

Report for Spotify (Capstone 1)

8/4/2022
Ashley Mersman

Background and Objective

Spotify has increased their annual revenue by 18-20% since 2018. They have revolutionized the use of data in the music industry and music streaming. With streaming services as the top way people listen to music, it's important that customers are hearing the tracks they want for customer retention and attraction. What trends are we seeing in track popularity and what features are most important in deciding what tracks to push to listeners? Can we predict song popularity to drive revenue and retention?

Data

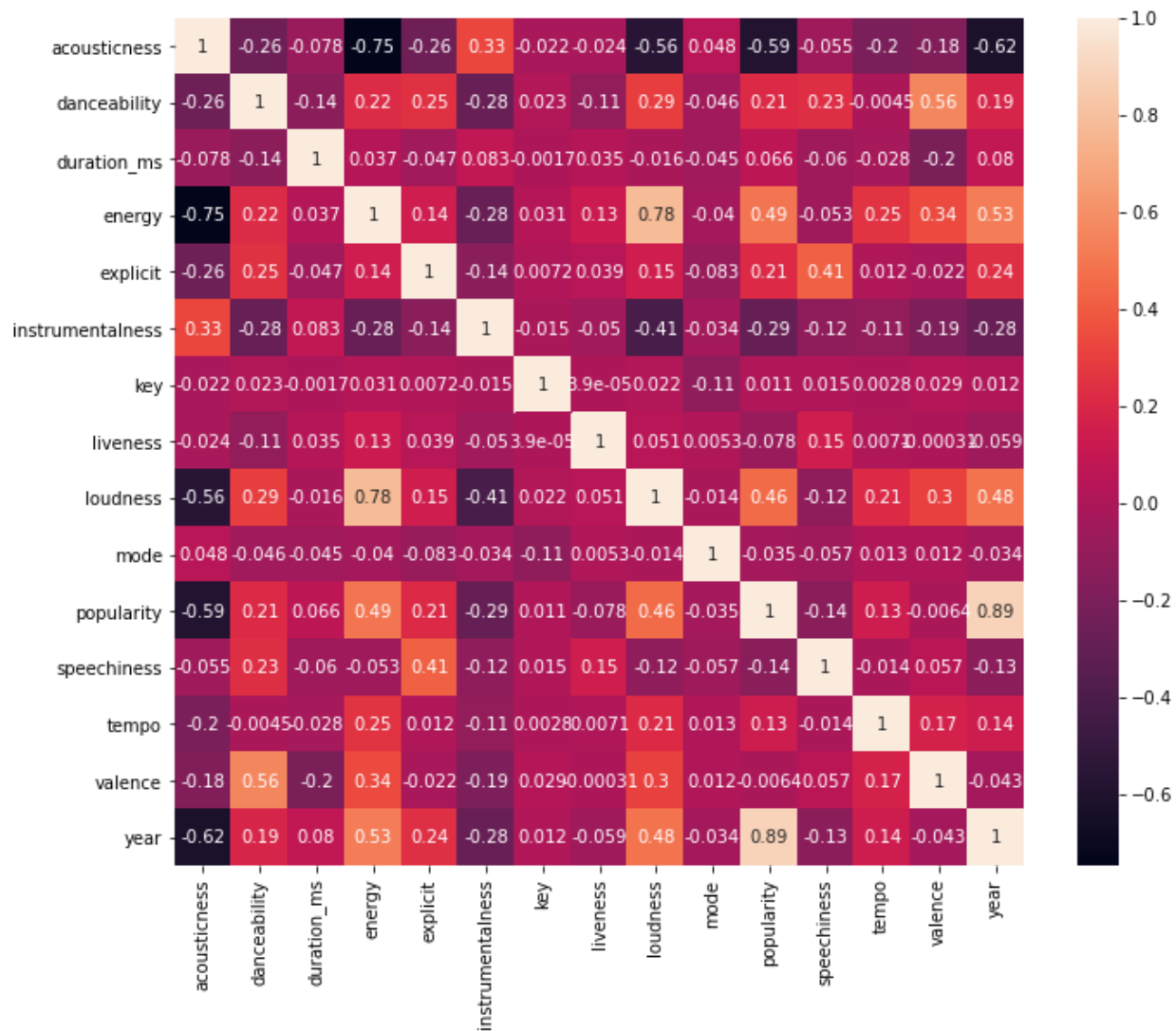
The data used for analysis and to train and test the model were found on Kaggle and included Spotify data from tracks released 1921- 2020. The data set was relatively clean with no missing values. There were some duplications from artists re-releasing the same track on different albums or as single versions. These were filtered out and columns that were not relevant to the analysis, such as Spotify ID, were dropped.

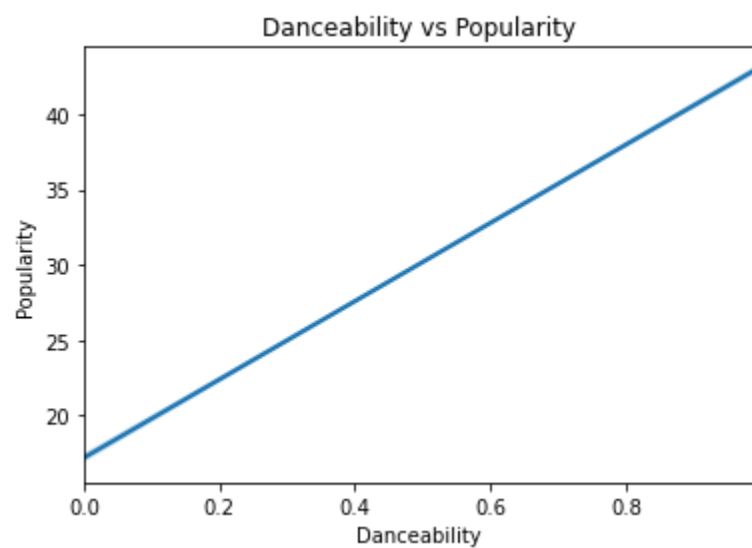
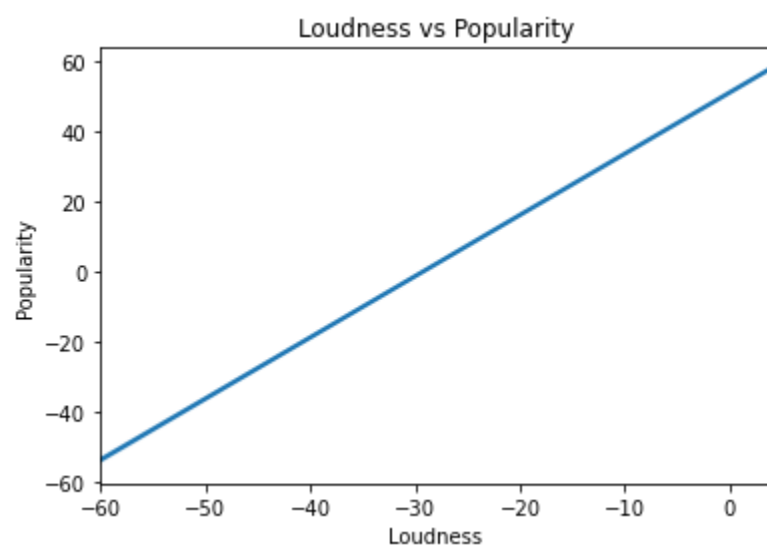
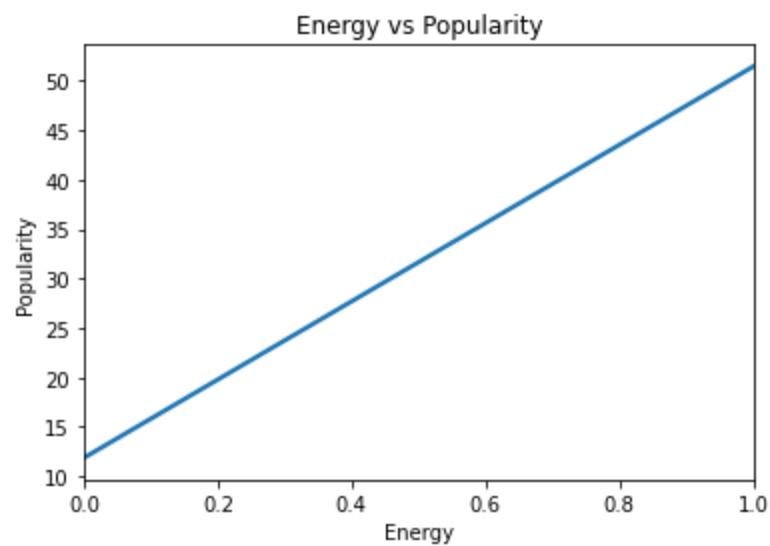
Spotify uses several features to quantify music; acousticness, danceability, duration, energy, instrumentalness, liveness, loudness, speechiness, tempo, valence, and popularity.

Exploratory Data Analysis

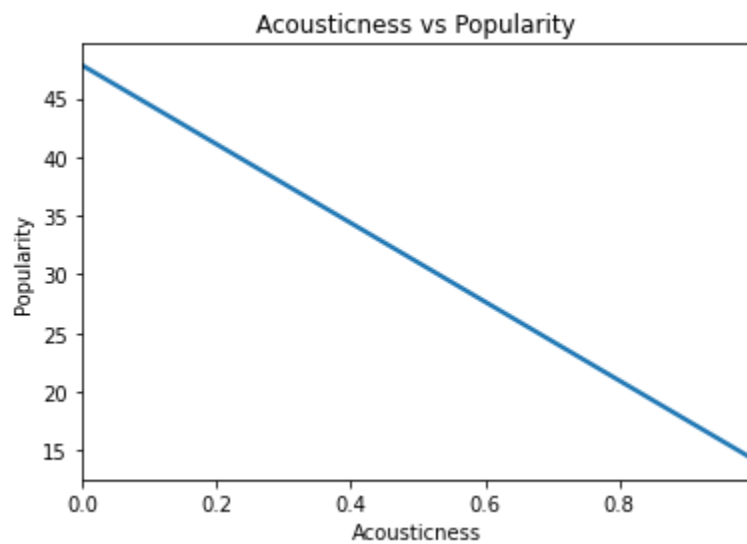
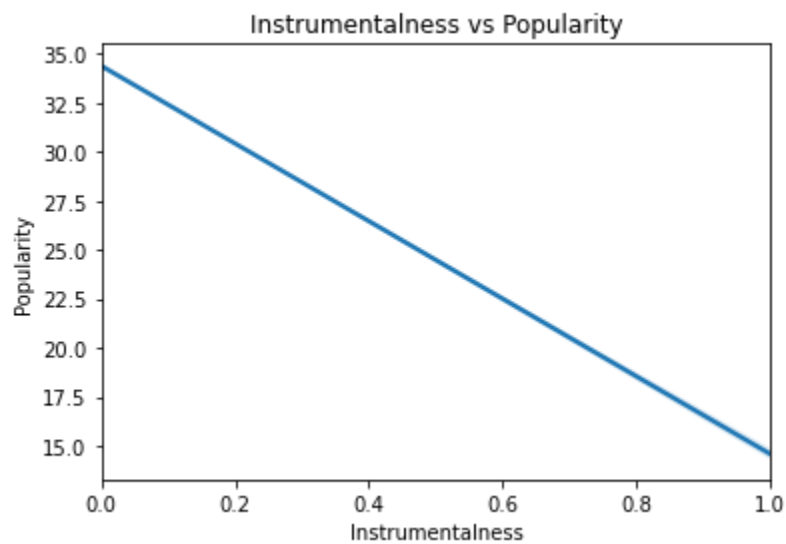
{Look at Popularity, Popularity vs Year, Popularity vs Energy, Popularity vs Acousticness, Popularity vs Loudness and include charts and analysis}

I used a heatmap to find the features with the highest correlation with 'popularity'. Year was highly correlated, which makes sense due to Spotify's popularity feature prioritizing recent plays. Loudness is relatively high in correlation, as well as energy and acousticness has a negative correlation with a high value.

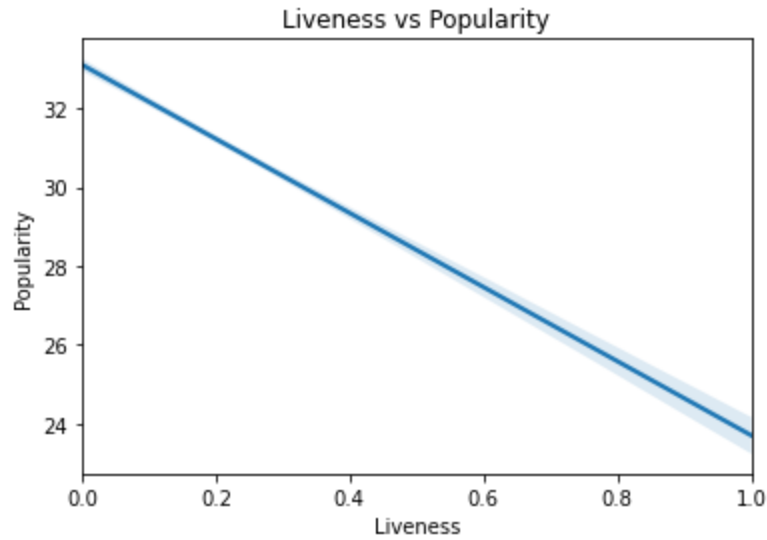




People seem to prefer tracks that are energetic and loud, tracks that have been marked as easy to dance to.

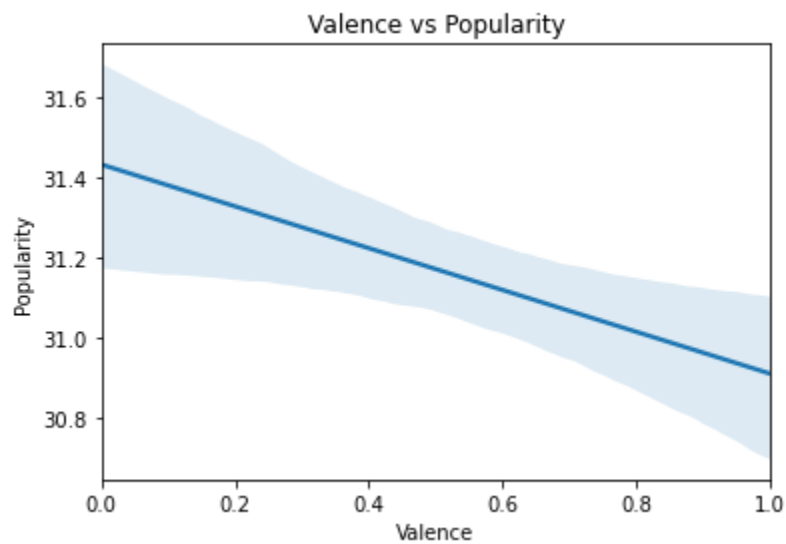


People don't seem to prefer instrumental tracks or acoustic tracks. These have a much more mellow sound so this makes sense when we know that energy and danceability had a positive correlation.



Tracks with a likelihood of being a “live” track fall into these categories and therefore have a negative correlation.

An unexpected correlation was valence. Increased measures of happiness (as defined by Spotify) was inversely related to popularity.



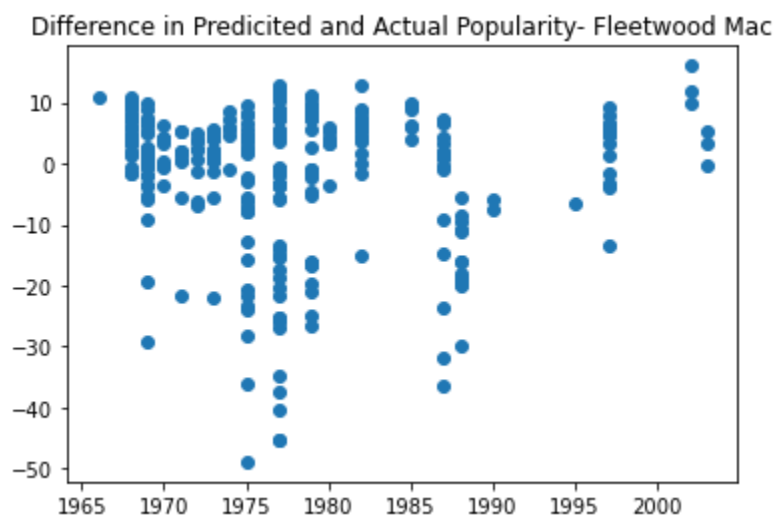
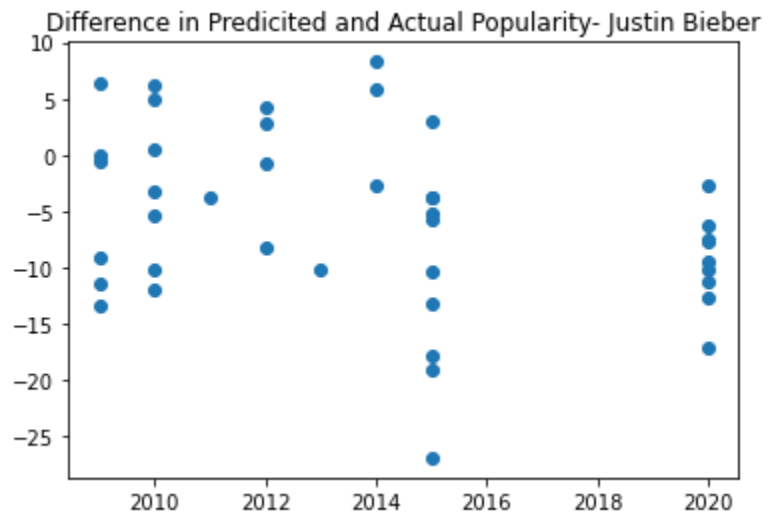
Preprocessing and Model Selection

I started with a 70/30 training/testing split. I then scaled the data using StandardScaler. I tested both a linear regression model and a random forest model with various hyperparameter functions. Hyperparameter grid_search found the best linear model was not using scaled data. The model with the best accuracy was a random forest regressor with n_estimators of 124. However the time it took to fit the model to fit and predict was a huge hurdle for only a 4% improvement over the much faster linear regression.

Model	Best Parameters	CV Score
Linear Regression	StandardScaler: None	0.794
Decision Tree Regressor		0.650
Random Forest Regressor	N_estimators: 195 StandardScaler: Standard Scaler	0.833

Model Analysis

I used the model to predict popularity and check accuracy on both a recent artist, Justin Bieber, and an artist from the 70's, Fleetwood Mac. For each occurrence I fit the model using all the data except that of the artist, did a 5-fold cross-validation and used the average. The difference in release date did not make a difference on the accuracy of the popularity prediction. 0.774 for Justin Bieber vs 0.775 for Fleetwood Mac. I plotted the difference between the actual popularity and the predicted popularity for both artists.



Conclusions

There's some improvements and further research to be done but the model is about 77.5% accurate in predicting popularity on the overall listening audience. This can be used to look at which songs to push to Discover playlists and combined with personalized data can help create playlists, which artists to invest in for Spotify Concerts, which songs Spotify can use in their ad materials to attract listeners, what types of music independent artists can tailor their albums to in order to maximize market share.

The popularity was based on Spotify's own metric which skews toward more recent plays. What is the correlation between this metric of popularity and others, for instance the Billboard top 100 for each year? How have the trends on popular music changed over the past 100 years and can we predict the trends in these attributes for the upcoming decade? These answers could help us invest in new artists.

With the new popularity and boom in podcasts we could also apply these same questions to that media and use that to set investment strategies.