# MUSIC POPULARITY PREDICTION

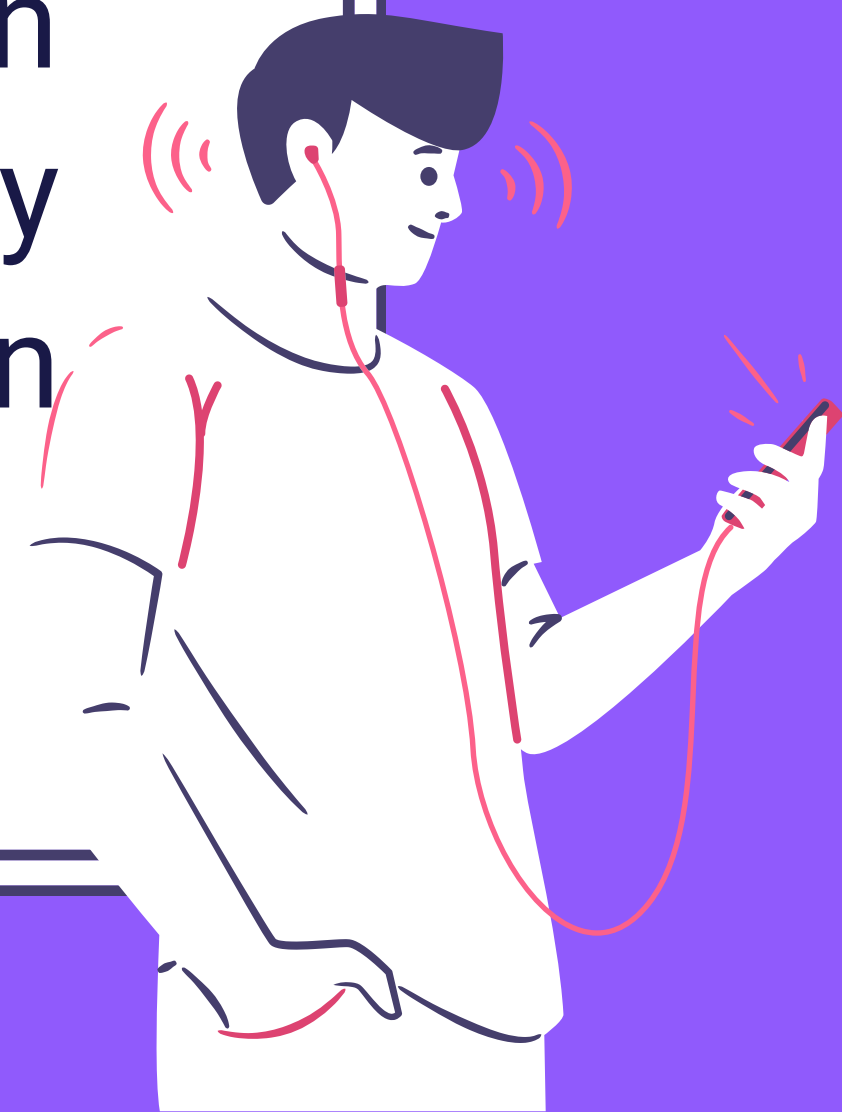Presented By:

Group 4

# Table Of Content

# OVERVIEW

**Why certain songs are popular?**

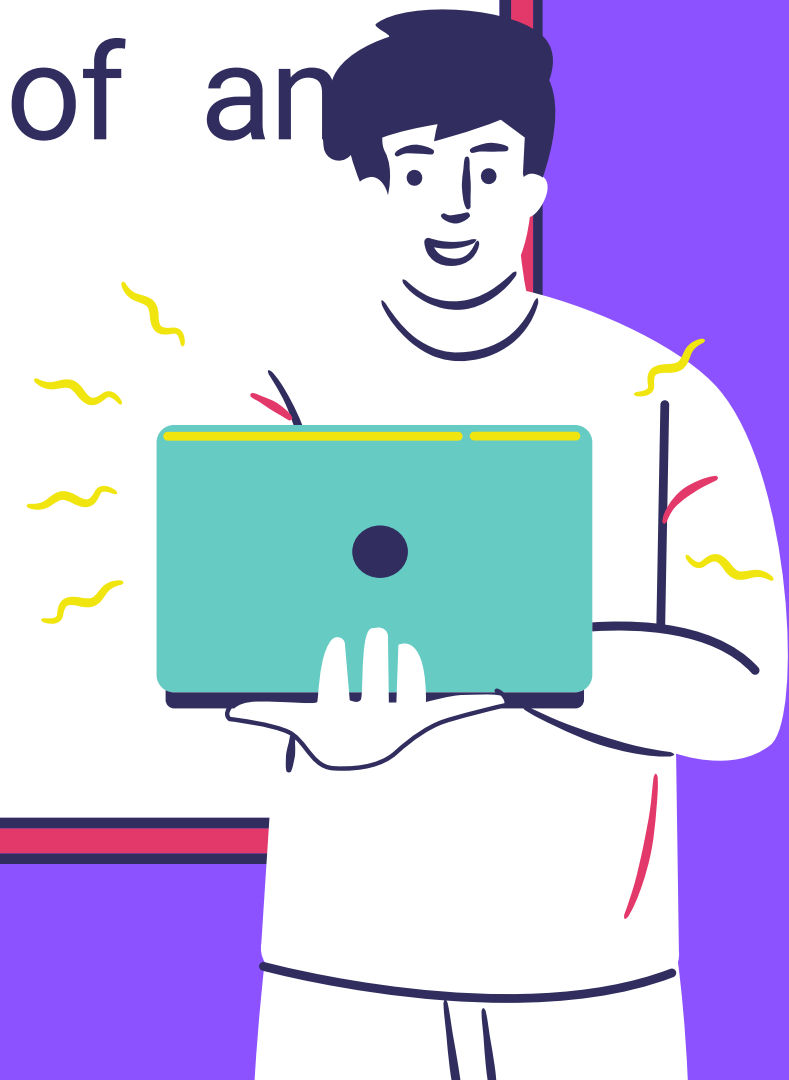**What makes them popular amongst millions of other songs?**

**Music has become integral of many people's life**

- We should know that there are some specific attributes of a song that makes them popular.

- With our analysis we are trying ascertain the factors responsible for the popularity of the songs by uncovering the trends in the data.

- Our goal here is to determine if a song makes its place in the top 40% even before the release.

- Music companies can be highly benefitted by this model in predicting the popularity of an upcoming song.

# BUSINESS UNDERSTANDING

## Porter's Five Forces framework

- **Bargaining Power of Suppliers - (High)**

A music service is nothing without actual music to play, as there are four major music labels that hold monopoly on most music. Artists and their catalog of music are highly distinct in customer's mind, so the supplier has the higher power of bargaining.

- **Bargaining Power of Customers - (Low)**

Buyers are negotiating the digital music market as individuals rather than large groups. As a result, individual buyers/consumers have very little power over the music providers.

- **Threat of New Entrants - (Low)**

It's low for well established music industries and it's notoriously hard for newcomers to break into the industry

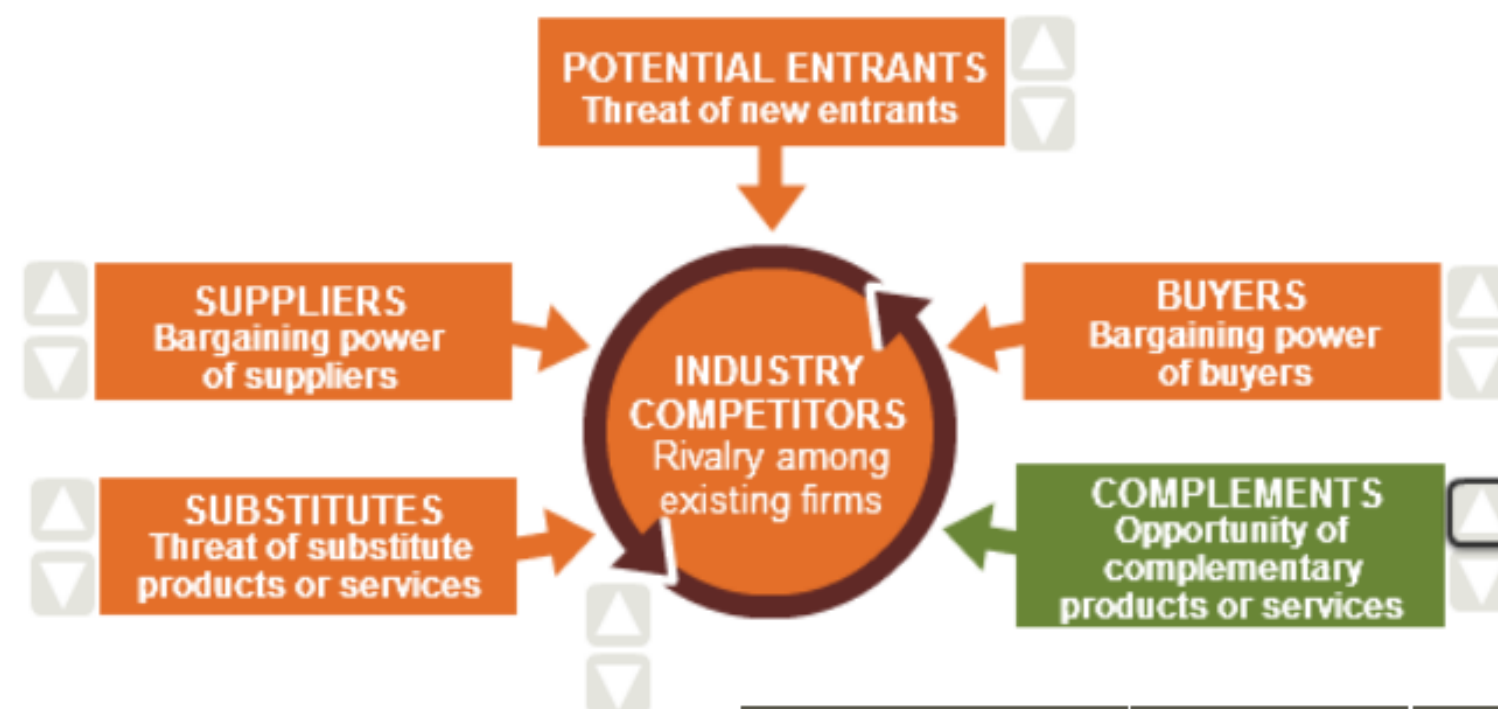- **Availability of Substitute Products - (High)**

A person can consume music many ways—with a number of them free. So much choice makes it difficult for digital music providers to drive profit margins higher.

- **Intensity of Competitive Rivalry in the Industry - (High)**
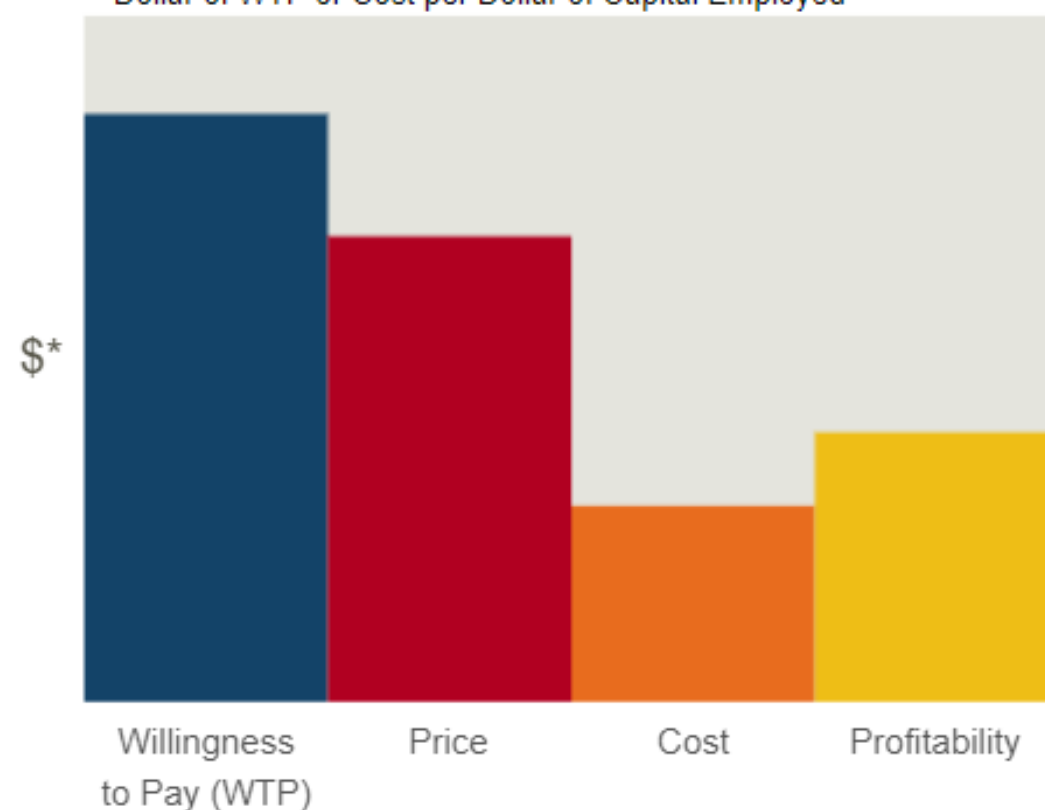
The user experience of consuming music can vary with each provider. Since there is no clear leader(s) in the digital music market, companies are fighting hard to secure customer loyalty.

## Porter's Forces Framework

Legend:
- Positive
- Neutral
- Negative

POTENTIAL ENTRANTS
Threat of new entrants

SUPPLIERS
Bargaining power of suppliers

BUYERS
Bargaining power of buyers

INDUSTRY COMPETITORS
Rivalry among existing firms

SUBSTITUTES
Threat of substitute products or services

COMPLEMENTS
Opportunity of complementary products or services

* Dollar of WTP or Cost per Dollar of Capital Employed

Bar chart ($*):
- Willingness to Pay (WTP)
- Price
- Cost
- Profitability

| The Force | Impact | | | Root Causes | | | | |
|---|---|---|---|---|---|---|---|---|
| IF threat of entry ⬆ | | | | | | | | |
| IF threat of entry ⬇ | profitability | ⬆ | because | WTP ⬆ | Prices ⬆ | Costs ⬇ | | |
| IF supplier power ⬆ | profitability | ⬇ | because | Costs ⬆ | | | | |
| IF supplier power ⬇ | | | | | | | | |
| IF buyer power ⬆ | | | | | | | | |
| IF buyer power ⬇ | profitability | ⬆ | because | Prices ⬆ | | | | |
| IF substitutes ⬆ | profitability | ⬇ | because | WTP ⬇ | Prices ⬇ | | | |
| IF substitutes ⬇ | | | | | | | | |
| IF rivalry ⬆ | | | | | | | | |
| IF rivalry ⬇ | | | | | | | | |
| IF complements ⬆ | profitability | ⬆ | because | WTP ⬆ | Prices ⬆ | | | |
| IF complements ⬇ | | | | | | | | |

Sources

HARVARD BUSINESS PUBLISHING

# ♪ SONY ENTERTAINENT

## Target Firm - Sony Entertainment
- American entertainment company established in 2012.
- Headquarters is located at New York, 18000 employees.

## Major Products/ Service lines :
- Sony Pictures Entertainment(7.097 billion USD)
- Sony Music Group(8.86 billion USD)

## Geographic Footprints:
- United States of America
- Latin America and the Caribbean
- Europe
- Singapore
- Japan

# SWOT Analysis

## Strength

- High Market Share
- Strong Brand Awareness
- Social Media Marketing

## Weakness

- Bad and Low Quality music production
- Lousy investments on newly signed artists

## Opportunities

- AI and ML has notched some notable wins in the field of music in line of popularity prediction.

## Threat

- Wrong selection of an album released may result badly to the company
- Piracy

# DATA ANALYSIS AND UNDERSTANDING

- Dataset - Spotify Data. https://developer.spotify.com/
- Dataset Source - Spotify.
- It has 34080 records and 23 columns.
- There are 24128 unique tracks with 10379 artists.
- 23 attributes for each song, 12 of them numerical.
- Time Period - 1956 - 2022.
- Target Variable - POPULARITY.

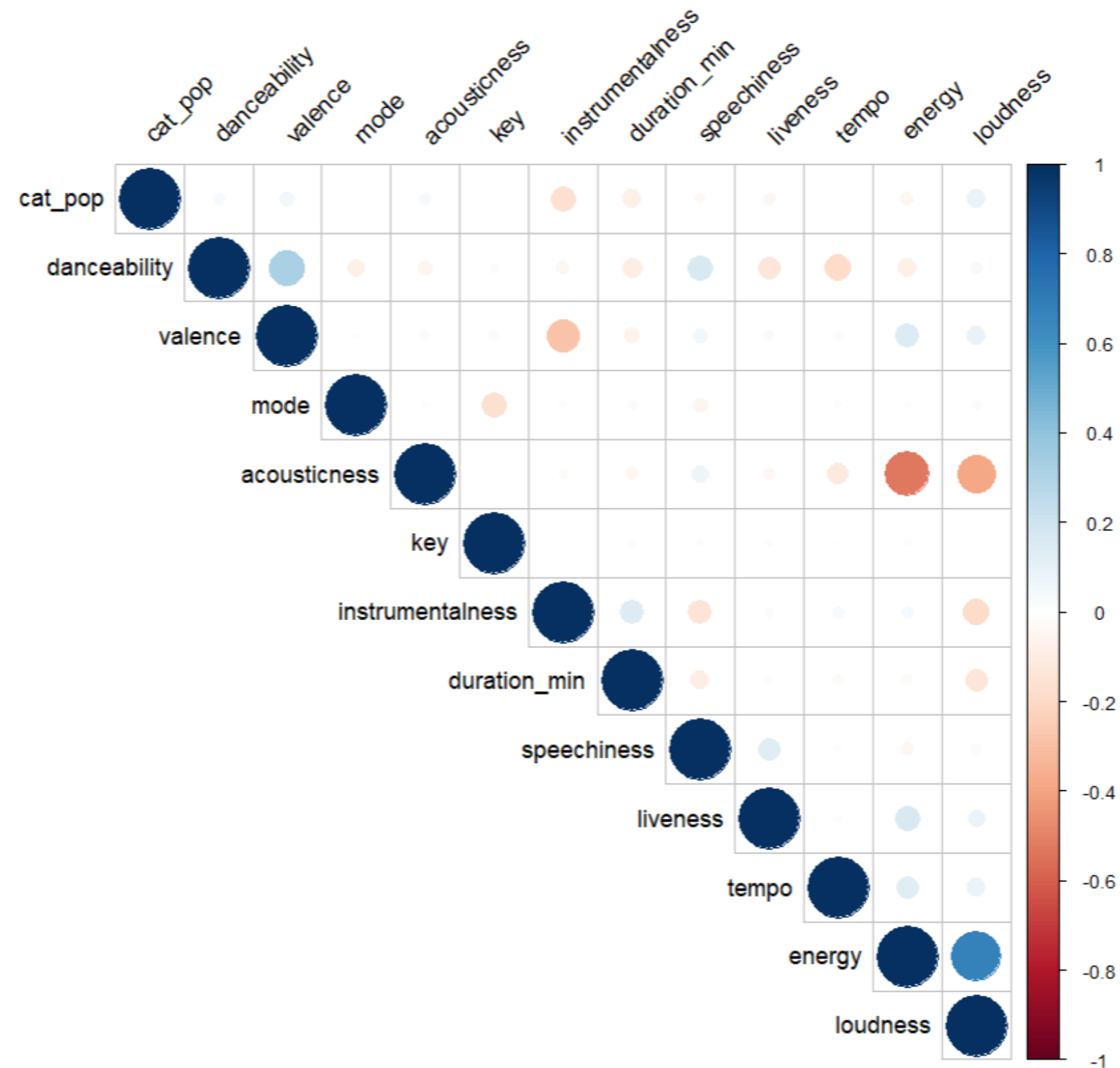| ATTRIBUTE | DESCRIPTION |
|---|---|
| Acousticness | A confidence measure from 0.0 to 1.0 of whether the track is acoustic. |
| Danceability | Danceability ranges from 0.0 - 1.0 and describes how suitable a track is for dancing based on a combination of musical elements including tempo, rhythm stability, beat strength, and overall regularity. |
| duration_ms | Duration of the track in ms |
| 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. |
| instrumentalness | Predicts whether a track contains no vocals. "Ooh" and "aah" sounds are treated as instrumental in this context. Values ranges from 0 - 1 |
| key | The estimated overall key of the track. Values ranges from 0 - 11 |
| liveness | Detects the presence of an audience in the recording. Higher liveness values represent an increased probability that the track was performed live. |

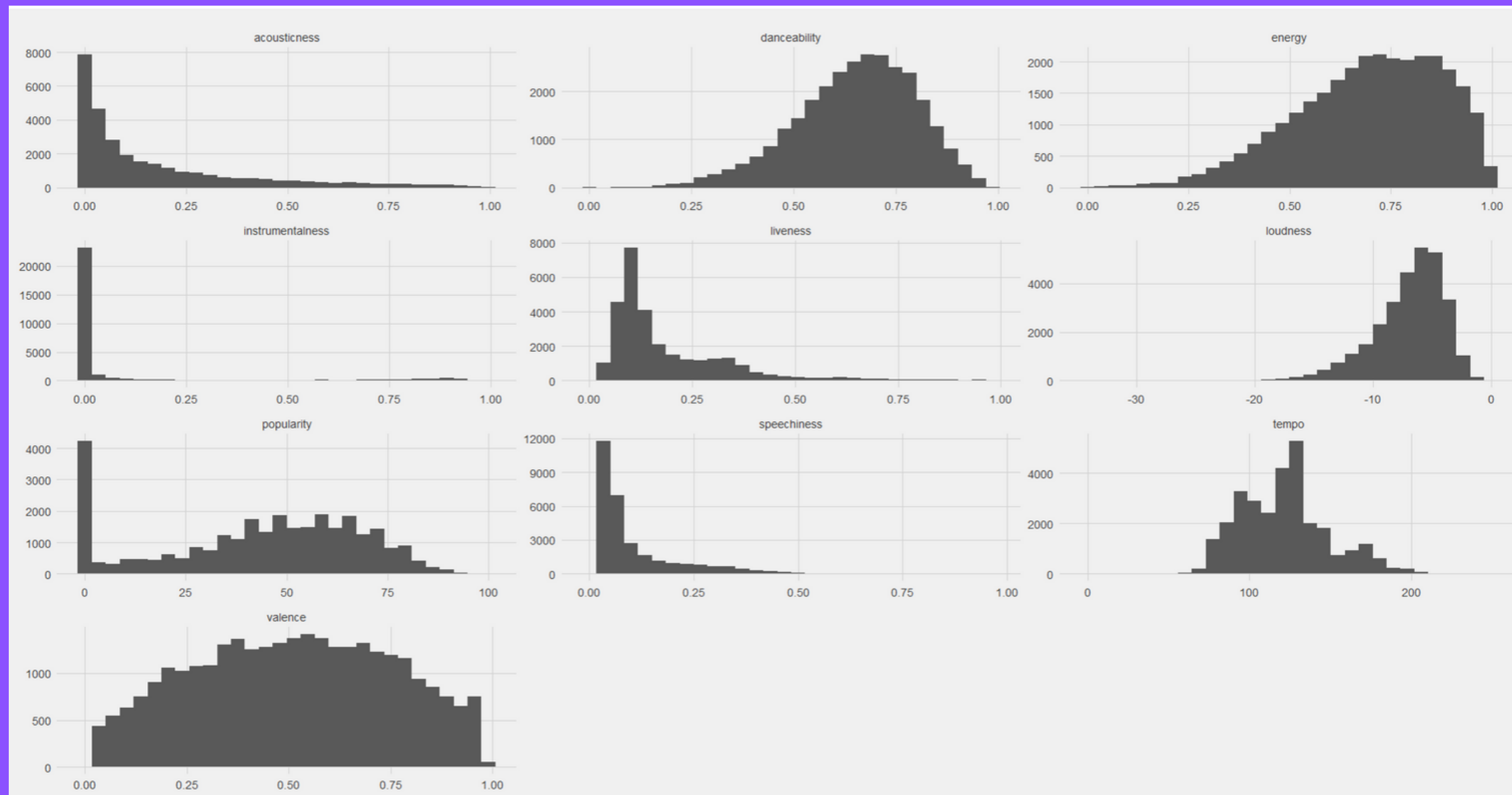| ATTRIBUTE | DESCRIPTION |
| --- | --- |
| loudness | The overall loudness of a track in decibels (dB). Values ranges from -60 - 0db. |
| mode | Mode indicates the modality (major or minor) of a track, the type of scale from which its melodic content is derived. |
| speechiness | Speechiness detects the presence of spoken words in a track. The more exclusively speech-like the recording (e.g. talk show, audio book, poetry), the closer to 1.0 the attribute value. |
| tempo | The overall estimated tempo of a track in beats per minute (BPM). In musical terminology, tempo is the speed or pace of a given piece and derives directly from the average beat duration. |
| 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), while tracks with low valence sound more negative (e.g. sad, depressed, angry). |

# Exploratory Data Analysis

Descriptive Statistics

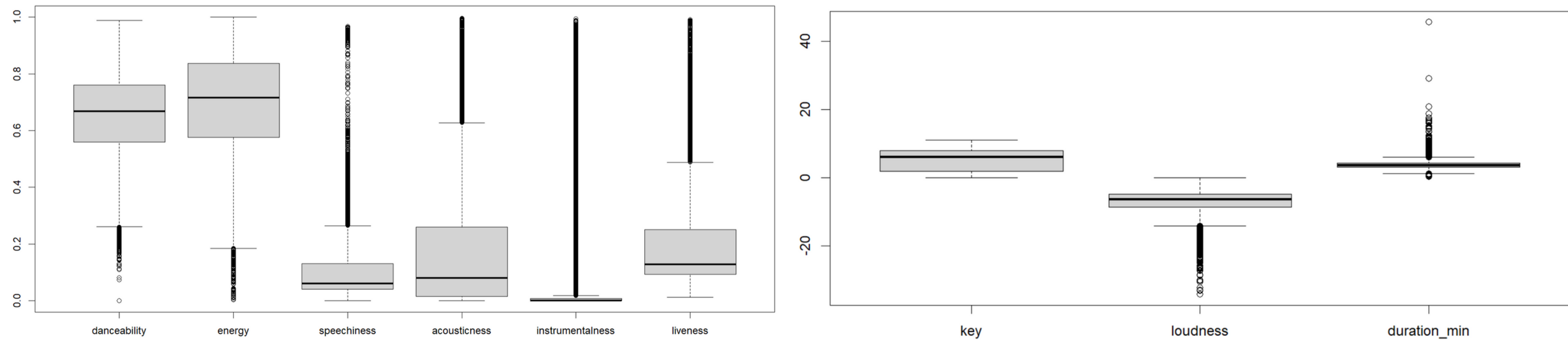|  | mean | median | sd | variance | min | max | count | miss.val |
|---|---|---|---|---|---|---|---|---|
| popularity | 42.97342865 | 47.0000000 | 25.2138004 | 635.73572962 | 0.00000000 | 100.00000 | 29656 | 0 |
| danceability | 0.65261971 | 0.6670000 | 0.1459697 | 0.02130715 | 0.00000000 | 0.98800 | 29656 | 0 |
| energy | 0.69440247 | 0.7150000 | 0.1796107 | 0.03226001 | 0.00251000 | 1.00000 | 29656 | 0 |
| key | 5.33136633 | 6.0000000 | 3.5926388 | 12.90705321 | 0.00000000 | 11.00000 | 29656 | 0 |
| loudness | -6.94016057 | -6.3630000 | 3.0443399 | 9.26800542 | -34.11400000 | 0.00000 | 29656 | 0 |
| mode | 0.57519558 | 1.0000000 | 0.4943216 | 0.24435386 | 0.00000000 | 1.00000 | 29656 | 0 |
| speechiness | 0.10918757 | 0.0598000 | 0.1136165 | 0.01290870 | 0.00000000 | 0.96700 | 29656 | 0 |
| acousticness | 0.17659792 | 0.0788000 | 0.2209015 | 0.04879749 | 0.00000132 | 0.99600 | 29656 | 0 |
| instrumentalness | 0.09907276 | 0.0000192 | 0.2446509 | 0.05985407 | 0.00000000 | 0.99500 | 29656 | 0 |
| liveness | 0.19061995 | 0.1280000 | 0.1535439 | 0.02357573 | 0.01140000 | 0.99200 | 29656 | 0 |
| valence | 0.51086714 | 0.5155000 | 0.2423133 | 0.05871572 | 0.00000000 | 0.99000 | 29656 | 0 |
| tempo | 121.39733669 | 122.0210000 | 27.1897991 | 739.28517578 | 0.00000000 | 249.43800 | 29656 | 0 |
| duration_min | 3.79708943 | 3.6069083 | 1.1912893 | 1.41917018 | 0.40488333 | 45.70093 | 29656 | 0 |
| Bin_pop | 0.28678851 | 0.0000000 | 0.4522696 | 0.20454776 | 0.00000000 | 1.00000 | 29656 | 0 |

# Correlation Matrics



- **We can see that there is a strong linear relationship between energy and loudness.**
- **Inverse linear relationship between acousticness and energy**
- **Inverse linear relationship between acousticness and loudness**

# Feature Distribution



- Most of the songs are not acoustic.
- Most of the songs are not recorded during live concert as liveness distribution is close to 0.
- Songs are pretty fast as the mean of the temp is ~ 120.

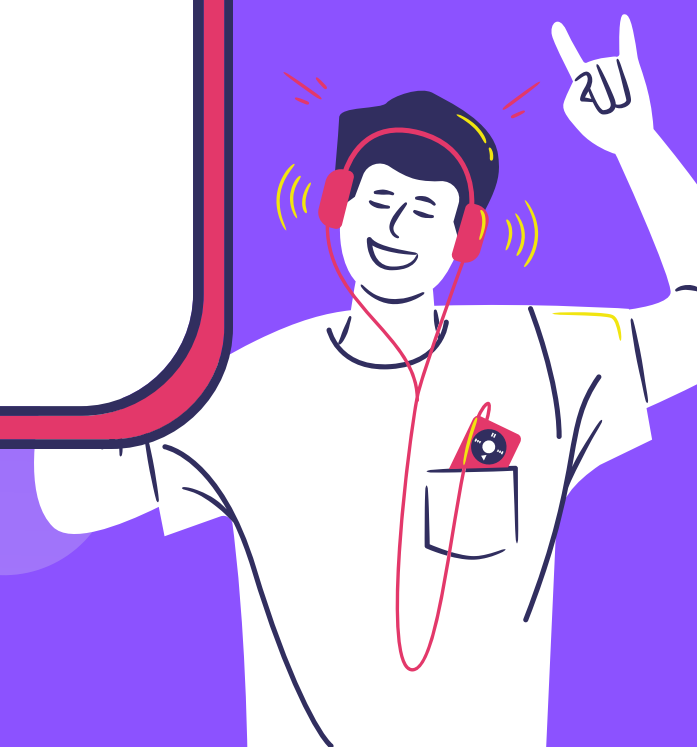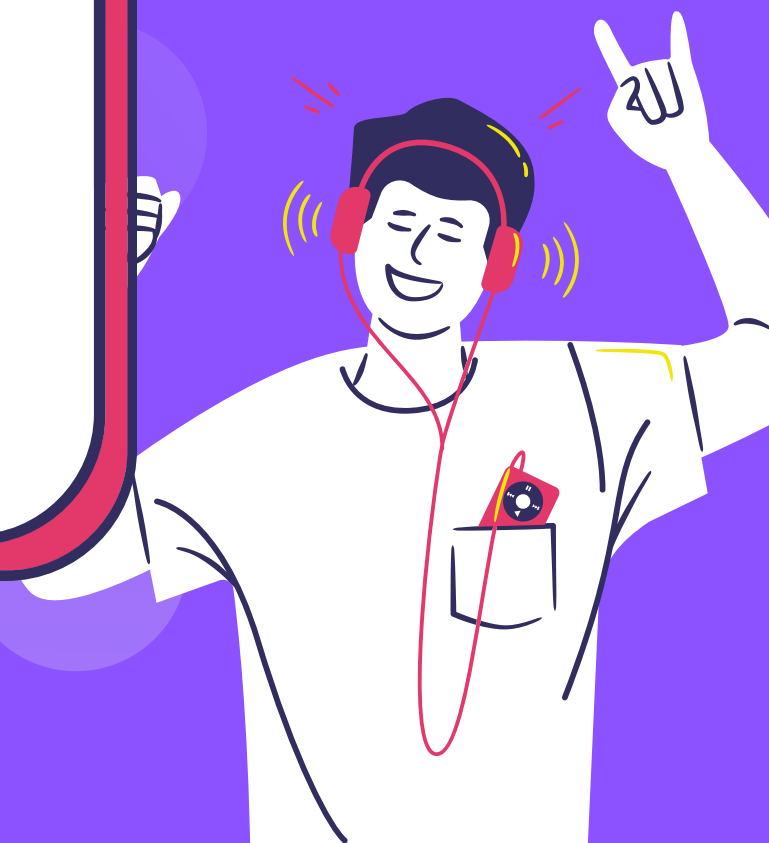# Box Plots of the Independent Variables

# Data Modeling

- The goal of our project is to predict the popularity of the song.
- The target variable in our analysis is popularity and it's value ranges from 0 - 100.
- For the purpose of our analysis we have categorised the popularity into binary values 0 and 1 with the threshold value of 60.
- With this analysis one can determine if a song lies in the top 40% of the songs.
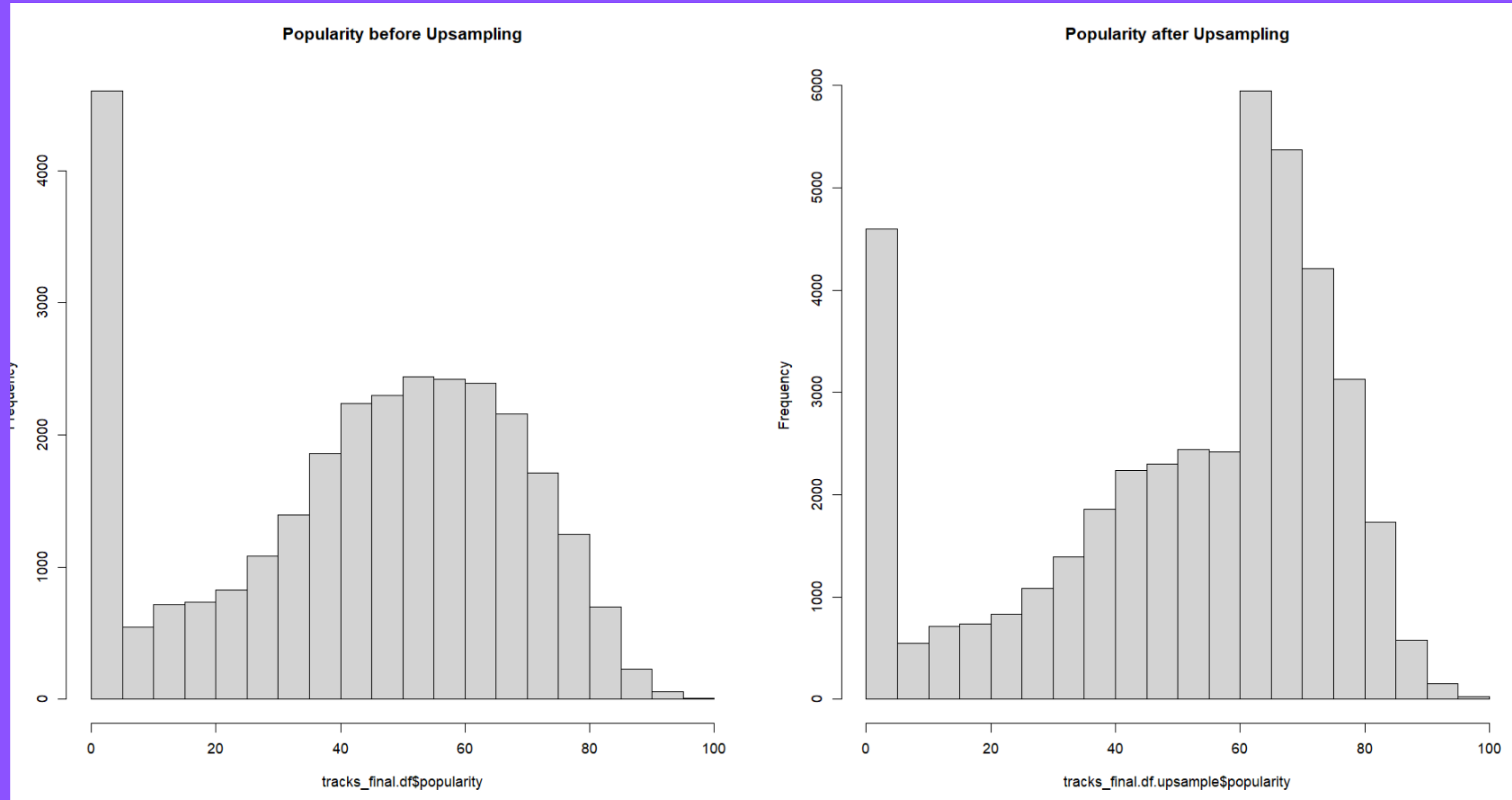
# Data Transformations

- Removed unnecessary columns such as track.album.id, track.album.name, playlist_name, playlist_id.
- Renamed few columns for easy understanding.
- Eliminated the duplicate records using id.
- Removed NA's present in artist and track.name columns.
- Loudness values > 0 is set to 0.
- Converted duration_ms into minutes.
- Categorised the popularity into binary values 0 and 1 with the threshold value of 60.
- Factorized artist and genre.
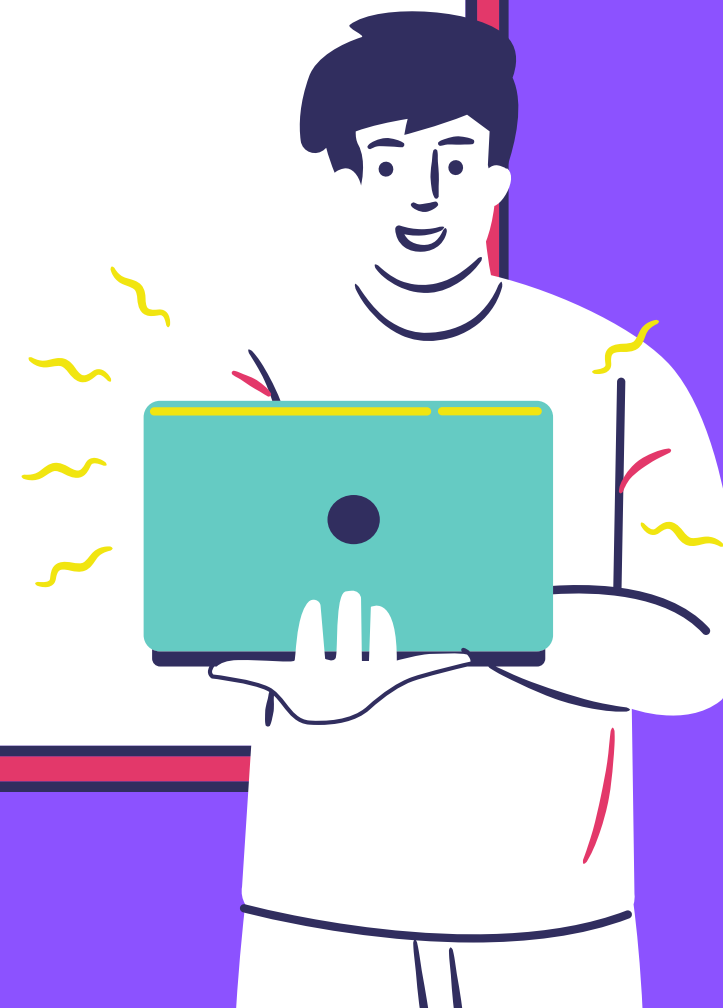- Upsampled our data to create a balance between 0's and 1's.

# Distribution of Popularity(target variable)

# K Nearest Neighbours

- The model has the highest accuracy of 79.97% when k = 1.
- Type 1 error - 13.4%
- Type 2 error - 28.14%
- Hence we can conclude that we are 79.97% sure that the target song will be in top 40% of the songs.

```
Confusion Matrix and Statistics

                 Reference
Prediction    0     1
         0 7432 2389
         1 1000 6100

              Accuracy : 0.7997
                95% CI : (0.7936, 0.8057)
   No Information Rate : 0.5017
   P-Value [Acc > NIR] : < 2.2e-16

                 Kappa : 0.5996

Mcnemar's Test P-Value : < 2.2e-16

           Sensitivity : 0.8814
           Specificity : 0.7186
        Pos Pred Value : 0.7567
        Neg Pred Value : 0.8592
            Prevalence : 0.4983
        Detection Rate : 0.4392
  Detection Prevalence : 0.5804
     Balanced Accuracy : 0.8000

      'Positive' Class : 0
```
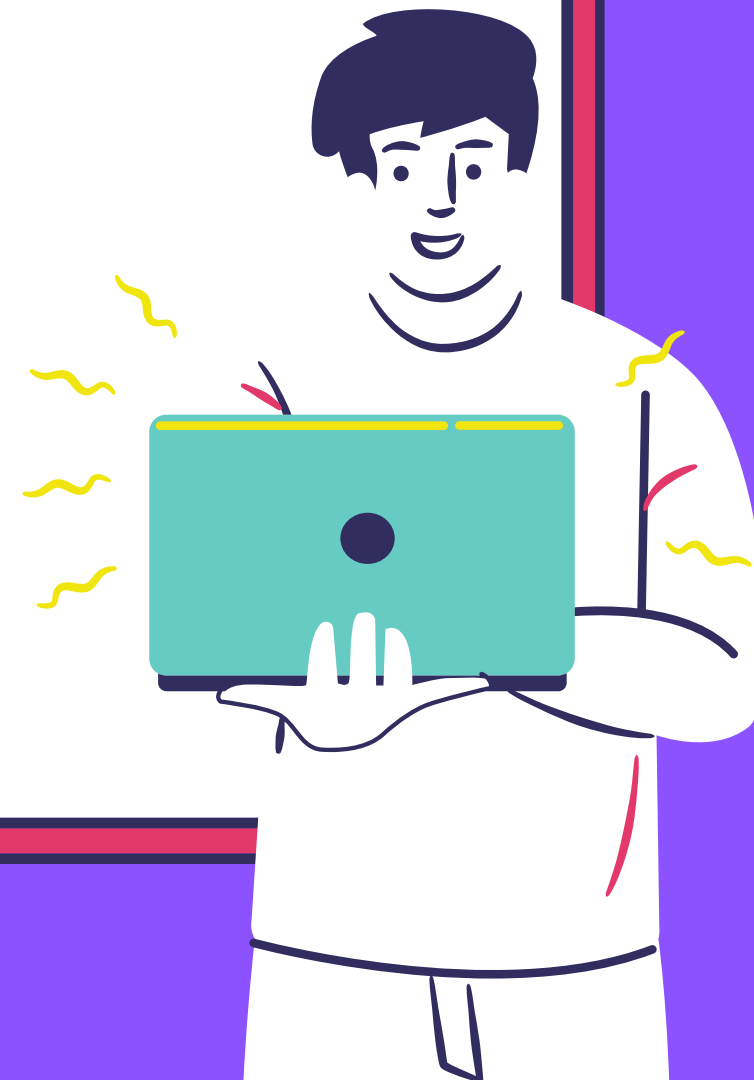
# Decision Tress

- This model has an accuracy of 77.14%.
- To reduce the complexity of the decision trees we have used pruning technique and reduced the size of the tree.
- CP - 0.00167
- Best Split - 10
- Type 1 error - 9.7%
- Type 2 error - 35.8%

```
Confusion Matrix and Statistics

                Reference
Prediction     0     1
        0 7609 3045
        1  823 5444

              Accuracy : 0.7714
                95% CI : (0.765, 0.7777)
    No Information Rate : 0.5017
    P-Value [Acc > NIR] : < 0.0000000000000022

                 Kappa : 0.5432

 Mcnemar's Test P-Value : < 0.0000000000000022

           Sensitivity : 0.9024
           Specificity : 0.6413
        Pos Pred Value : 0.7142
        Neg Pred Value : 0.8687
            Prevalence : 0.4983
        Detection Rate : 0.4497
  Detection Prevalence : 0.6296
     Balanced Accuracy : 0.7718

      'Positive' Class : 0
```
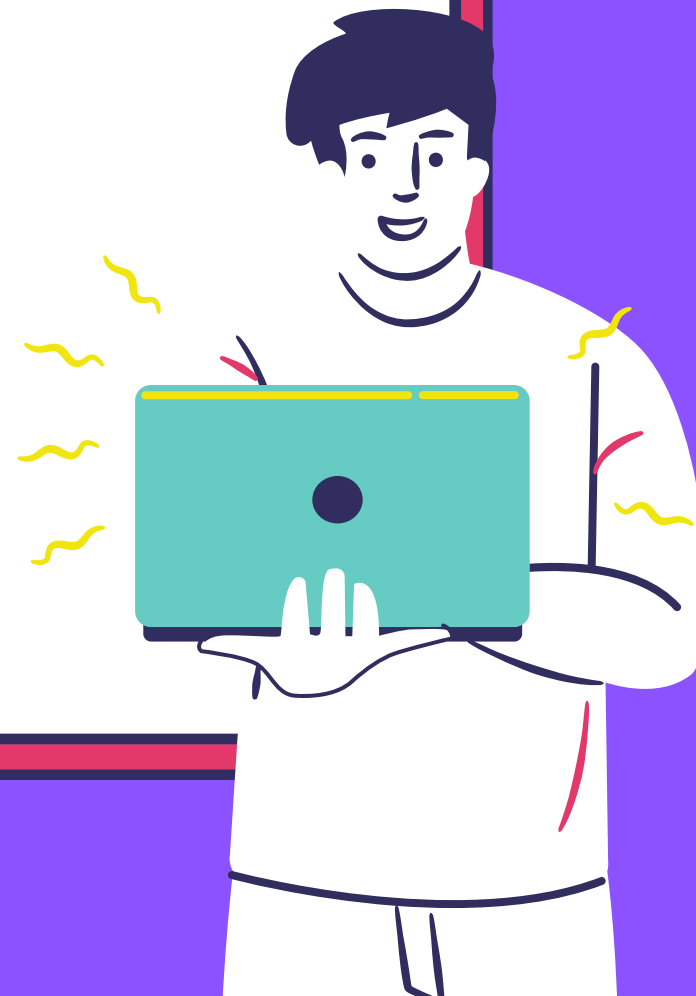
# Random Forest

- This model has the accuracy of 59.69%.
- ntree - 1000
- Type 1 error - 2.6%
- Type 2 error - 77.67%

```
Confusion Matrix and Statistics

               Reference
Prediction    0     1
         0  8205 6594
         1   227 1895

              Accuracy : 0.5969
                95% CI : (0.5895, 0.6043)
   No Information Rate : 0.5017
   P-Value [Acc > NIR] : < 2.2e-16

                 Kappa : 0.1958

Mcnemar's Test P-Value : < 2.2e-16

           Sensitivity : 0.9731
           Specificity : 0.2232
        Pos Pred Value : 0.5544
        Neg Pred Value : 0.8930
            Prevalence : 0.4983
        Detection Rate : 0.4849
  Detection Prevalence : 0.8746
     Balanced Accuracy : 0.5982

      'Positive' Class : 0
```
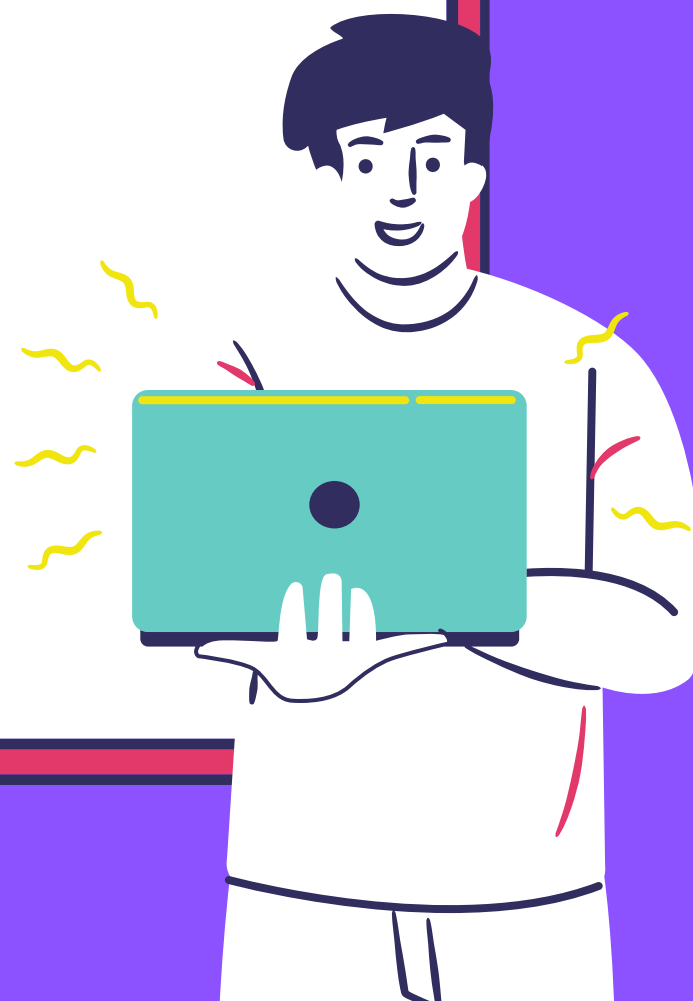
# Logistic Regression

- This model has the accuracy of 54.1%.
- Threshold value - 0.285.
- Independent Variables used are - danceability, speechiness, valence, duration_min.
- Type 1 error - 46.67%
- Type 2 error - 45.14%

```
Confusion Matrix and Statistics

          FALSE TRUE
FALSE     4497 3832
TRUE      3935 4657

               Accuracy : 0.541
                 95% CI : (0.5334, 0.5485)
    No Information Rate : 0.5017
    P-Value [Acc > NIR] : <0.0000000000000002

                  Kappa : 0.0819

 Mcnemar's Test P-Value : 0.2471

            Sensitivity : 0.5333
            Specificity : 0.5486
         Pos Pred Value : 0.5399
         Neg Pred Value : 0.5420
             Prevalence : 0.4983
         Detection Rate : 0.2658
   Detection Prevalence : 0.4922
      Balanced Accuracy : 0.5410

       'Positive' Class : FALSE
```

# Summary

| Model | Score |
|---|---|
| K Nearest Neighbours | 79.97% |
| Decision Tree | 77.14% |
| Random Forest | 59.69% |
| Logistic Regression | 54.1% |

# Model Evaluation

- KNN is the best models with accuracy 79.97% and Type 1 error of 13.4%.

**Model limitations**
- Currently we have a sample dataset of around 34k records from a metadata of 5M rows. The KNN algorithm doesn't work well with the large dataset.
- KNN is sensitive to outliers and missing values as it works on euclidian distance.
- We can separate the tracks in different types such as Instrumental music, live performed, podcasts which will give us the clustered data and remove the outliers.
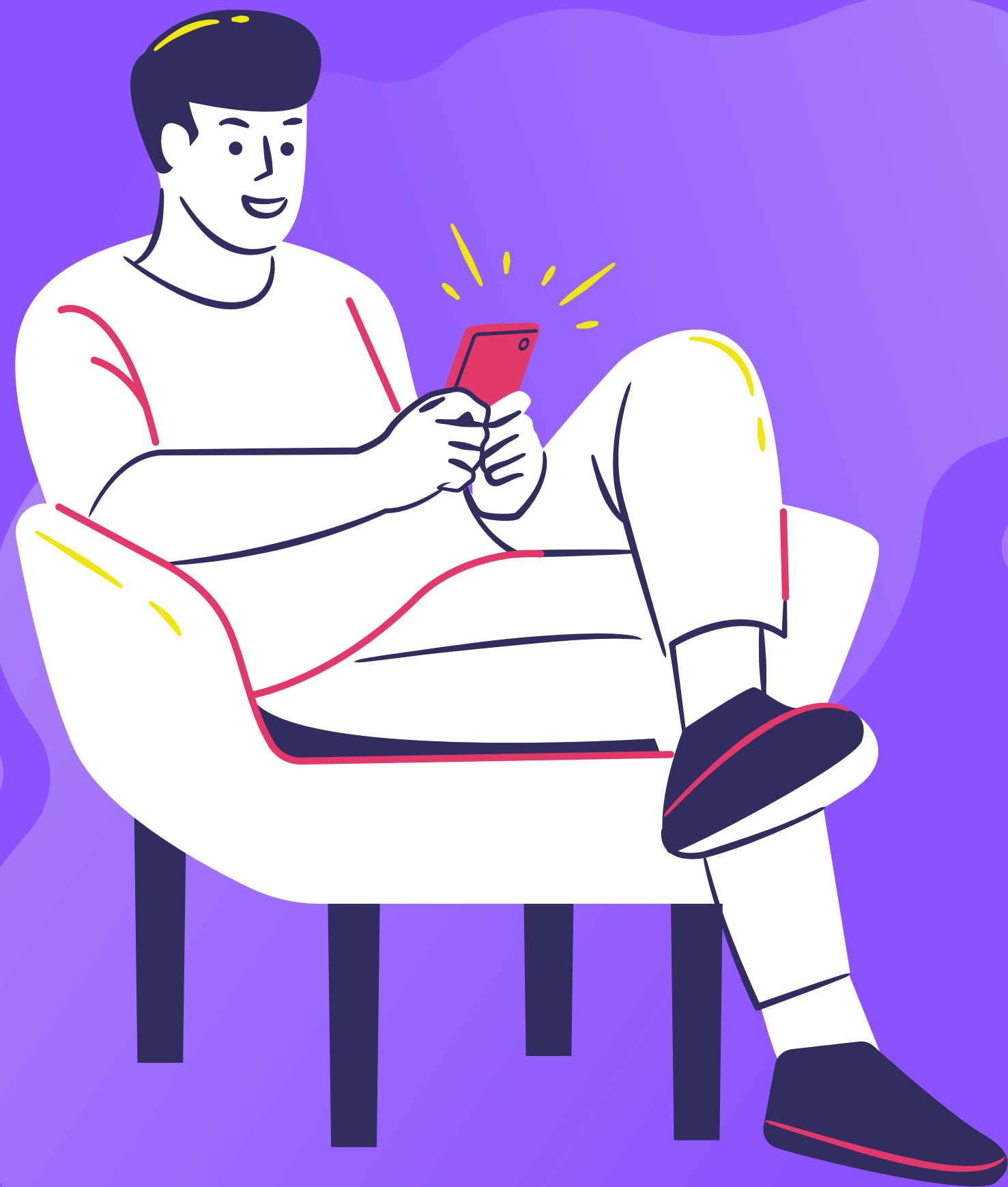
# Recommendations

Can be used to predict if a song can hit the top 40% when a new artist comes up with a new song to our record label.

Using this model people who work in music industry can craft their musical works to better suit with the market needs.

Music industry can also have the capability to monitor and shape peoples music listening behaviour using this model.

# Thank You

Music is the color for the world and everything in it