**Project Proposal**
SYS 6018: Data Mining
Prof. Matthew S. Gerber
TA: Peter Wu

**Team Members**
Sakshi Jawarani
Aman Shrivastava
sj8em, as3ek

## Problem Statement

The aim of the project is to build an unsupervised framework to automatically generate playlists from a given set of songs.

Major music streaming platforms continually look to improve their products and develop features that lead to a more personalized user experience. Playlists are an integral part of their services. With respect to automated playlist generation, currently music platforms generate playlists based on similar artists or genres. This project aims to generate playlists automatically only from a given set of songs without using existing playlists as training data, using a set of diverse qualitative and quantitative features, ranging from song lyrics to features that parameterize the musical properties of songs. Users want a premium service that excels at building playlists, hence premium music streaming services would care about this problem.

## Objectives And Metrics

We approach this unsupervised non parametric problem by using clustering methods in which songs are grouped based on their attribute similarity. The objective is to cluster songs such that the songs in the same cluster are as similar as possible, and the songs in different clusters are highly distinct. The average distance within the cluster should be as small as possible and the average distance between clusters to be as large as possible. The Dunn Index[1] is one such metric that assesses the goodness of a clustering, by measuring the maximal diameter of clusters and relating it to the minimal distance between clusters. If the data set contains compact and well-separated clusters, the diameter of the clusters is expected to be small and the distance between the clusters is expected to be large. Thus, Dunn index should be maximized.

## Alternative Solutions

The data set that we will be using for this problem will be an amalgamation of multiple music related datasets. Further relevant data will be collected using public APIs offered by music streaming services. Existing research in this domain is majorly focused on approaching this as a supervised classification problem. The songs in the data are tagged based on existing playlists and the model is built to classify a new song to one of those existing playlists based on historical tagging[2] or playlists are generated using seed songs[3]. The various supervised learning techniques that have been explored for this problem include linear and tree based classification methods. We believe that this problem is better dealt as an unsupervised problem since that would enable us to generate playlists that contain similar songs more organically.

## Hypotheses and Approach

Looking at other approaches in this domain, we realize that using existing playlists to classify songs into one of them inadvertently introduces a certain amount of bias in the method stemming from the historically tagged data. Additionally, some features in the metadata of the songs

(specifically the genre) get disproportionate importance while classifying songs into playlists as the training data available largely is clustered into playlists based on one of these features. We aim to develop a more robust and data-driven approach that minimizes these biases by staging this problem as an unsupervised clustering problem.

We will be extracting features that holistically capture the musical and lyrical properties of songs and use these abstract features to cluster given songs into a user-specified number of playlists based on their intrinsic similarity. We hypothesize that this will deliver better results that are free of biases in the training data that other supervised methods use.

We will be using several public APIs (namely Spotify, musiXmatch and metrolyrics) to assemble a data set that contains the metadata, lyrical and musical properties of songs. Which will then be clustered into playlists using various clustering algorithms like k-means[4], fuzzy k-means[5], hierarchical[6] and mixture of Gaussian clustering algorithms[7]. The validity and performance of these methods will be evaluated using internal measures for cluster purity, also since the quality of a playlist is a largely subjective matter, some amount of human assessment will be required to tune these models.

# Execution and Results

## 0.1   Data Collection and Feature Extraction

Audio features documented in `Table 1` were extracted from spotify for all the songs.Additionally, it was imperative that we numerically capture the representation of the lyrics of the song for a more effective clustering performance. Word vectors trained on the entire text of all the lyrics for this purpose using the `GloVe`[9] algorithm, GloVe is an unsupervised learning algorithm for obtaining vector representations for words. Training is performed on aggregated global word-word co-occurrence statistics from a corpus, and the resulting representations showcase interesting linear substructures of the word vector space.

To capture the sentiment expressed in the lyrics of the song as a numerical feature, we use a lexicon and rule-based sentiment analysis framework called `VADER`[8] (Valence Aware Dictionary and sEntiment Reasoner), a lexicon and rule-based sentiment analysis tool that is specifically attuned to sentiments expressed in popular media. We obtain the amount of positive, negative and neutral sentiment along with a compound score indicating overall sentiment polarity.

These features together gave a holistic and a comprehensive numerical representation of each song and were then used as parameters in the clustering algorithm. Songs that were grouped in the same cluster were the ones that were the closest to one another in the high dimensional vector space defined by the parameters that we extracted. To make sure that all the features were given equal importance, all the parameters were scaled using the min-max scaling method.

## 0.2   Exploratory Data Analysis

To explore relationships and trends in the numerical data collected, a correlation matrix was plotted. We observed that features like danceability, energy and loudness were highly correlated with each other. The correlations observed were in sync with our intuitive understanding of the features and reinforced our hypothesis that human understanding of musical properties were comprehensively represented in the numerical features we extracted. We also plotted trends in

all of these features on both the song and the artist level. Please refer to the jupyter notebooks for EDA for all the plots and corresponding analysis.



Figure 1: Correlation heatmap b/w audio features

## 0.3 Playlist generation

We used clustering algorithms to cluster similar songs based on all the numerical features and the sentiment of the songs to generate a user defined number of playlists from the songs selected for analysis. As hypothesised, songs similar to each other were automatically grouped together in a playlist. `Kmeans` and `Agglomerative` (a form of hierarchical clustering) algorithms were implemented to achieve this. Both the algorithms were successful in creating playlists that capture the essence of the songs.

## 0.4 Evaluation

For evaluation of the clusters generated we performed Human assessment on the playlists. These playlists were compared to randomly generated playlists, evaluators were asked to pick

the playlist that they think makes more sense based on song similarity. 17 songs were selected and clustered into 3 playlists.
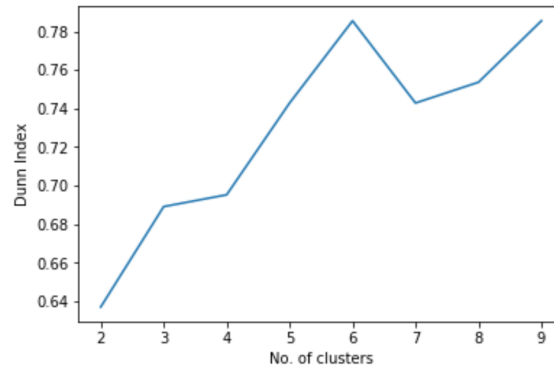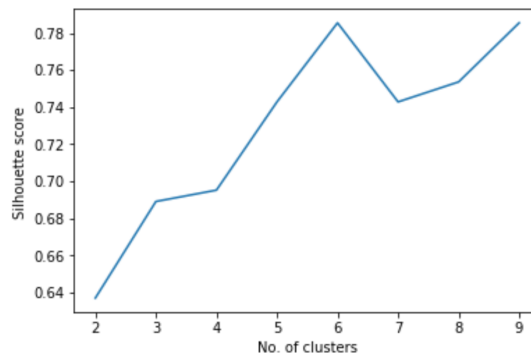


Figure 2: Trend of Dunn Index



Figure 3: Trend of Silhouette Scores

Of the three playlists generated the playlists generated by the algorithm fair better in comparison to randomly generated playlist in all three cases as seen in `Figure 1`. We found that around 80% of the audience found that the playlist generated by the framework we developed were better than random clustering of songs.



Figure 4: Human Assessment result — Playlist 1

In addition to Human assessment we used within cluster validation parameters to judge the

4

Figure 5: Human Assessment result — Playlist 2



Figure 6: Human Assessment result — Playlist 3

clusters created. `Dunn index` and `Silhouette` score were calculated for all approaches to gauge the relative goodness of the clustering algorithm for all the approaches and the trend was captured with the number of playlists to be generated.

The number clusters can be defined based on analyzing the plot for Dunn Index(`Figure 2`)/Silhouette Score(`Figure 3`) vs The number of clusters as shown in figures . We select the cluster size for which the metrics are maximized. The user can also input the desired number of playlists to be generated.

For more personalised generation of playlists, the user can define the attributes of the song that should be given more weightage. These parameters can now be used to generate clusters using weighted Kmeans algorithm.

# References

[1] Dunn, Joseph C. "Well-separated clusters and optimal fuzzy partitions." Journal of cybernetics 4.1 (1974): 95-104

[2] Gong, Xingting, and Xu Chen. "Automatic Playlist Generation." (1989).

[3] Platt, John C. "Auto playlist generation with multiple seed songs." U.S. Patent No. 6,987,221. 17 Jan. 2006.

[4] Hartigan, John A., and Manchek A. Wong. "Algorithm AS 136: A k-means clustering algorithm." Journal of the Royal Statistical Society. Series C (Applied Statistics) 28.1 (1979): 100-108.

[5] Li, Mark Junjie, et al. "Agglomerative fuzzy k-means clustering algorithm with selection of number of clusters." IEEE transactions on knowledge and data engineering 20.11 (2008): 1519-1534.

[6] Johnson, Stephen C. "Hierarchical clustering schemes." Psychometrika 32.3 (1967): 241-254.

[7] Banfield, Jeffrey D., and Adrian E. Raftery. "Model-based Gaussian and non-Gaussian clustering." Biometrics (1993): 803-821.

[8] Hutto, C.J. & Gilbert, E.E. (2014). VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. Eighth International Conference on Weblogs and Social Media (ICWSM-14). Ann Arbor, MI, June 2014.

[9] Pennington, Jeffrey, Richard Socher, and Christopher Manning. "Glove: Global vectors for word representation." Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP). 2014.