# Spotify song clustring analysis

## Problem Description :

**Problem Description:**

As music lovers, we often find ourselves with a vast collection of songs from various genres, artists, and moods. However, organizing these songs into different clusters based on their characteristics can be a challenging task. That's where Spotify Song Cluster Analysis comes to the rescue!

**Background Information:**

Spotify is a popular music streaming service that provides access to millions of songs from various genres and artists. Each song in Spotify's database is associated with several audio features like energy, danceability, loudness, and more. These features represent different aspects of the song's musical composition and are essential for understanding its overall vibe and mood.

**Dataset Information:**

**For this project, we have a dataset containing information about the top 50 songs on Spotify. The dataset includes the following columns:**

1. **Track.Name**: The name of the song.
2. **Artist.Name:** The name of the artist who performed the song.
3. **Genre:** The genre of the song.
4. **Beats.Per.Minute:** The number of beats per minute in the song.
5. **Energy:** Represents the energy level of the song (ranging from 0 to 100).
6. **Danceability:** Indicates how suitable the song is for dancing (ranging from 0 to 100).

7. **Loudness:** The loudness of the song in decibels.
8. **Liveness:** Reflects the presence of an audience in the recording (ranging from 0 to 100).
9. **Valence:** Describes the musical positiveness of the song (ranging from 0 to 100).
10. **Length:** The duration of the song in seconds.
11. **Acousticness:** Represents the acoustic quality of the song (ranging from 0 to 100).
12. **Speechiness:** Indicates the presence of spoken words in the song (ranging from 0 to 100).
13. **Popularity:** A popularity score assigned by Spotify to the song.

**Project Goal:**

The goal of this project is to analyze the top 50 Spotify songs and perform clustering to group similar songs together based on their audio features. Clustering will help us create playlists with songs that share similar characteristics, making it easier for users to find songs that match their current mood or preferences. This way, users can enjoy a more personalized and organized music listening experience on Spotify.

By using data visualization techniques and clustering algorithms, we can uncover hidden patterns and insights within the songs, providing valuable information to both music enthusiasts and Spotify's recommendation system.

# Possible Framework :

The code provided is for performing Spotify Song Cluster Analysis on a dataset containing information about the top 50 songs on Spotify. The goal is to analyze the audio features of these songs and group them into clusters based on their similarities. Let's break down the code outline step-by-step:

**1. Importing Necessary Libraries:**

- Import the required Python libraries such as pandas, numpy, seaborn, matplotlib, scipy, scikit-learn, statsmodels, and plotly for data manipulation, visualization, and analysis.

**2. Loading the Dataset:**

- Read the dataset containing information about the top 50 Spotify songs using pandas read_csv() function.
- Check the dimensions of the dataset and remove any rows with missing values if necessary.

**3. Exploratory Data Analysis (EDA):**

- Perform EDA to understand the dataset and its features.
- Visualize the correlation between different audio features using a heatmap.
- Use scatter plots and pair plots to identify relationships between selected features.

**4. Feature Engineering:**

- Group some genres together to create a new feature called GeneralGenre.
- Standardize the numeric features to ensure they are on the same scale for clustering purposes.

**5. K-Means Clustering:**

- Use the K-means clustering algorithm to group similar songs into clusters.
- Determine the optimal number of clusters using the Elbow method and PCA variance.
- Assign each song to a cluster based on the K-means results.

## 6. Data Visualization of Clusters:

- Visualize the clusters in 3D space using a scatter plot with different colors representing each cluster.
- Merge the cluster information with the original dataset to examine the relationship between clusters, genres, and general genres.

## 7. Hierarchical Clustering (Optional):

- Perform hierarchical clustering to explore additional clustering possibilities.

## 8. Conclusion and Insights:

- Summarize the findings from the clustering analysis.
- Provide insights into the characteristics of songs in each cluster and how they differ from one another.
- Discuss the potential use cases for these clusters in music recommendation and playlist generation on Spotify.

## 9. Future Work and Improvements:

- Suggest potential improvements or additional analyses that can be performed to enhance the clustering results and gain more insights.
- Discuss the significance of incorporating user preferences and feedback in refining the clustering model.

## 10. Final Thoughts:

- Wrap up the project by emphasizing the importance of data-driven music recommendation systems and how this clustering analysis can benefit both users and the Spotify platform.

# Code Explanation :

*If this section is empty, the explanation is provided in the .ipynb file itself.

**Step 1: Importing Necessary Libraries**

- We start by importing various Python libraries that will help us with data manipulation, visualization, and analysis. Libraries like pandas, numpy, seaborn, matplotlib, and plotly are common tools for working with data in Python.

**Step 2: Loading the Dataset**

- The next step is to load the dataset containing information about the top 50 songs on Spotify. This dataset holds details about each song, such as its title, artist, genre, and various audio features like beats per minute, energy, loudness, and more. We use the pandas read_csv() function to read the data from a CSV file into a pandas DataFrame, which is a data structure used to hold tabular data.

**Step 3: Exploratory Data Analysis (EDA)**

- In this step, we perform Exploratory Data Analysis (EDA) to get a better understanding of our dataset. EDA helps us to explore the data, identify patterns, and spot any potential issues or interesting insights. We use various visualization techniques like heatmaps and scatter plots to analyze the relationships between different audio features of the songs.

**Step 4: Feature Engineering**

- Feature Engineering involves creating new features or modifying existing ones to make our data more informative. In this project, we group similar genres into broader categories, creating a new feature called GeneralGenre. This simplifies our analysis and makes it easier to identify similarities between different genres.

**Step 5: K-Means Clustering**

- This is the main part of our project where we apply the K-means clustering algorithm. Clustering is a technique that groups similar items together based on their characteristics. In our case, we group songs into clusters based on their

audio features. The K-means algorithm works by repeatedly dividing the dataset into K clusters until each song is assigned to the closest cluster.

## Step 6: Data Visualization of Clusters

- After performing K-means clustering, we visualize the results to better understand how the songs are grouped into clusters. We use a 3D scatter plot to represent the songs in a three-dimensional space, where each axis represents a different audio feature, and each song is colored according to its cluster.

## Step 7: Hierarchical Clustering (Optional)

- In addition to K-means, the code also provides an optional step for hierarchical clustering. Hierarchical clustering is another way of grouping data into clusters by building a hierarchy of clusters. Although optional, it can be a useful alternative or complement to K-means in some cases.

## Step 8: Conclusion and Insights

- In the end, we summarize our findings and insights from the clustering analysis. We discuss how songs within each cluster are similar and how they differ from songs in other clusters. These insights can be valuable for music recommendation and playlist generation on Spotify.

## Step 9: Future Work and Improvements

- The code suggests potential improvements and additional analyses that can be performed to enhance the clustering results. It highlights the importance of incorporating user preferences and feedback to improve the clustering model and provide more personalized music recommendations.

## Step 10: Final Thoughts

- The project concludes with a reflection on the importance of data-driven music recommendation systems and how this clustering analysis can benefit both users and the Spotify platform. It encourages further exploration and application of data science techniques in the field of music and entertainment.

In summary, the code takes us through the process of clustering Spotify songs based on their audio features, allowing us to discover patterns and similarities among the top 50

songs. This analysis can provide valuable insights for both music enthusiasts and music streaming platforms like Spotify, making the music listening experience more enjoyable and personalized.

# Future Work :

Clustering analysis is an ongoing process, and there are several potential areas for future work and improvements. Here's a step-by-step guide on how to implement the future work for the project:

**Step 1: Data Collection and Expansion**

- Gather a more extensive and diverse dataset of songs from various genres and artists. Expanding the dataset will provide a broader representation of different music styles, which can lead to more accurate and meaningful clusters.

**Step 2: Feature Selection and Engineering**

- Conduct a more in-depth analysis of the audio features and explore additional attributes that could contribute to the clustering process. Experiment with different feature engineering techniques to extract more relevant information from the audio data.

**Step 3: Scaling and Normalization**

- Apply different scaling and normalization methods to the audio features. Experiment with different normalization techniques and observe their impact on the clustering results.

**Step 4: Dimensionality Reduction**

- Implement dimensionality reduction techniques like Principal Component Analysis (PCA) or t-distributed Stochastic Neighbor Embedding (t-SNE) to visualize high-dimensional data and reduce computational complexity.

**Step 5: Optimal Number of Clusters**

- Determine the optimal number of clusters using various clustering validation metrics such as the Elbow Method, Silhouette Score, or Davies-Bouldin Index. Finding the right number of clusters can improve the interpretability of the results.

**Step 6: Advanced Clustering Algorithms**

- Explore other clustering algorithms like DBSCAN, Agglomerative Hierarchical Clustering, or Gaussian Mixture Models (GMM) to compare their performance with K-means.

## Step 7: Evaluate Cluster Quality

- Assess the quality of the clusters by performing a qualitative analysis of the songs within each cluster. Identify common patterns, genres, or artist associations to validate the effectiveness of the clustering.

## Step 8: Cluster Labeling

- Assign meaningful labels to each cluster to describe the characteristics of the songs within it. This will make the interpretation of the clusters more intuitive and user-friendly.

## Step 9: User Feedback Integration

- Incorporate user feedback and preferences into the clustering process. Collect user ratings or feedback on songs to fine-tune the clustering model and provide more personalized music recommendations.

## Step 10: Real-Time Clustering and Recommendation

- Implement real-time clustering and recommendation systems for streaming platforms like Spotify. Continuously update the clustering model based on user interactions and provide up-to-date song recommendations.

## Step 11: Music Playlist Generation

- Use the clustering results to create thematic playlists that group similar songs together. Implement an algorithm to generate playlists based on user preferences and cluster assignments.

## Step 12: Evaluation and A/B Testing

- Conduct A/B testing to evaluate the impact of the clustering-based music recommendation system on user engagement and satisfaction. Compare the performance of the clustering model with traditional recommendation methods.

## Step 13: Interpretability and Visualization

- Develop interactive visualizations to showcase the clustering results to users. Allow users to explore the clusters and discover new songs within each category.

**Step 14: Deployment and Productionization**

- Deploy the clustering model and recommendation system in a production environment. Optimize the system for scalability and efficiency to handle large-scale data and user requests.

**Step 15: Continuous Improvement**

- Continuously monitor and evaluate the clustering model's performance over time. Keep updating the model with new data and feedback to improve the accuracy and relevance of song clusters.

In conclusion, the future work for the Spotify Song Clustering Analysis involves data expansion, feature engineering, advanced clustering techniques, and user feedback integration to create a more personalized and engaging music recommendation system. By following these steps, we can enhance the clustering model and provide Spotify users with better music discovery and playlist generation experiences.

# Concept Explanation :

Alright, party animals, buckle up for a fun ride into the world of K-Means Clustering! Imagine you are hosting a wild party, and your friends come in with all their different dance moves and party vibes. You want to group them based on their dance styles and energy levels so that you can create the ultimate dance floor experience! That's where K-Means Clustering comes to the rescue!

**What is K-Means Clustering?** K-Means Clustering is like a DJ for your party who takes all your friends and divides them into different dance crews based on their moves. The "K" in K-Means represents the number of dance crews you want to create. Each dance crew is called a "cluster," and the DJ tries to make sure that friends with similar dance styles end up in the same cluster.

**How K-Means Works:** Here's how the DJ (K-Means) works the magic:

1. **Picking the Dance Crews (Centroids):** The DJ starts by randomly picking "K" friends to be the dance crew leaders. These leaders are called "centroids." They will be the ones leading each dance crew.
2. **Grooving with the Beat (Data Points):** Next, the DJ looks at each friend's dance style and energy level. Each friend represents a "data point" in the dance floor space.
3. **Joining the Crews (Assigning Data Points):** The DJ then looks at each friend and decides which dance crew (centroid) they should join. The DJ does this by measuring the distance between each friend and the centroids. Friends will join the dance crew with the closest centroid. It's like joining the crew with the dance style that matches their moves!
4. **Shaking It Up (Updating Centroids):** Now, the DJ asks each dance crew to huddle up and calculates the average dance style and energy level of all the friends in the crew. This new average becomes the new position of the centroid. The DJ keeps shaking things up until the dance crews' leaders (centroids) stop moving.
5. **Ultimate Dance Floor! (Final Clustering):** Finally, when the dance crews' leaders (centroids) stop moving, all your friends have found their ultimate dance floor buddies! Each crew represents a cluster of friends with similar dance styles and

energy levels. The DJ has created the ultimate dance floor experience for everyone!

**Example** - Let's Party with K-Means! Let's say you have ten friends at your party, and you want to create three dance crews (clusters). Each friend has two qualities: dance style and energy level. Now, the DJ (K-Means) starts the magic! The DJ randomly picks three friends as the initial dance crew leaders (centroids). Then, the DJ assigns each friend to the dance crew with the closest dance style and energy level. The crews huddle up, and the DJ recalculates the average dance style and energy level for each crew. The dance crews adjust their moves and positions, and the DJ keeps repeating this process until the crews' leaders (centroids) stop moving. In the end, each friend ends up in one of the three dance crews based on their dance style and energy level! Voila! We have our awesome dance floor experience thanks to K-Means Clustering!

So, the next time you are hosting a party and want to group your friends based on their dance moves and vibes, remember the DJ (K-Means) is there to help you create the ultimate dance floor experience! Happy clustering and partying! 🎉🎉🎉

# Exercise Questions :

**Exercise 1: Understanding the Dataset**

1. **What is the size of the dataset, and how many features (columns) are present?**

**Answer**: The dataset contains 50 rows (instances) and 14 columns (features).

2. **Are there any missing values in the dataset? If yes, how are they handled?**

**Answer:** Yes, there were missing values in the dataset, but they were dropped using the dropna() function.

3. **How many different genres are present in the dataset, and what are their counts?**

**Answer:** There are several different genres present in the dataset, and the counts for each genre have been calculated using the groupby() function.

4. **What is the correlation between different numeric features in the dataset?**

**Answer:** The correlation heatmap and pairplot have been created using sns.heatmap() and sns.pairplot() functions to visualize the correlation between numeric features.

**Exercise 2: Visualizing Data and Feature Engineering**

5. **How are house prices distributed across postal codes?**

**Answer**: The distribution of house prices across postal codes has been visualized using a box plot with the sns.boxplot() function.

6. **What is the relationship between house size and other features such as bedrooms, bathrooms, etc.?**

**Answer:** The relationship between house size (living area) and other features has been visualized using a scatterplot matrix with the sns.pairplot() function.

7. **How does the age of the house and the renovation year affect house prices?**

**Answer:** The age and renovation year's effect on house prices have been visualized using a scatterplot with the sns.scatterplot() function.

## Exercise 3: Clustering Similar Properties

8.  **What features were selected for clustering similar properties, and why were they chosen?**

**Answer:** Features such as living area, number of views, bedrooms, bathrooms, etc., were selected for clustering similar properties as they play a significant role in determining property similarity.

9.  **How was the optimal number of clusters determined for K-Means?**

**Answer:** The optimal number of clusters (dance crews) was determined using the elbow method, visualized with a line plot and the KElbowVisualizer from the yellowbrick.cluster module.

## Exercise 4: Modeling & Prediction

10. **What models were used for predicting house prices, and how was their performance compared?**

**Answer:** Random Forest, Support Vector Machine (SVM), and XGBoost models were used to predict house prices. Their performance was compared by calculating the root mean squared error (RMSE) on the test set.