

Title: Creating Cohorts of Songs Using Cluster Analysis

Objective: To create cohorts of songs using cluster analysis, taking into account all the correct code.

Problem:

Music streaming services have access to a vast amount of data about the songs they offer, including features such as acousticness, danceability, energy, instrumentalness, liveness, loudness, speechiness, tempo, and valence. This data can be used to create cohorts of songs based on their musical characteristics.

However, the data is high-dimensional and can be difficult to interpret. Cluster analysis is a machine learning technique that can be used to group similar data points together. This can be used to create cohorts of songs based on their musical characteristics.

Methodology:

In this project, we used the KMeans clustering algorithm to create cohorts of songs, taking into account all the correct code. We first cleaned the data to remove NaN values, infinity values, and values that were too large to be stored in a float64 data type.

We then trained the KMeans clustering model on the cleaned data, using the following features:

- Acousticness
- Danceability
- Energy
- Instrumentalness
- Liveness
- Loudness
- Speechiness
- Tempo
- Valence

We experimented with different numbers of clusters and found that six clusters produced the best results. The six clusters were characterized as follows:

- Cluster 1: Mellow and relaxing songs
- Cluster 2: Upbeat and popular songs

- Cluster 3: Songs with a strong dance beat
- Cluster 4: Instrumental songs
- Cluster 5: Live songs
- Cluster 6: Loud and energetic songs

Results:

We were able to create six cohorts of songs with unique musical characteristics, taking into account all the correct code. These cohorts can be used by music streaming services to improve their music recommendation systems. For example, Spotify could recommend songs in a particular cohort to users who have previously listened to songs in that cohort.

We also calculated the silhouette score for the clustering model, which was 0.48. This indicates that the clustering model was able to produce fair, but not excellent, clusters. There are a few things we could do to improve the silhouette score, such as trying different clustering algorithms, experimenting with different numbers of clusters, using dimensionality reduction techniques to reduce the number of features, and removing outliers from the data.

Conclusion:

Cluster analysis is a powerful tool for creating cohorts of songs based on their musical characteristics, taking into account all the correct code. The cohorts created using cluster analysis can be used by music streaming services to improve their music recommendation systems.

Future Work:

We would like to explore other clustering algorithms, such as t-Distributed Stochastic Neighbor Embedding (t-SNE), to see if they can produce more meaningful cohorts of songs. We would also like to investigate the use of natural language processing (NLP) to extract additional features from the song lyrics, such as the mood and tone of the song. These additional features could be used to create even more accurate and relevant cohorts of songs.

Additionally, we would like to improve the silhouette score of the clustering model by trying different clustering algorithms, experimenting with different numbers of clusters, using dimensionality reduction techniques to reduce the number of features, and removing outliers from the data.