

Automated Playlist Generation

Kade Keith, *Student, Stanford University - Computer Science Department*,
Demetrios Fassois, *Student, Stanford University - Computer Science Department*

Abstract—Our project generates music playlists based on a song or set of seed songs, using diverse features ranging from lyrical sentiment to song popularity. We approach the problem as both a graph problem and as a classification problem, and evaluate our results based on real human-curated playlists.

Index Terms—Playlist Generation, Song Recommendation, Sentiment Analysis

1 MOTIVATION

WITH the growth of music streaming services, there are now more songs than ever at music listeners' fingertips. Because of this growth, the art of constructing playlists has become increasingly challenging, and discovering new music in the expanse of choices is a daunting task. For this reason, we seek to build an automatic playlist generator, that can take a few songs as a seed set, and generate a complete playlist for the listener.

1.1 Goal

Using popular, human-curated playlists as our training data and test data, our system should construct playlists of similar quality. A novel aspect of our project is that we incorporate lyrical analysis in our model, as we believe that lyrical content plays an important role when creating playlists.

2 METHOD

We combine data from a number of sources in our project. The primary source is the Million Song Dataset (MSD) [2], and the corresponding lyrics dataset, which provides lyrics for roughly a quarter of those songs in a bag-of-words format. In addition to that we use Spotify [5] as the source of our playlist data, as well as using their's and Last.fm's [1] song info to augment the data from the MSD.

In total, that gives us the following group of attributes:

Feature	Source
Year	MSD, Spotify
Tempo	MSD, Spotify
Timbre	MSD
Tags	Last.fm
Danceability	MSD, Spotify
Energy	MSD, Spotify
Loudness	Spotify
Popularity	Spotify
Speechiness	Spotify
Acousticness	Spotify
Instrumentalness	Spotify
Liveness	Spotify
Valence	Spotify

With these features, we approach the task two ways; first as a graph problem, and second as a classification problem.

2.1 Graph Problem

The first approach is to think of songs as nodes in a graph. With this you can apply k-nearest neighbors to find most similar songs given a seed or set of seeds. The cluster of songs geometrically close together form a playlist.

2.2 Classification

The second is to think of deciding whether or not a song belongs on a playlist as a classification problem. Positive training examples are a subset of songs on the playlist. Negative examples are a random selection of songs not on the playlist. Then, presented with a previously unseen song, the model classifies it as either belonging on the playlist or not.

3 EVALUATION METHOD

TODO - HOW WE SCORE RESULTS. CLASSIFICATION IS EASY - DOES SYSTEM GET IT RIGHT OR NOT. TALK ABOUT CROSS VALIDATION.

Given a subset of a playlist as a seed, evaluation is based on the comparison of our results with the actual remainder of the playlist.

4 PRELIMINARY EXPERIMENTS

4.1 Data Gathering

The bulk of our work thus far has been in data collection and processing. We have set up the complete pipeline for our model using a small subset of the MSD. We first take the MSD and filter out the songs for which we do not have lyrics data. We also remove the fields that we are not considering, such as 'sections', and the fields for which the data is too sparse to be useful, such as artist location (latitude and longitude). Doing this filtering and elimination reduces the size on disk by almost a factor of 10, which is very important considering the whole dataset is 280 GB.

For the remaining songs, we perform preliminary sentiment analysis. As a baseline we are just using Naive Bayes

(with the NLTK movie review corpus as training data [3]) to score each song as either positive or negative, which then gets included in the features. After that, we enhance our data with audio features from Spotify and Last.fm.

For gathering playlists, we rely on searching for a particular term, such as “summer” or “love”, and saving the top playlists for that search. We have automated this process so that all we need to input is the term. Deciding which terms to search for provides an interesting sub-problem. In general, we want to avoid “trending” or “hits” playlists, because even though those are popular, they are not cohesive. We focus instead on themes that we believe users will make cohesive playlists about, such as emotions or seasons.

Unfortunately the overlap of these two sets (the MSD with lyrics and the sampling of popular playlists) is not as high as we would like. Based on preliminary estimates, only 3-5% of the MSD songs with lyrics are showing up in our popular playlists. We attribute this largely to the fact that Spotify surfaces modern songs over older ones, and the MSD is a few years old. Since that would only give us a dataset on the order of a thousand, more gathering of lyrics data will be necessary. We plan to write scripts to automate this process. In addition, we will likely need to dig a little deeper to find more playlists that contain less modern songs.

4.2 Method: Graph Problem

The simplest graph approach is k nearest neighbors, which we have implemented. We represent each song as a point in 11-dimensional space according to our normalized features (all those listed above, excluding Timbre and Tags). Then we select the next songs for that playlist based on proximity to the seed. This is done repeatedly to construct the whole playlist.

4.3 Method: Classification

As a classification baseline, we rely on SciKit Learn’s [4] linear regression implementation. For now we have hand-picked a subset of the MSD and fed that in as our playlist.

5 RESULTS

5.1 Graph Problem

TODO - HOW'D IT DO? HAND-PICKED EXAMPLES?

5.2 Classification

TODO - HOW'D IT DO? HAND-PICKED EXAMPLES?

6 NEXT STEPS

Besides what has already been mentioned (gathering more lyrics/playlists) we have a number of tasks remaining to enhance our model. First is utilizing Timbre information in a meaningful way (TODO - DESCRIBE IN MORE DETAIL). Second is more advanced NLP methods on the song lyrics to provide more meaningful features than just sentiment, and third is expanding the two models we are using to tackle the problem.

6.1 Graph Problem

For the graph version problem, a particularly interesting way to think about a playlists is as a path through a graph [6]. With this approach, you can pick a start song and an end song, and let the algorithm find the shortest path of some predetermined length between them. We believe this and similar approaches will create compelling playlists. In particular we think these playlists will be dynamic and able to “tell a story” more so than playlists that are just clusters of songs.

6.2 Classification

The initial effort for improving the classification version of the task will be in improving the features as previously mentioned. Beyond that, the obvious improvement is to use a non-linear regression model since the assumption our data is linearly separable is problematic.

7 CONTRIBUTIONS

TODO:

KADE - REMOVING SONGS WE DON'T HAVE LYRICS FOR AND FIELDS WE DON'T CARE ABOUT, BASELINE SENTIMENT ANALYSIS, KNN, AND REGRESSION, WROTE REPORT

DEMETRIOS - AUGMENTING MSD DATA WITH SPOTIFY DATA, CROSS VALIDATION

REFERENCES

- [1] Last.fm. <https://www.last.fm/>.
- [2] Million song dataset. <https://labrosa.ee.columbia.edu/millionsong/>.
- [3] Nltk text corpora. <http://www.nltk.org/book/ch02.html>.
- [4] Scikit learn. <http://scikit-learn.org/stable/>.
- [5] Spotify. <https://www.spotify.com/>.
- [6] Masoud Alghoniemy and Ahmed H. Tewfik. A network flow model for playlist generation. In *In Proc IEEE Intl Conf Multimedia and Expo*, 2001.