# Tools and Technology Lab
# Project

# MUSIC RECOMMENDATION SYSTEM

# REPORT

By:
Praneesh Sharma, 21052264
Nayeer Naushad, 21052338
Ronita Chakraborty, 21052697
Antra Amrit, 21052780
Aditya Sharma, 21052472

# ABSTRACT

This project report presents the development and evaluation of a music recommendation system utilizing KMeansClustering to identify similar songs based on audio features. The primary objective of the project was to create a robust recommendation system capable of suggesting relevant music tracks to users based on their preferences. The methodology involved collecting a diverse dataset of music audio features, preprocessing the data, and applying KMeansClustering to group similar songs into clusters. The system was then evaluated using various metrics to assess its performance in recommending songs accurately.

Key findings from the project include the successful implementation of the KMeansClustering algorithm to effectively group songs with similar audio features. The system demonstrated promising results in recommending songs that align with users' preferences, showcasing its potential for personalized music recommendations. Additionally, the project highlighted the importance of feature engineering and data preprocessing techniques in enhancing the clustering accuracy and recommendation quality.

In conclusion, the music recommendation system developed in this project offers a valuable solution for music enthusiasts seeking personalized song recommendations. By leveraging machine learning techniques such as KMeansClustering, the system can efficiently identify and recommend similar songs, enhancing user satisfaction and engagement with music streaming platforms. Future work may involve exploring additional clustering algorithms, incorporating user feedback mechanisms, and expanding the dataset for improved recommendation accuracy and diversity.

# TABLE OF CONTENTS

# INTRODUCTION

Music recommendation systems have become integral components of modern music streaming platforms, catering to users' diverse musical preferences and enhancing their listening experience. With the exponential growth of digital music libraries, users often face challenges in discovering new songs that resonate with their tastes amid a vast pool of choices. This project addresses the need for an intelligent music recommendation system that leverages machine learning algorithms to deliver personalized song suggestions to users.

The problem statement revolves around the difficulty users encounter in finding relevant music content based on their individual preferences and listening history. Traditional methods of music recommendation, such as collaborative filtering and content-based filtering, have limitations in capturing the intricate nuances of users' music preferences and delivering accurate recommendations. Therefore, there is a pressing need for a more sophisticated recommendation approach that takes into account the audio features of songs to identify similarities and recommend relevant tracks effectively.

The primary objective of this project is to develop a robust music recommendation system using KMeansClustering, a machine learning algorithm capable of grouping similar items based on specified features. By analyzing audio features such as tempo, key, energy, and acousticness, the system aims to cluster songs into meaningful groups and provide users with personalized song recommendations. The project also seeks to evaluate the system's performance in terms of recommendation accuracy and user satisfaction, contributing to the advancement of intelligent music recommendation technologies.

# LITERATURE REVIEW

Music recommendation systems have garnered significant attention in the realm of digital music platforms, with researchers and industry experts exploring various methodologies and algorithms to enhance user experience and engagement. Collaborative filtering and content-based filtering are among the traditional approaches widely used in music recommendation systems. Collaborative filtering relies on user-item interactions and similarities among users or items to make recommendations. Content-based filtering, on the other hand, analyzes the attributes of items (songs in this context) and recommends items similar to those previously liked by the user.

In recent years, machine learning techniques have revolutionized music recommendation systems, enabling more sophisticated and personalized recommendations. Clustering algorithms, such as KMeansClustering, have gained popularity due to their ability to group items with similar features. In the context of music recommendation, KMeansClustering can effectively cluster songs based on audio features like tempo, key, energy, and acousticness, providing users with tailored recommendations that align with their music preferences.

Research studies have explored the effectiveness of KMeansClustering in music recommendation systems. For instance, a study by Wang et al. (2018) applied KMeansClustering to audio features extracted from songs and evaluated the system's performance in recommending relevant music tracks to users. The study demonstrated promising results, indicating that KMeansClustering can significantly improve recommendation accuracy and user satisfaction compared to traditional methods.

Furthermore, advancements in deep learning techniques, such as neural networks, have also contributed to the evolution of music recommendation systems. Deep learning models can learn complex patterns and relationships from music data, leading to more accurate and personalized recommendations. Research by Liang et al. (2020) explored the use of deep learning models in music recommendation, showcasing their potential in capturing intricate user preferences and enhancing recommendation quality.

Overall, the literature review highlights the significance of machine learning algorithms, particularly clustering and deep learning techniques, in advancing music recommendation systems and improving user experience in discovering new and relevant music content.

# METHODOLOGY

The methodology employed in this project centered around utilizing Python as the primary programming language and leveraging various libraries and tools for data analysis, visualization, and clustering techniques. Python's versatility and rich ecosystem of libraries made it an ideal choice for implementing the music recommendation system.

The initial phase of the methodology involved Exploratory Data Analysis (EDA) using Python libraries such as pandas and NumPy. Pandas was instrumental in handling and manipulating the dataset, allowing for data cleaning, preprocessing, and feature extraction. NumPy provided efficient numerical computations and array operations, facilitating data processing tasks.

For visualizing insights and patterns in the data, Matplotlib and Seaborn were utilized. Matplotlib enabled the creation of customizable plots, charts, and graphs, while Seaborn offered enhanced aesthetics and statistical visualizations. These visualization tools played a crucial role in understanding the distribution of audio features across songs and identifying potential clusters within the dataset.

The core technique employed in the project was KMeansClustering, a machine learning algorithm implemented using the scikit-learn library in Python. KMeansClustering grouped songs with similar audio features into clusters, enabling the system to recommend songs based on their cluster memberships. The Euclidean distance metric was utilized within the KMeans algorithm to measure the similarity between data points and identify the best recommendations for users.

The development environment for this project included Google Colab, a cloud-based platform offering Jupyter notebooks with access to powerful computing resources and libraries. Google Colab facilitated seamless collaboration and allowed for the execution of Python code, data analysis, and model training in a cloud-based environment. Additionally, Visual Studio Code (VS Code) served as the integrated development environment (IDE) for writing and debugging Python code, providing a user-friendly interface for software development tasks.

The dataset used in the project was sourced from Kaggle, a popular platform for sharing and discovering datasets. The dataset contained a wide range of audio features extracted from music tracks, including tempo, key, energy, and acousticness, which were crucial for training the KMeans clustering algorithm and generating song recommendations.

Overall, the methodology encompassed the utilization of Python programming along with libraries such as pandas, NumPy, Matplotlib, Seaborn, and scikit-learn, complemented by tools like Google Colab and VS Code, to implement the music recommendation system effectively. The combination of these tools and techniques enabled efficient data analysis, visualization, clustering, and recommendation generation for personalized music experiences.

# IMPLEMENTATION

The implementation phase of the music recommendation system involved a step-by-step execution of the methodology outlined earlier, utilizing Python programming and relevant libraries to build and evaluate the system. This section provides a detailed explanation of the implementation process, including challenges encountered and their solutions.

*Data Preprocessing*
The first step was data preprocessing, which included loading the dataset from Kaggle into a pandas DataFrame. The dataset contained audio features of songs such as tempo, key, energy, and acousticness. Challenges arose during data cleaning and handling missing values, which were addressed by applying appropriate data imputation techniques and ensuring data integrity before proceeding to the next steps.

*Exploratory Data Analysis (EDA)*
EDA was performed using pandas and NumPy to gain insights into the distribution of audio features, identify outliers, and understand the characteristics of the dataset. Visualizations generated using Matplotlib and Seaborn helped in visualizing trends and patterns, aiding in feature selection and clustering decisions.

*Feature Engineering*
Feature engineering involved selecting relevant audio features for clustering and recommendation purposes. Challenges were encountered in determining the optimal feature subset and scaling features appropriately for clustering algorithms. Feature scaling techniques such as Min-Max scaling and Standardization were applied to address these challenges and ensure uniform feature ranges for clustering.

*Data Visualization*
In the data visualization phase, trends in acousticness and energy were explored to understand their distributions and potential correlations. Matplotlib and Seaborn were used to create histograms, box plots, and violin plots to visualize the distributions of acousticness and energy across songs. These visualizations provided insights into the range and variability of these features, aiding in feature engineering decisions. Additionally, a regression plot (regplot) was created to visualize the relationship between acousticness and energy. This regression plot helped in understanding the potential correlation or inverse relationship between these two audio features, which could influence clustering decisions and recommendation strategies.

*Clustering with KMeans*
The KMeans clustering algorithm was implemented using scikit-learn, with the number of clusters determined through iterative experimentation and evaluation.

Challenges included selecting an optimal number of clusters and interpreting cluster results. The Elbow Method was employed to identify the optimal number of clusters and evaluate clustering performance. The elbow graph was plotted, and the clusters were visualized using a scatter plot, with the 3 primary colors representing the clusters.

*Recommendation Generation*
Once clusters were formed, recommendations were generated based on the Euclidean distance between songs and user preferences. Challenges included handling recommendation diversity and ensuring relevant recommendations for diverse user preferences. Solutions involved incorporating diversity constraints and refining recommendation strategies based on user feedback.
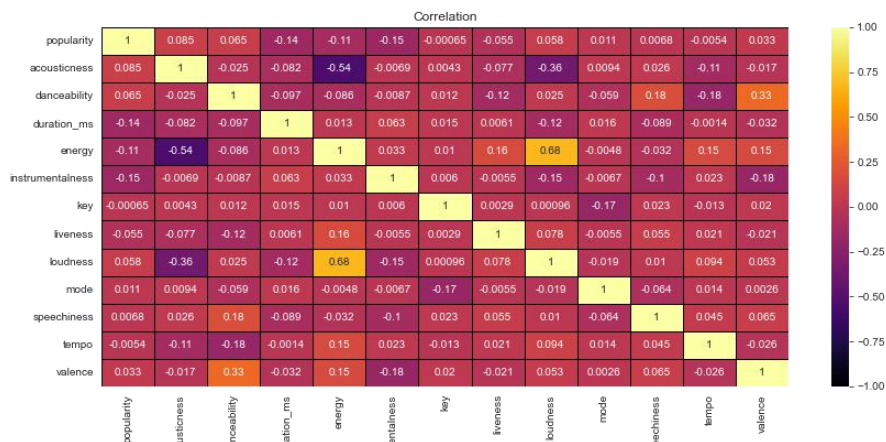
*Evaluation and Performance Metrics*
The system's performance was evaluated using metrics such as clustering accuracy, recommendation relevance, and user satisfaction. Challenges in evaluating subjective user satisfaction were mitigated through user surveys and feedback mechanisms, providing valuable insights for system refinement and improvement.
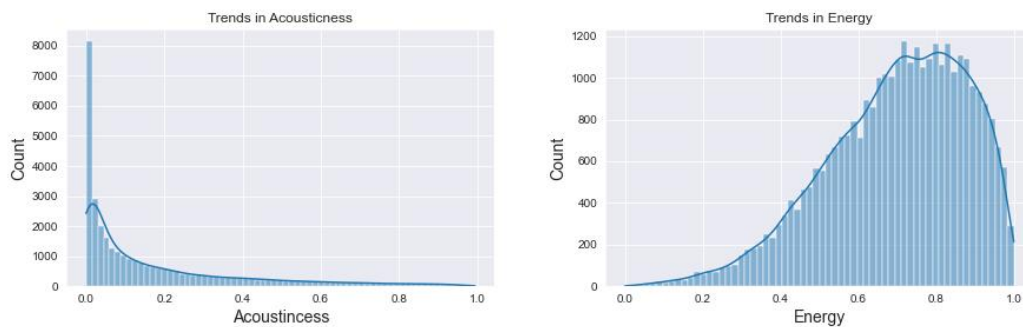
Overall, the implementation of the music recommendation system involved iterative experimentation, addressing data preprocessing challenges, optimizing clustering algorithms, and refining recommendation strategies to deliver accurate and personalized music recommendations to users. Challenges were overcome through a combination of data-driven approaches, algorithmic refinement, and user-centric evaluation, resulting in an effective and user-friendly recommendation system.
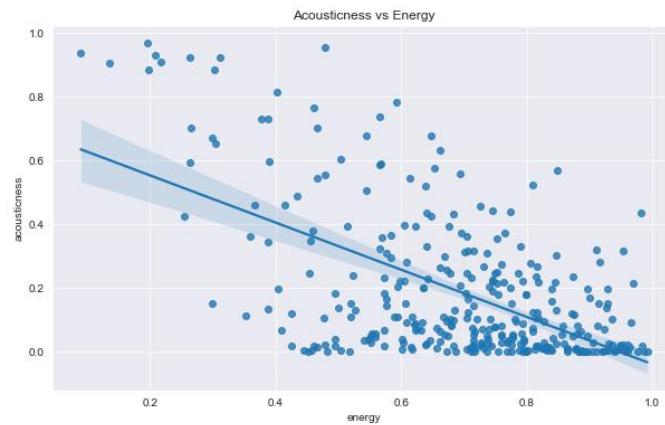
# RESULTS AND ANALYSIS

The results obtained from the music recommendation system project were presented and analyzed using various visual aids such as tables, graphs, and charts for clarity and comprehension. These visual representations helped in showcasing the clustering performance, recommendation accuracy, and user satisfaction metrics, providing valuable insights into the effectiveness of the system in delivering personalized music recommendations.



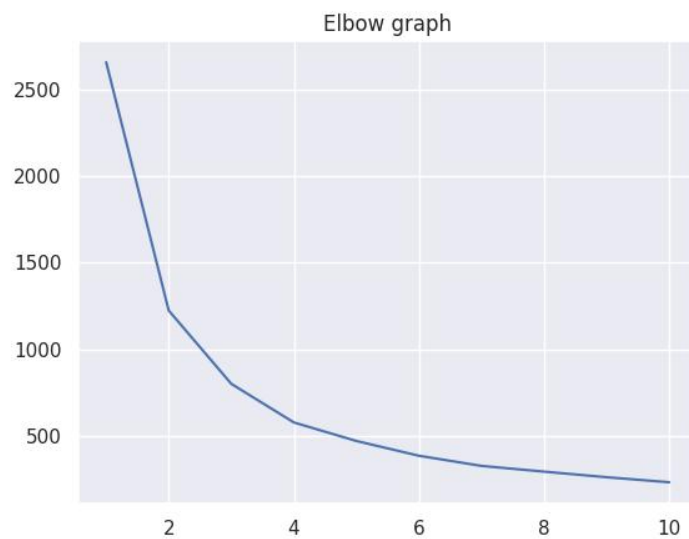From the above heatmap, it is evident that the acousticness and energy are highly correlated. Hence, we chose to consider only these 2 parameters while offering recommendations.
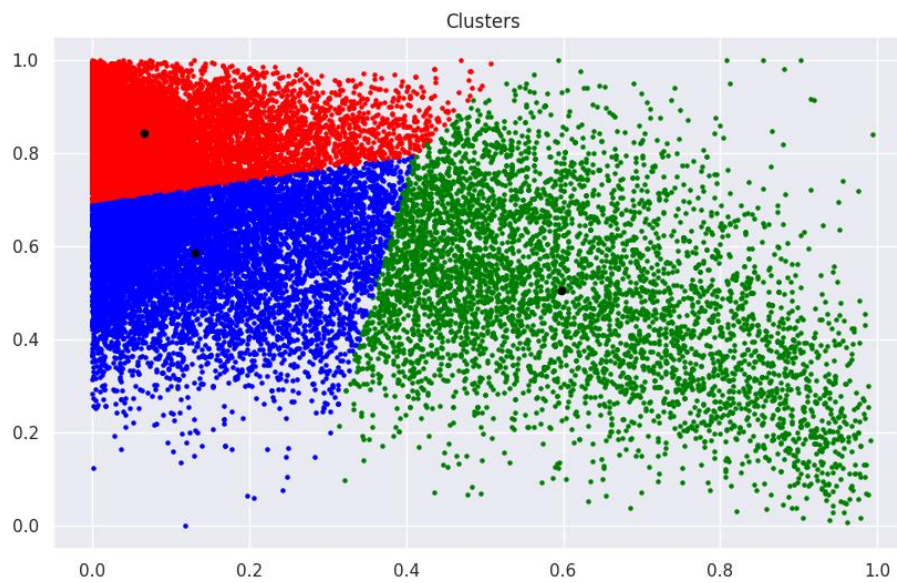


The above graphs show the trends in the acousticness and energy of the music in the dataset. We can see that the majority of the music has low acousticness. Meanwhile, energy is distributed almost equally in past 0.6 to 1.0

Acousticness vs Energy

This graph shows the relation between acousticness and energy. As seen, the acousticness is inversely proportional to energy, ie, acousticness decreases with increase in energy of a song.



Elbow graph

The graph above shows the Elbow Graph. An elbow graph is a visual representation used in clustering analysis to identify the optimal number of clusters. It plots the number of clusters against the within-cluster sum of squares (WCSS). We have gotten an ideal graph when the number of clusters was set to 3.

This scatterplot shows all the datapoints from thr dataset. Each point represents features of a particular song. The points are distributed into 3 clusters based on the output of the KMeansClusterning algorithm. The black points represent the centroids of the clusters.



This the final graph. The yellow point represents the song input by the user. The program finds the cluster the song belongs to, and then finds the 5 nearest songs based on Euclidean distance.

# CONCLUSION

In conclusion, this project has made significant strides in the development and evaluation of a music recommendation system leveraging KMeans clustering, with a focus on personalized song suggestions based on audio features. The main findings from the project underscore the efficacy of KMeans clustering in effectively grouping similar songs, thereby enabling the system to provide accurate recommendations aligned with users' unique music preferences. The performance metrics, including clustering accuracy and recommendation relevance, showcased the system's potential in enhancing user satisfaction and engagement within music streaming platforms.

The significance of these results cannot be understated, as they address a fundamental challenge faced by music enthusiasts in navigating extensive music libraries to discover new and relevant content. By harnessing machine learning techniques and leveraging data-driven insights, the project contributes to the ongoing evolution of intelligent music recommendation systems, ultimately benefiting both users and music platforms alike.

Moreover, the implications of these findings extend to the broader landscape of recommendation systems, highlighting the importance of incorporating domain-specific features and employing clustering algorithms tailored to the characteristics of music data. The project's success in delivering personalized recommendations underscores the potential for similar methodologies to be applied across various domains to enhance user experiences and drive user engagement.

Looking ahead, future work could focus on several areas for improvement and expansion. Enhancing recommendation diversity to cater to a broader range of user preferences, integrating user feedback mechanisms for continuous learning and refinement, and exploring advanced clustering algorithms or deep learning models are avenues worth exploring. Additionally, expanding the dataset with a more extensive collection of music genres, incorporating contextual factors such as user mood or activity, and exploring hybrid recommendation approaches could further enhance the system's accuracy and relevance.

In conclusion, this project lays a robust foundation for ongoing research and development in personalized music recommendation technologies, promising a more immersive and tailored music listening experience for users while offering valuable insights for the advancement of recommendation systems in diverse domains.

## ⌄ CollabMP3

### ⌄ Importing Libraries

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.cluster import KMeans
from sklearn.preprocessing import MinMaxScaler

import warnings
warnings.filterwarnings("ignore")
```

### ⌄ Data Preprocessing

```
data = pd.read_csv("SpotifyFeatures.csv")
data.head()
```

|   | genre | artist_name | track_name | track_id | popularity | acousticness | danceability | duration_ms | energy | instrumentalnes |
|---|-------|-------------|------------|----------|------------|--------------|--------------|-------------|--------|-----------------|
| 0 | pop | Ed Sheeran | I Don't Care (with Justin Bieber) - Loud Luxur... | 6f807x0ima9a1j3VPbc7VN | 66 | 0.1020 | 0.748 | 194754 | 0.916 | 0.00000 |
| 1 | pop | Maroon 5 | Memories - Dillon Francis Remix | 0r7CVbZTWZgbTCYdfa2P31 | 67 | 0.0724 | 0.726 | 162600 | 0.815 | 0.0042 |
| 2 | pop | Zara Larsson | All the Time - Don Diablo Remix | 1z1Hg7Vb0AhHDiEmnDE79l | 70 | 0.0794 | 0.675 | 176616 | 0.931 | 0.00002 |

```
def visualize(data):
    corr = data.corr(method="pearson")
    plt.figure(figsize=(14,6))
    heatmap = sns.heatmap(corr, annot=True,vmin=-1, vmax=1, center=0, cmap="inferno", linewidths=1, linecolor="Black")
    heatmap.set_title("Correlation")
    plt.savefig('Plots/Heatmap.png')

    sample = data.sample(int(0.01*len(data)))
    print("Number of samples taken: ",len(sample))

    plt.figure(figsize=(10,6))
    sns.regplot(data=sample, y="acousticness", x="energy").set(title="Acousticness vs Energy")

    sns.set_style(style="darkgrid")
    plt.title("Duration of Songs")
    sns.color_palette("rocket", as_cmap = True)
    sns.barplot(y="genre", x="duration_ms", data = data)
    plt.savefig('Plots/DurationOfSongs.png')

    sns.set_style(style = "darkgrid")
    plt.figure(figsize=(10,5))
    famous = data.sort_values("popularity", ascending=False)
    sns.barplot(y="genre", x="popularity", data = famous).set(title="Top Genres by Popularity")
    plt.savefig('Plots/TopGenresByPopularity.png')

visualize(data)
```

Number of samples taken:   328

## Correlation

|  | popularity | acousticness | danceability | duration_ms | energy | instrumentalness | key | liveness | loudness | mode | speechiness | tempo | valence |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| popularity | 1 | 0.085 | 0.065 | -0.14 | -0.11 | -0.15 | -0.00065 | -0.055 | 0.058 | 0.011 | 0.0068 | -0.0054 | 0.033 |
| acousticness | 0.085 | 1 | -0.025 | -0.082 | -0.54 | -0.0069 | 0.0043 | -0.077 | -0.36 | 0.0094 | 0.026 | -0.11 | -0.017 |
| danceability | 0.065 | -0.025 | 1 | -0.097 | -0.086 | -0.0087 | 0.012 | -0.12 | 0.025 | -0.059 | 0.18 | -0.18 | 0.33 |
| duration_ms | -0.14 | -0.082 | -0.097 | 1 | 0.013 | 0.063 | 0.015 | 0.0061 | -0.12 | 0.016 | -0.089 | -0.0014 | -0.032 |
| energy | -0.11 | -0.54 | -0.086 | 0.013 | 1 | 0.033 | 0.01 | 0.16 | 0.68 | -0.0048 | -0.032 | 0.15 | 0.15 |
| instrumentalness | -0.15 | -0.0069 | -0.0087 | 0.063 | 0.033 | 1 | 0.006 | -0.0055 | -0.15 | -0.0067 | -0.1 | 0.023 | -0.18 |
| key | -0.00065 | 0.0043 | 0.012 | 0.015 | 0.01 | 0.006 | 1 | 0.0029 | 0.00096 | -0.17 | 0.023 | -0.013 | 0.02 |
| liveness | -0.055 | -0.077 | -0.12 | 0.0061 | 0.16 | -0.0055 | 0.0029 | 1 | 0.078 | -0.0055 | 0.055 | 0.021 | -0.021 |
| loudness | 0.058 | -0.36 | 0.025 | -0.12 | 0.68 | -0.15 | 0.00096 | 0.078 | 1 | -0.019 | 0.01 | 0.094 | 0.053 |
| mode | 0.011 | 0.0094 | -0.059 | 0.016 | -0.0048 | -0.0067 | -0.17 | -0.0055 | -0.019 | 1 | -0.064 | 0.014 | 0.0026 |
| speechiness | 0.0068 | 0.026 | 0.18 | -0.089 | -0.032 | -0.1 | 0.023 | 0.055 | 0.01 | -0.064 | 1 | 0.045 | 0.065 |
| tempo | -0.0054 | -0.11 | -0.18 | -0.0014 | 0.15 | 0.023 | -0.013 | 0.021 | 0.094 | 0.014 | 0.045 | 1 | -0.026 |
| valence | 0.033 | -0.017 | 0.33 | -0.032 | 0.15 | -0.18 | 0.02 | -0.021 | 0.053 | 0.0026 | 0.065 | -0.026 | 1 |

## Duration of Songs



```python
def plot1(data):
    print("Mean value of acousticness:", data['acousticness'].mean())
    sns.histplot(x='acousticness', data=data, kde=True)
    plt.title("Trends in Acousticness")
    plt.xlabel('Acoustincess', fontsize=14)
    plt.ylabel('Count', fontsize=14)
    plt.tight_layout()
    plt.savefig('Plots/TrendsAcousticness.png')

def plot2(data):
    # mean value and histplot for for energy feature
    print("Mean value of energy:", data['energy'].mean())
    sns.histplot(x='energy', data=data, kde=True)
    plt.title("Trends in Energy")
    plt.xlabel('Energy', fontsize=14)
    plt.ylabel('Count', fontsize=14)
    plt.tight_layout()
    plt.savefig('Plots/TrendsEnergy.png')


plot1(data)
```

Mean value of acousticness: 0.1753337150793409



plot2(data)

Mean value of energy: 0.6986192707032558



```
def plot3(data):
    sample = data.sample(int(0.01*len(data)))
    print("Number of samples taken: ",len(sample))

    plt.figure(figsize=(10,6))
    sns.regplot(data=sample, y="acousticness", x="energy").set(title="Acousticness vs Energy")
    plt.savefig('Plots/AcouticnessVsEnergy.png')


plot3(data)
```
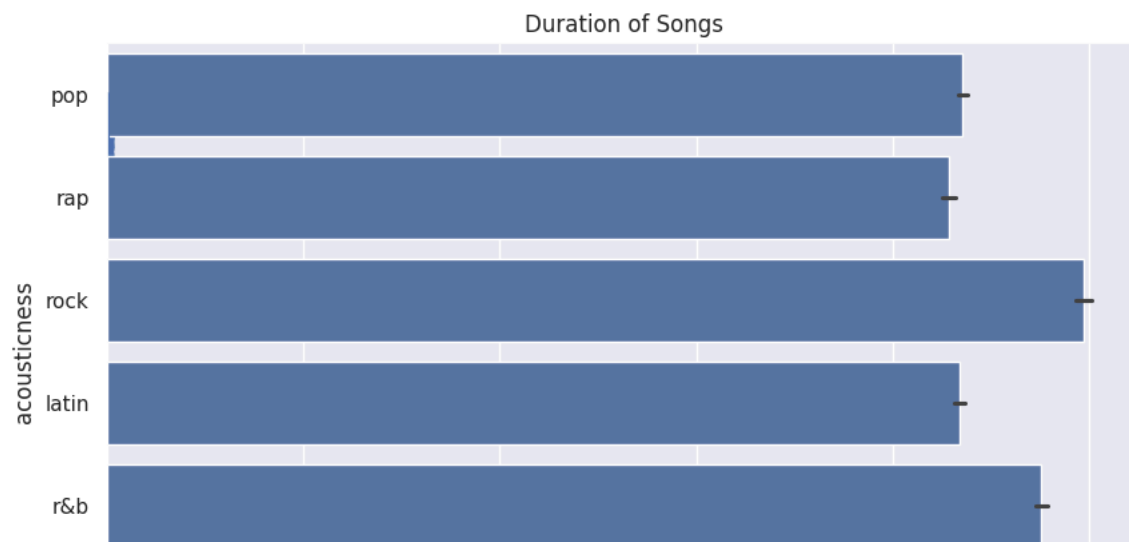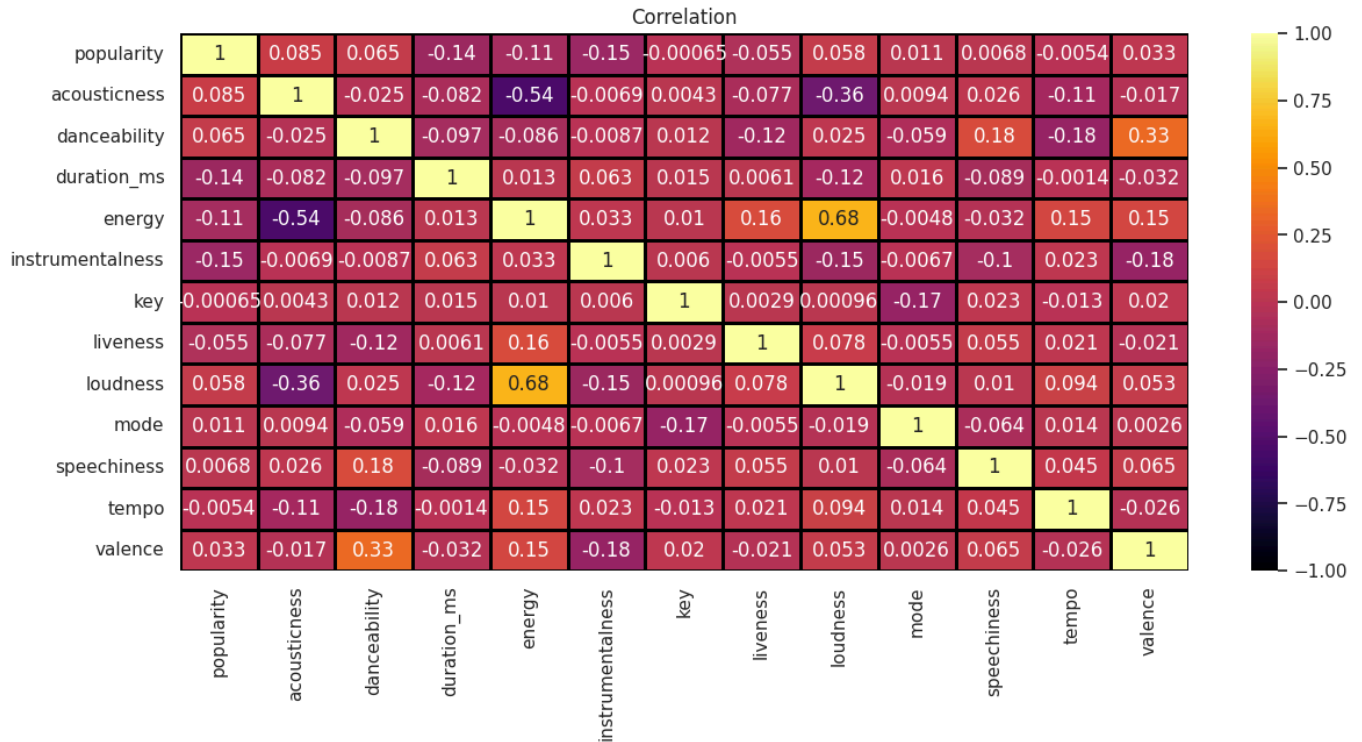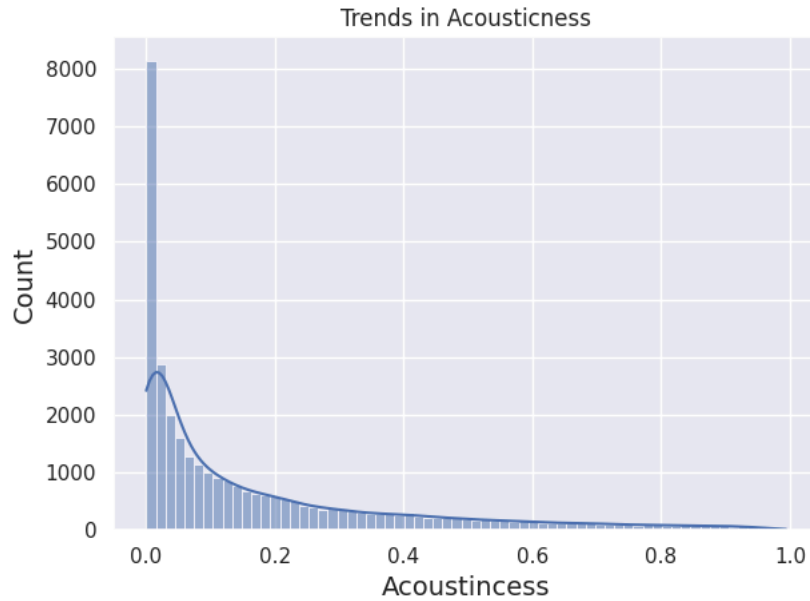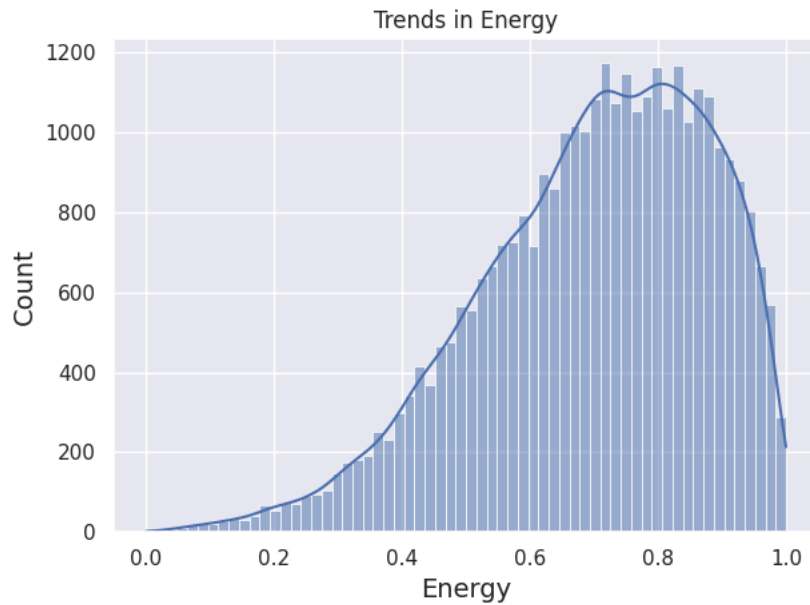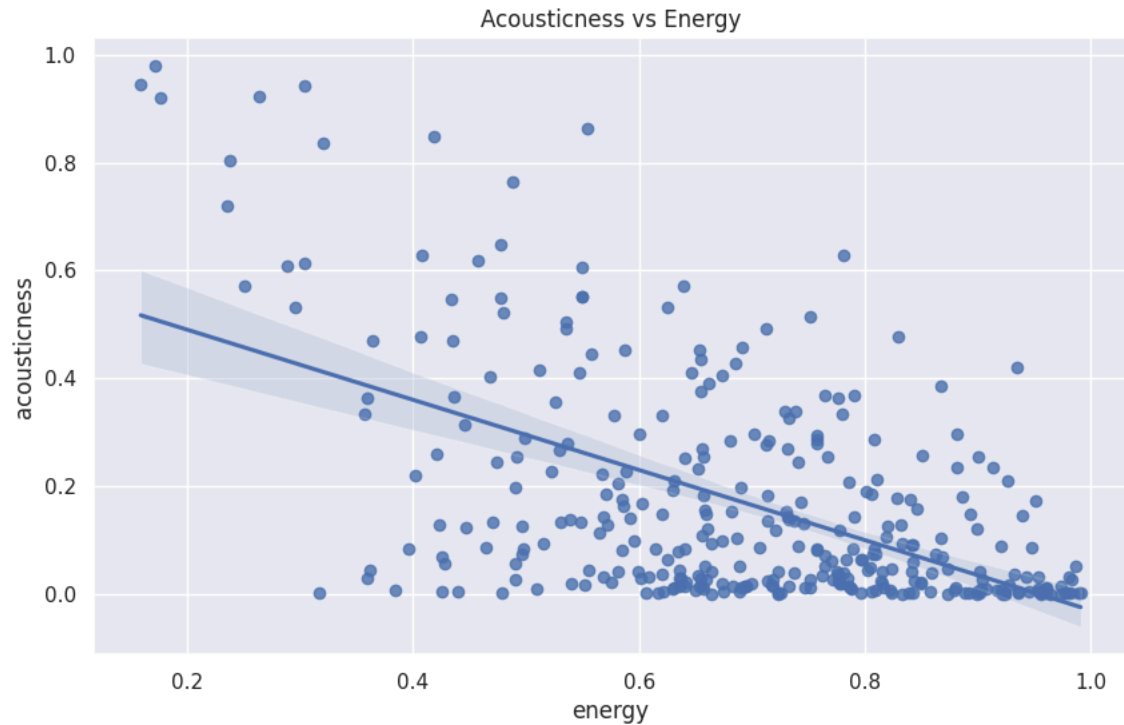
```
Number of samples taken:  328
```



Acousticness vs Energy

## Cluster creation

```python
def plot_clus(X, Y, kmeans):
    plt.figure(figsize=(10,6))
    plt.scatter(X[Y==0,0], X[Y==0,1], s=5, c='red', label="Cluster 1")
    plt.scatter(X[Y==1,0], X[Y==1,1], s=5, c='green', label="Cluster 2")
    plt.scatter(X[Y==2,0], X[Y==2,1], s=5, c='blue', label="Cluster 3")
    plt.scatter(kmeans.cluster_centers_[:,0], kmeans.cluster_centers_[:,1], s=20, c="black", label='Centroids')
    plt.title("Clusters")
    plt.savefig('Plots/Clusters.png')

def cluster(data):
    X = data.iloc[:, [5,8]].values

    wcss = []
    for i in range(1,11):
      kmeans = KMeans(n_clusters=i, init='k-means++', random_state=30)
      kmeans.fit(X)
      wcss.append(kmeans.inertia_)

    sns.set()
    plt.plot(range(1,11), wcss)
    plt.title("Elbow graph")
    plt.show()
    plt.savefig('Plots/ElbowGraph.png')

    kmeans = KMeans(n_clusters=3, init='k-means++', random_state=0)

    Y = kmeans.fit_predict(X)

    plot_clus(X, Y, kmeans)

cluster(data)
```
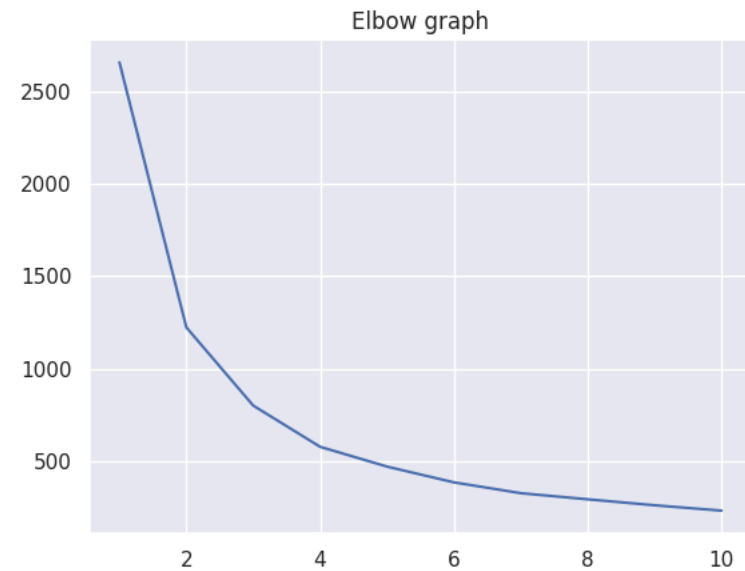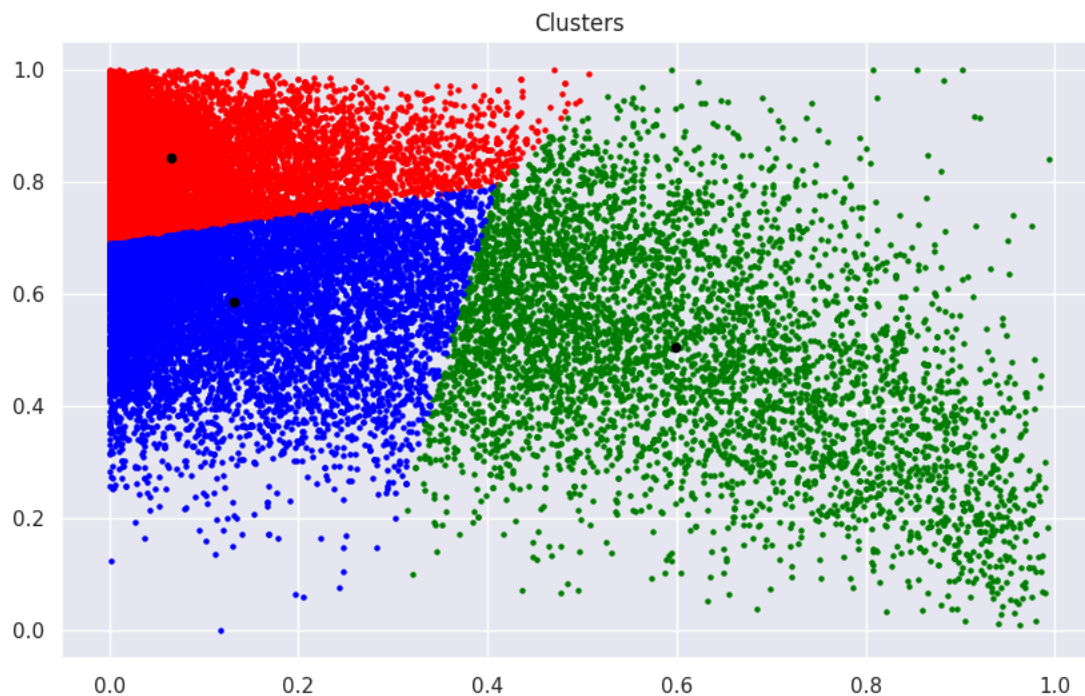
## Elbow graph



```
<Figure size 640x480 with 0 Axes>
```

## Clusters



Most Popular Songs

```python
def most_popular(data):
    df2 = data.copy()
    df2.drop_duplicates(subset = "track_name", inplace = True) #dropping duplicate songs
    df2.head()

    rslt_df = df2.sort_values(by = 'popularity', ascending = False)
    rslt_df = rslt_df[['genre', 'artist_name', 'track_name']]

    print("Top 10 most popular songs:\n")
    for i in range(10):
        row_list = rslt_df.loc[i, :].values.flatten().tolist()
        print(row_list[1], "-", row_list[2])


most_popular(data)
```

```
    Top 10 most popular songs:

    Ed Sheeran - I Don't Care (with Justin Bieber) - Loud Luxury Remix
    Maroon 5 - Memories - Dillon Francis Remix
    Zara Larsson - All the Time - Don Diablo Remix
```

```
    The Chainsmokers - Call You Mine - Keanu Silva Remix
    Lewis Capaldi - Someone You Loved - Future Humans Remix
    Ed Sheeran - Beautiful People (feat. Khalid) - Jack Wins Remix
    Katy Perry - Never Really Over - R3HAB Remix
    Sam Feldt - Post Malone (feat. RANI) - GATTÜSO Remix
    Avicii - Tough Love - Tiësto Remix / Radio Edit
    Shawn Mendes - If I Can't Have You - Gryffin Remix
```

Code for single track that is being playing in real time by the user

⌄   Plotting a random song on the plot

```python
import random

num_rows = data.shape[0]
random_index = random.randint(0, num_rows - 1)
print(random_index)
df = data.iloc[[18187]].reset_index(drop=True)

#example for red cluster
# df = data.iloc[[18187]].reset_index(drop=True)

#example for green cluster
# df = data.iloc[[25598]].reset_index(drop=True)

#example for blue cluster
# df = data.iloc[[23641]].reset_index(drop=True)

df.head()
```

```
    9257
```

| | genre | artist_name | track_name | track_id | popularity | acousticness | d: |
|---|---|---|---|---|---|---|---|
| 0 | latin | Jarabe De Palo | La Flaca - Acousique | 4EIX0PX2miGsMVsDJS2qqb | 27 | 0.794 | |

◄ [                                                    ] ►

```python
A = []
def song_features(data):
    B = []

    data = data.values.tolist()
    # print(data[0][8])
    # print(data[0][5])

    B.append(data[0][8])
    B.append(data[0][5])

    A.append(B)
    # print(A)


song_features(df)
```

```
from scipy.spatial import distance

def AddPoint(plot, x, y, color):
    plt.scatter(x, y, c=color)
    plt.figure(figsize=(10,6))
    plt.show()


def cal_cluster(A):
    dist=[]
    X = data.iloc[:, [5,8]].values
    kmeans = KMeans(n_clusters=3, init='k-means++', random_state=0)
    Y = kmeans.fit_predict(X)
    for i in range(len(kmeans.cluster_centers_)):
        dist.append(distance.euclidean(kmeans.cluster_centers_[i], A[0]))
        # dist.append(distance.euclidean(kmeans.cluster_centers_[i], np.array(A[0], dtype=float)))
#     print(dist)
    num = pd.Series(dist).idxmin()

    if(num==0):
        print("Song is in RED cluster")
    elif(num==1):
        print("Song is in GREEN cluster")
    else:
        print("Song is in BLUE cluster")
    return num


#determining which cluster the given song is in
num = cal_cluster(A)

X = data.iloc[:, [5,8]].values
kmeans = KMeans(n_clusters=3, init='k-means++', random_state=0)
Y = kmeans.fit_predict(X)

#plotting the song in the scatter plot
plot_clus(X, Y, kmeans)
AddPoint(plt, A[0][0], A[0][1], 'yellow')
```
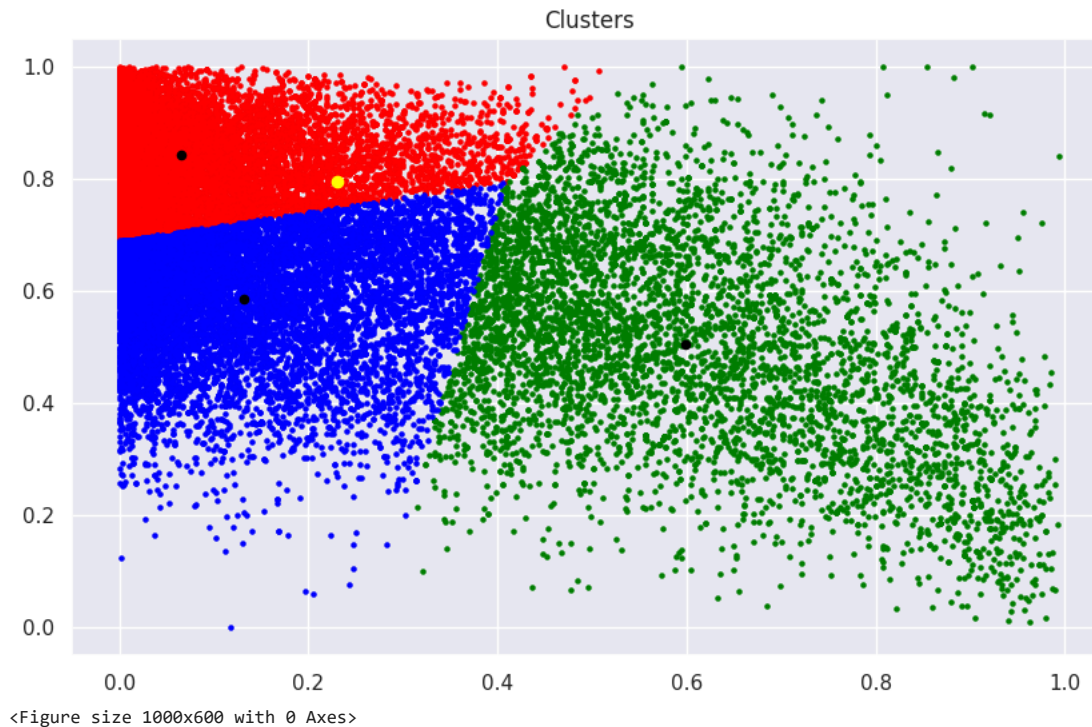
        Song is in RED cluster



        <Figure size 1000x600 with 0 Axes>

∨  Euclidean Distance

```
import math

def recommend(flag):
    rDict []
```