Artists per Year

# Results

- What patterns of repeated success among the artists?

Mean Weeks on Hot 100 per Year

# Results

- What is the temporal trend of the audio features?

**Successful Songs:**

**Unsuccessful Songs:**

# Results

- What is the temporal trend of the audio features?

**Successful Songs:**



**Unsuccessful Songs:**

# Results

- What is the temporal trend of the audio features?

**Successful Songs:**



Distribution of Duration over the Years

# Results

- Do tracks become more danceable ?

**Successful Songs:**



**Unsuccessful Songs:**

# Results

- Is there a seasonal mood change in the successful songs?

**Over the years:**

**Over the months:**

# Results

- Mood of successful songs:

- happy: valence > 0.5, energy > 0.5

- excited: valence <= 0.5, energy > 0.5

- sad: valence <= 0.5, energy <= 0.5

- peaceful: valence > 0.5, energy <= 0.5.

# Results

- What is the correlation between the audio features?

**Successful Songs:**



Correlation between the Features

# Results

- How has the sentiment of lyrics in popular music changed over time?

**Successful Songs:**



Mean Sentiment Score for Each Year

# Results

- How has the sentiment of lyrics in popular music changed over time?

**Unsuccessful Songs:**



Mean Sentiment Score for Each Year

# Results

- What are the top keywords in successful song lyrics over time?

**Successful Songs:**

# Results

- What are the top keywords in successful song lyrics over time?

**Successful Songs:**

# Results

- What are the top keywords in successful song lyrics over time?

**Unsuccessful Songs:**



Most Common Words in Lyrics

# Results

- What are the top keywords in successful song lyrics over time?

**Unsuccessful Songs:**

# Results

- Song lyrics' variance can be described using how many components?

**Successful Songs:**



Explained Variance Ratio by Number of Components

# Reflection

# Reflection

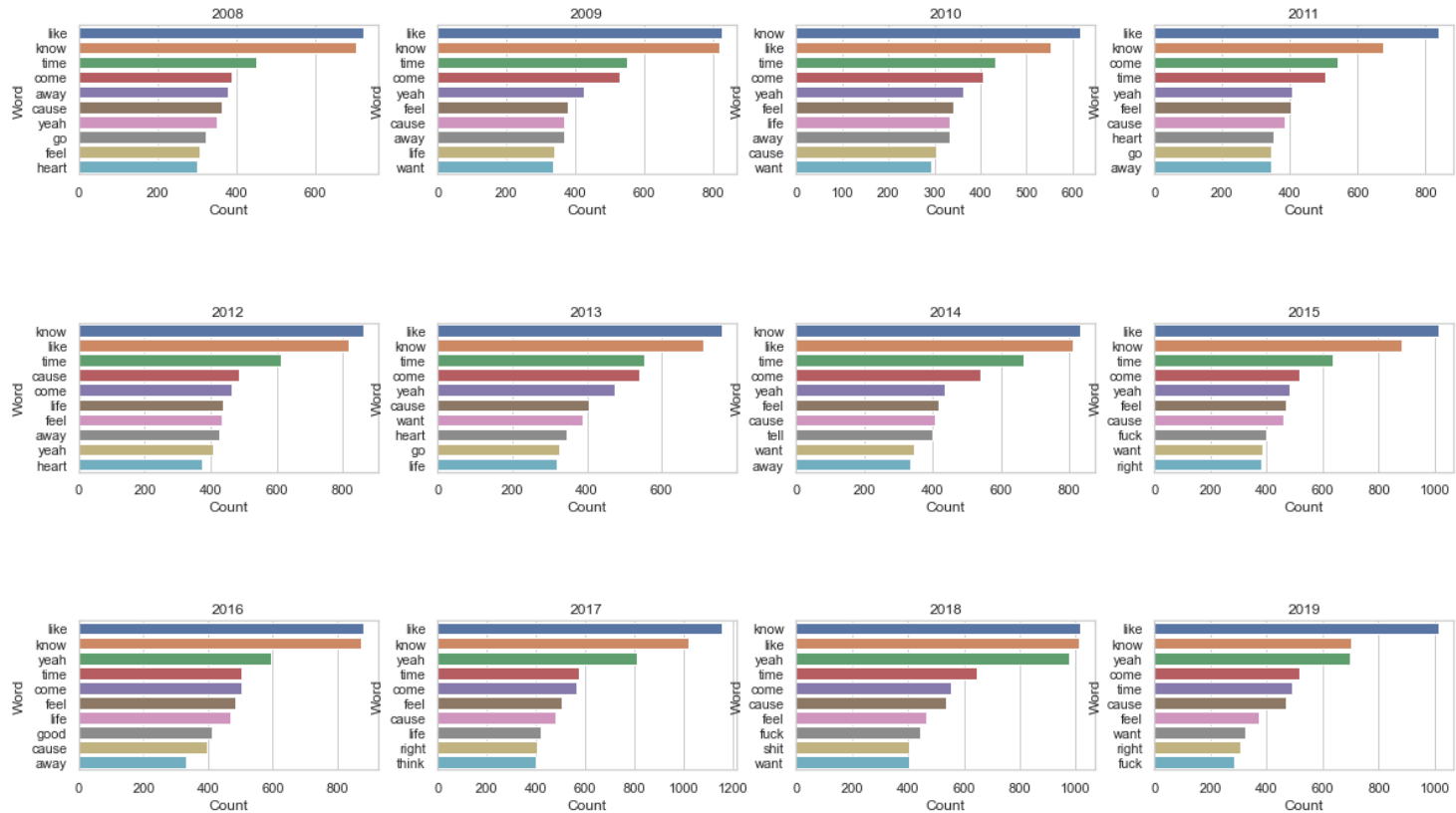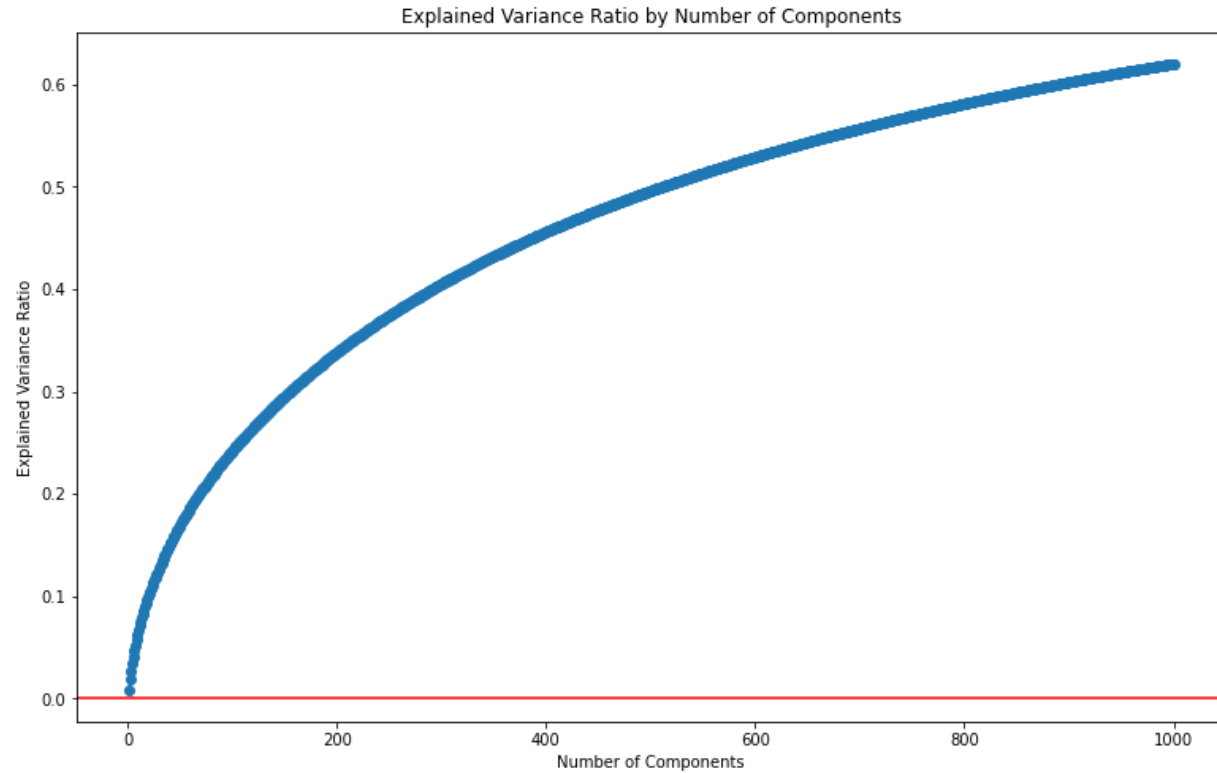- I performed a logistic regression analysis on the set of all songs (successful and unsuccessful) to assess the effect of the different variables (e.g., sentiment) on the probability that a song becomes a hit song.

```python
# create a list of independent variables
X = df[['Sentiment', 'danceability', 'energy', 'loudness', 'acousticness','instrumentalness', 'valence']]

# create a list of dependent variables
y = df['successful']

# split the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=0)

# fit the model
logreg = LogisticRegression()
logreg.fit(X_train, y_train)

# make predictions on the testing set
y_pred = logreg.predict(X_test)

# evaluate the model
cnf_matrix = metrics.confusion_matrix(y_test, y_pred)
```
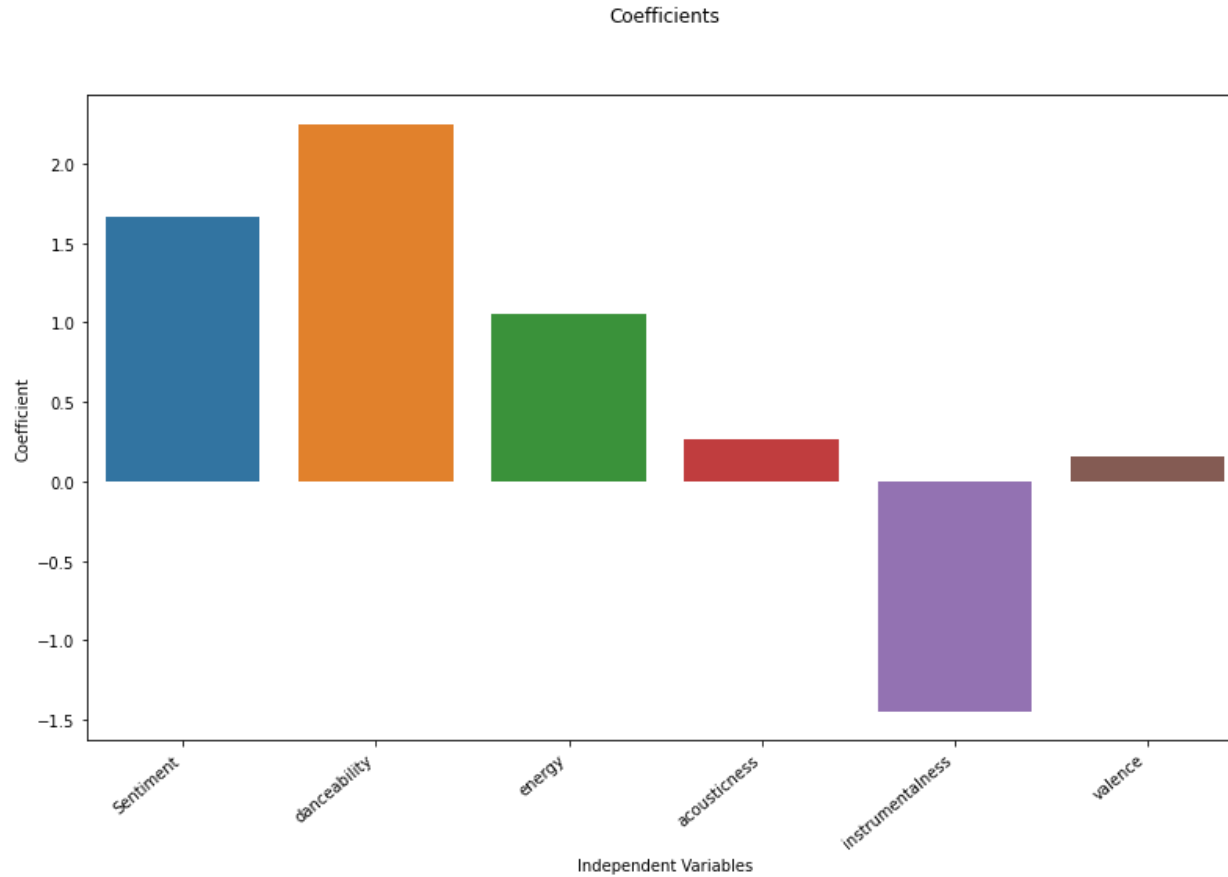
```
Accuracy: 0.6175563463819692
Precision: 0.6120985599260963
Recall: 0.6175563463819692
F1-score: 0.6111971509878914
```

# Reflection

# Limitations

# Limitations

**Weaknesses**

- **Time period limitations:** Sample does not include 4 years recent songs.

- **Generalizability:** The findings of your research may be specific to the Billboard Hot 100 chart and the U.S. market. It might not necessarily apply to other charts or international markets, which may have different criteria and preferences for success.

- **Theoretical limitations:** research questions cannot be answered by empirical data alone without theoretical frameworks for guidance.

# Future Work

# Future Work

**Future Steps:**

- **Expand the dataset:** Consider incorporating a larger and more diverse dataset that includes songs beyond the Billboard Hot 100 chart.

- **Cross-cultural analysis:** Explore the variations in success factors and trends across different cultural contexts. This could involve comparing the Billboard Hot 100 chart with charts from other countries or analyzing songs from different cultural backgrounds.

- **Machine learning-based prediction models:** Develop predictive models using machine learning techniques to forecast the potential success of a song based on its features, sentiment, and keywords.

# Thank you for your attention!

Sayedmahdi.raghib@rwth-aachen.de