

## more\_exploration

February 20, 2022

Game Plan: 1. Show that the choice of outcome variable does not matter very much because of the extremely high correlation between potential outcomes. - Create a nice looking correlation plot for the outcomes - Probably try to create the `sns.pairplot()` although that function was very very slow 2. Point out that the Spotify playlists are MUCH more successful than the typical user playlist - Can show that by plotting bar charts next to each other, group by whether owner == "spotify" 3. Have to make a decision on what outcome variable to choose (or most likely a combination of outcome variables!) 4. Look at univariate cuts of the outcome variable(s) with the predictors, just to try to give a sense of the relationship between the variables prior to modeling 5. Build a baseline model that we can use for prediction and then we can compare this to a more advanced model - will allow for simpler interpretations 6. Build an ML model that we can then go compare efficacy with the baseline model.

Note: 1. We will likely want to run separate analyses broken out by owner == 'spotify' vs owner != 'spotify' a. See if the importance of the other predictors is comparable for the two groups 2. If a user streams 2 or more songs consecutively, does that count as 2 streams or only 1 since the stream were consecutive?

Possible Options (no particular order) 1. Look into multivariate methods to accommodate the various potential outcome variables (or could just pick one/create a PCA) 2. Should probably think through how to use the 'skippers' data. Could this be at all helpful in dealing with potential biases in how Spotify may promote playlists?

Assumptions: 1. This is a random sample of playlists. If this a biased sample, then any generalizations that we make from the data is likely to be meaningfully inaccurate. 2. Spotify treats each non-Spotify playlist equally in terms of promotion. For examples, if the Spotify algorithms were promoting some genres above others at the time this data was collected, then we are unlikely to get a good read on how genre affects listenership.

3. (a) Spotify treats its own playlists differently than the non-Spotify playlists. If this assumption is correct, then it is likely that Spotify playlists are not particularly comparable to non-Spotify playlists.

(b) Spotify treats its own playlists equally with each other. Thus, an analysis with only Spotify playlists should be okay.

4. Each playlist included in the dataset has existed for at least two months. This ensures that the monthly average users in the previous month variable is not biased by how long the playlist has existed. 5. The Spotify algorithms do not amplify small variations in success. If playlist A was slightly more successful than Playlist B two months ago under 'fair' algorithmic treatment, then the algorithms will not amplify playlist A over playlist B, and thus widen the gulf between the success of the two playlists. In other words, there is a fair marketplace for the playlists to compete, where success does not necessarily beget success simply due to the algorithms. 6. For the categorical variables, genre\_1-genre\_3 and mood\_1-mood\_3, when the value is '-' this is not

a missing value, but is instead imparting the information that the given playlist does not easily fit into the predefined genre and mood types.

# 1 Introduction

The data under consideration for these analyses consists of 403,366 distinct playlists, with 314,899 distinct playlist owners. Of the 314,899 unique playlist owners, 261,040 (83%) have exactly one playlist in the data. Of the owners with more than one playlist, Spotify itself has the most, with 399. The data is composed of only playlists from US owners, and thus extrapolating any of these analyses to other countries is likely unwarranted or should be done with great caution. Each playlist is categorized by its top three genres and top three moods. There are 26 genres and 27 moods under consideration.

There are a number of potential measures of playlist success included in this dataset. Specifically, we have (1) The number of streams from the playlist today, (2) the number of streams greater than 30 seconds today, (3) the number of active users today, where an active user is defined as having a stream > 30 seconds, (4) the number of active users in the past week, (5) the number of active users in the past month, (6) the number of users who had a stream this playlist for any length of time in the past month, (7) the number of active users in the previous month, (8) the total number of > 30 second streams in the past month, (9) the number of users who were active this month and the previous month, and (10) the number of > 30 second streams by the playlist owner in the past month. The data also includes the number of users who skipped more than 90% of their total streams today who also used this playlist, which could be used as a reverse encoded outcome variable.

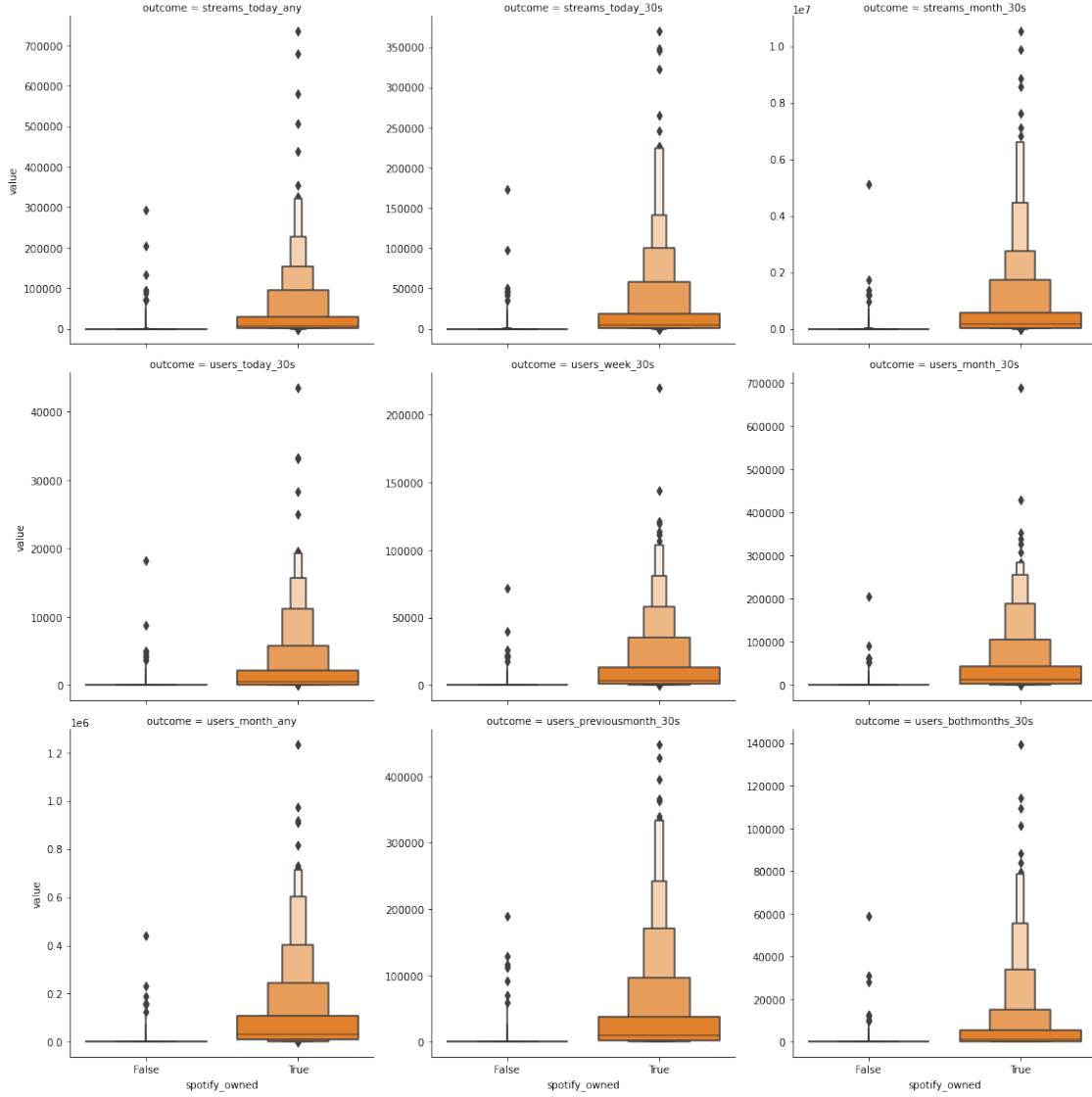
Some of the potential predictors of stream success include: (1) the number of tracks in the playlist, (2) the number of tracks that were added to the playlist today, (3) The number of unique artists in the playlist, (4) the number of unique albums in the playlist, (5-7) the first, second and third most common genre found in the playlist, (8-10) the first, second and third most common mood found in the playlist, and (11) the tokens associated with the playlist.

There are two Spotify-owned playlists that constitute extreme outliers across each of the potential outcome variables. The first is a pop, Dance & House, Indie Rock playlist with 100 tracks and has tokens ‘top’, ‘tracks’, ‘currently’, ‘spotify’. The second is a pop, R&B, Dance & House with 51 tracks and has tokens ‘top’, and ‘hits’. These two playlists have more than three times as many streams today as their nearest competitor and more than four times as many > 30 second streams in the past month as the nearest competitor. For the purposes of plotting, these playlists will be removed.

## 1.1 Comparing Spotify-Owned and Non-Spotify-Owned Playlists

Even after removing the two most successful Spotify-owned playlists, there is still a wide gulf between the Spotify-owned and non-Spotify-owned playlists. To illustrate this, consider the boxen plot below, which shows the distribution of each of the potential outcome variables stratified by whether the playlist is Spotify-owned. We observe that the Spotify-owned playlists are much more successful than the vast majority of non-Spotify-owned playlists. However, there are a few user-created playlists that can rival an average Spotify-created playlist.

```
<seaborn.axisgrid.FacetGrid at 0x1ab0fb1cb80>
```



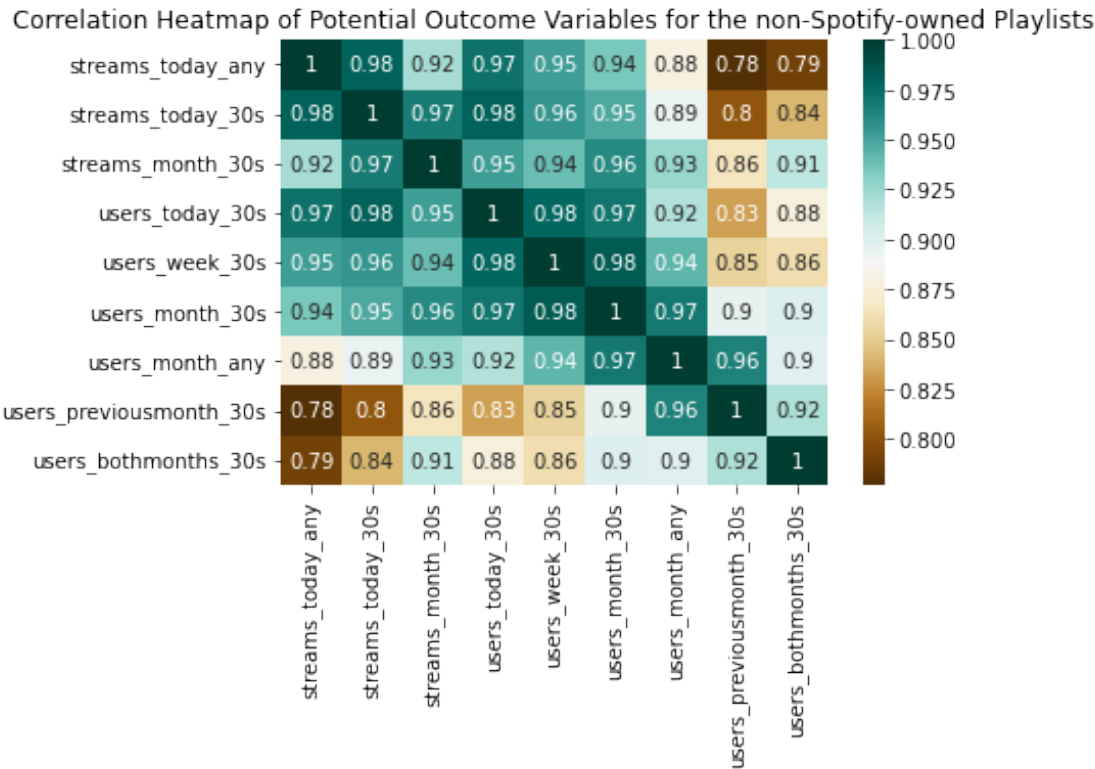
## 1.2 Exploration of Potential Outcomes

Because there is such an enormous difference between Spotify-created and user-created playlists, we will continue our data exploration stratifying by whether the playlist is Spotify-created. In results not shown, the correlation between the number of > 30 second streams by the playlist owner is **much** more weakly correlated with the other potential outcome variables than, and thus will not be included in the following analysis.

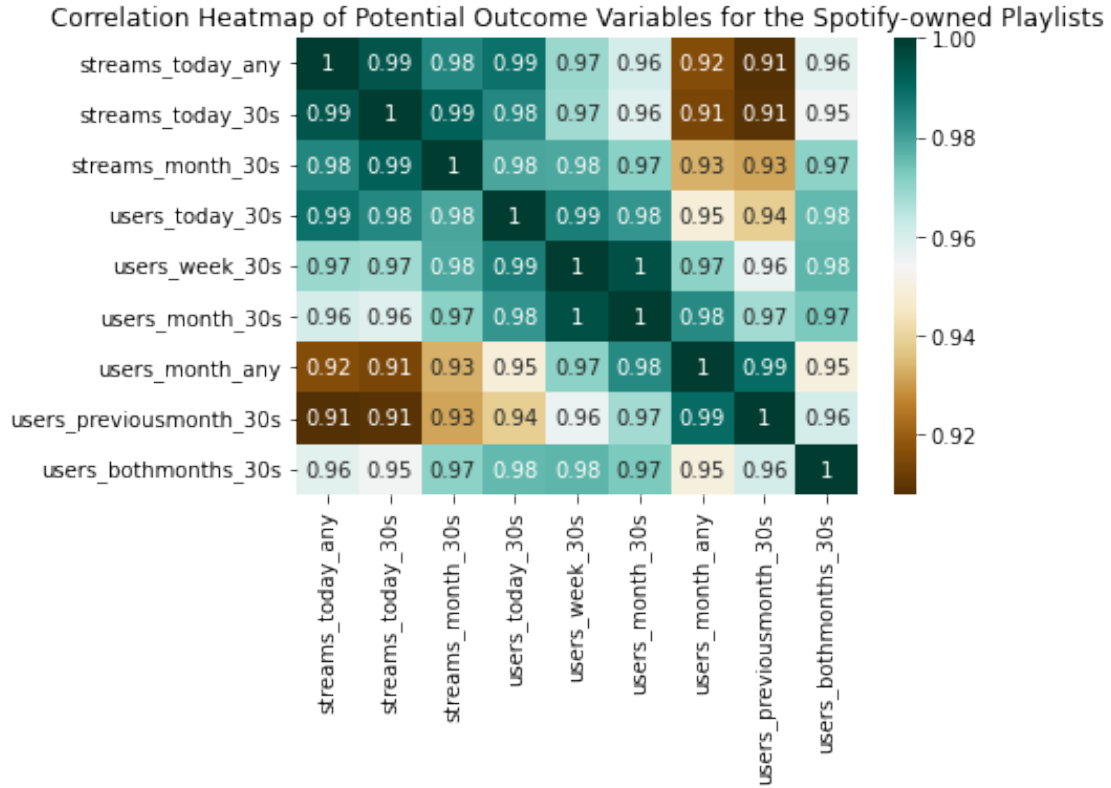
To begin understanding the relationship between the potential outcomes of interest, we present a heatmap of the Pearson correlation between each of the outcomes. We create one such heatmap for the non-Spotify-owned playlists and another for the Spotify-owned playlists. Amongst the non-Spotify-owned playlists, we see that the minimum correlation between any of the outcomes is 0.78, thus signifying a great deal of correlation between our potential outcomes. While all potential outcomes are highly correlated, the mostly weakly correlated outcomes are the outcomes

related to playlists' longer-term success (i.e. the monthly average users in the given month and the previous month) with the more recent measures of success (i.e. the number of total and >30 second streams today and the number of active users today and in the past week) Amongst the Spotify-owned playlists, the potential outcomes are even more highly correlated, with the smallest correlation being 0.91. The Spotify-owned playlists exhibit a similar general pattern to those of the non-Spotify-owned in that the weakest correlation between the outcomes is between the more long-term measures of success and the measures of success in the more recent past. However, the degree of correlation is still immense between monthly active users over the past two months and the number of streams that occurred today, thus indicating that Spotify playlists tend to have considerable 'staying power'. Of course, if more successful playlists in the past are algorithmically pushed to users, then this could become a self-fulfilling prophecy rather than a true indication of how 'intrinsically good' the palylist is.

Text(0.5, 1.0, 'Correlation Heatmap of Potential Outcome Variables for the non-Spotify-owned Playlists')



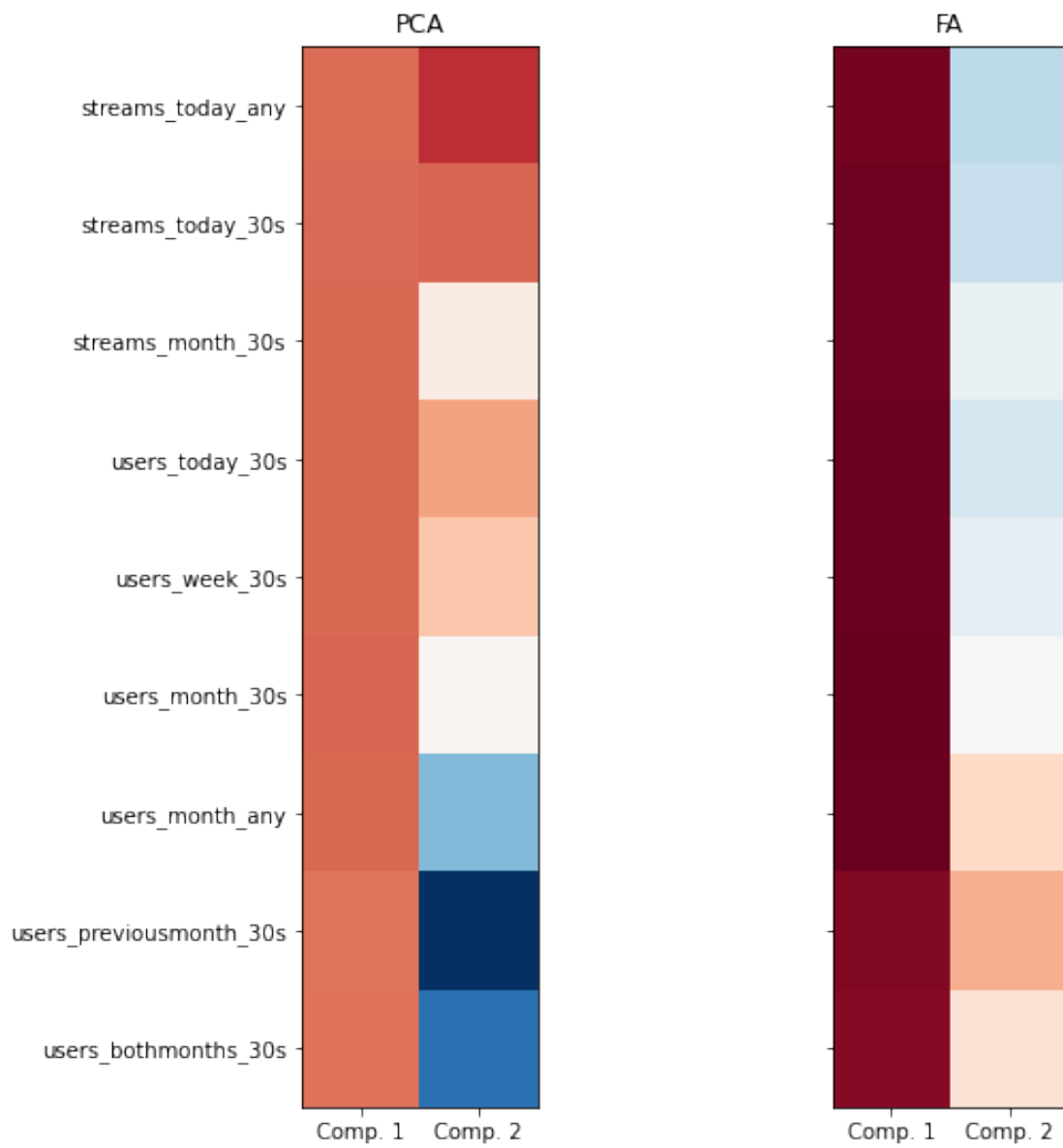
Text(0.5, 1.0, 'Correlation Heatmap of Potential Outcome Variables for the Spotify-owned Playlists')

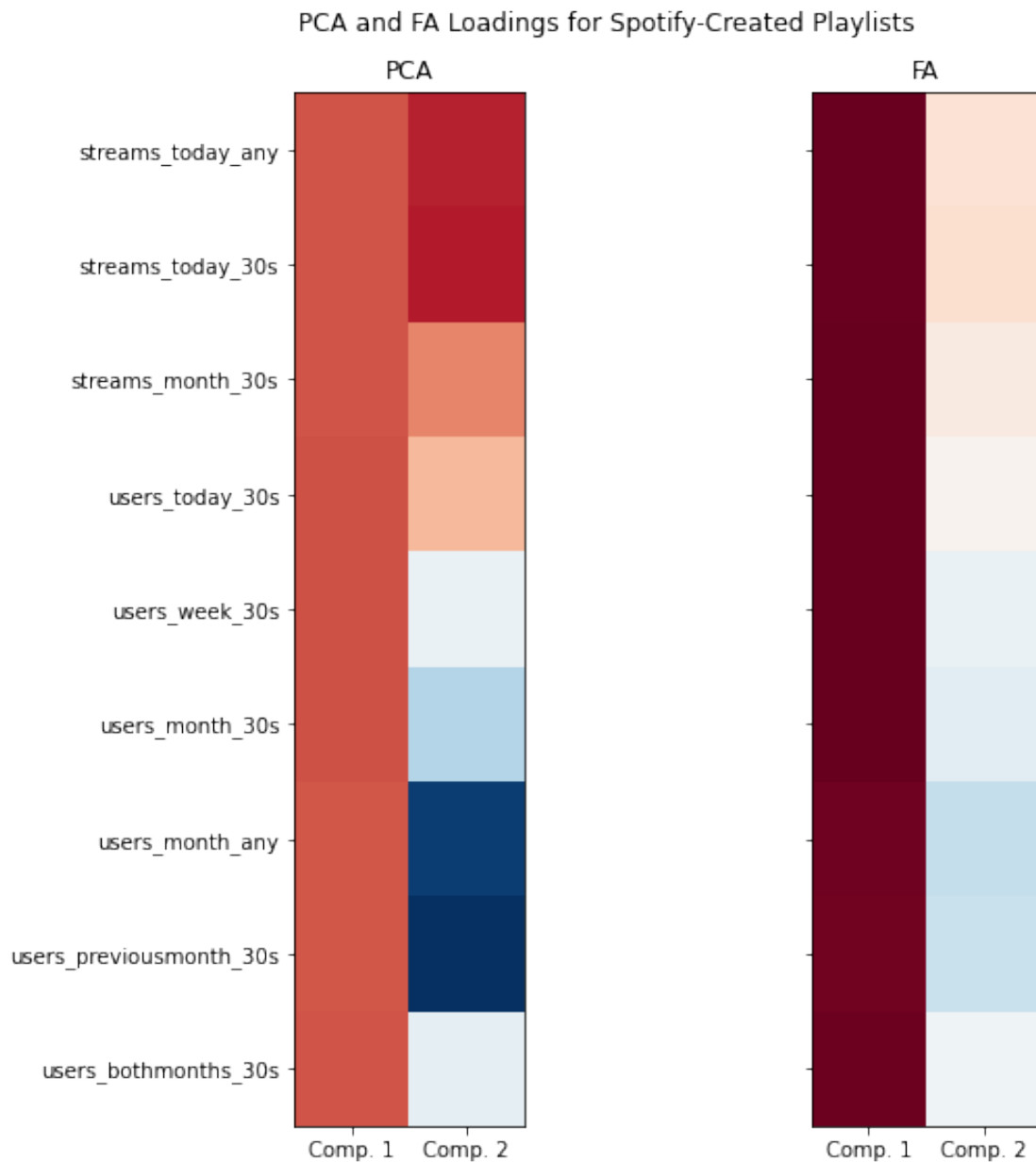


### 1.3 Factor Analysis of Potential Outcomes

We will continue our exploration of the relationship between the potential outcomes by performing a factor analysis.

PCA and FA Loadings for User-Created Playlists



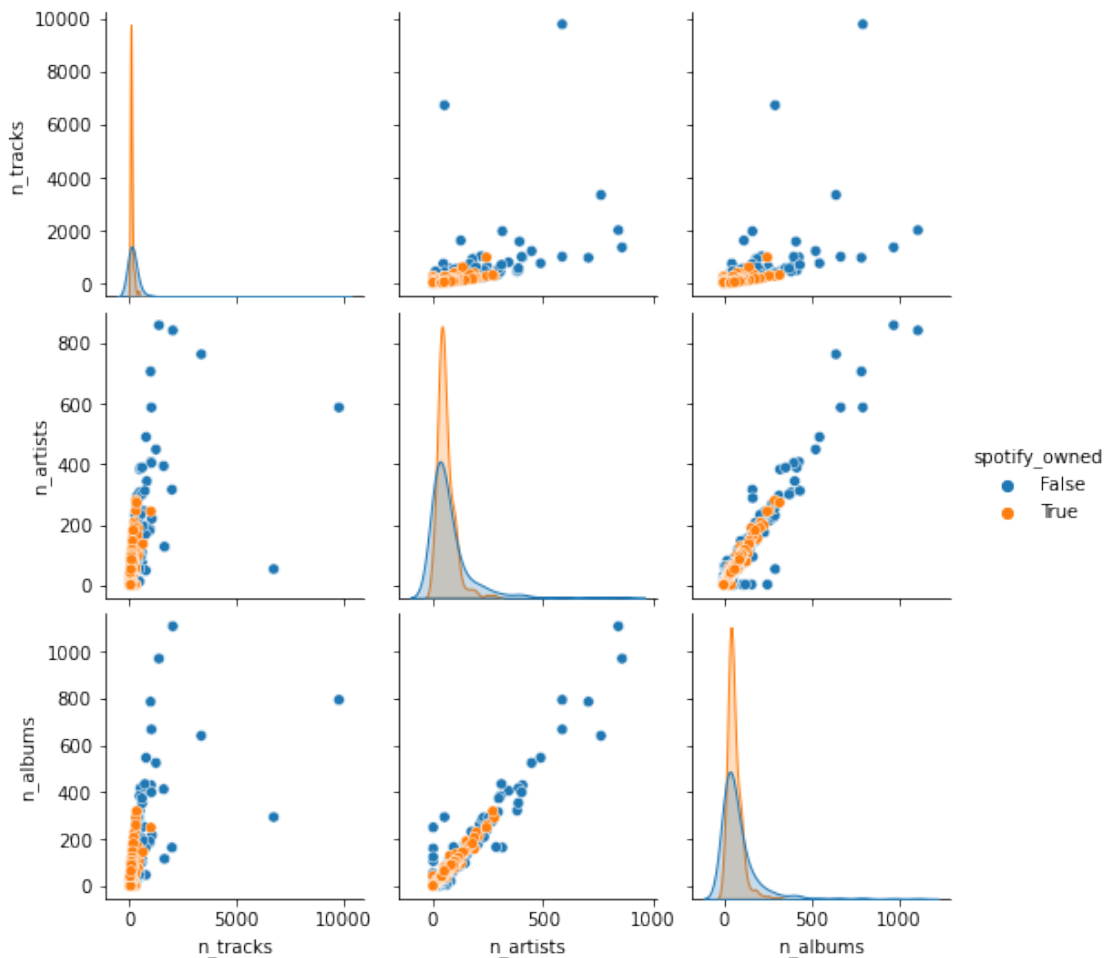


## 1.4 Exploration of Predictors of Interest

Note that for the plots below, we drew a random sample of user-created playlists for these plots in order for them not to visually swamp the Spotify-created playlists. We can see from the plot below that, unsurprisingly, as a playlist's number of tracks increases so do the number of artists and albums. We will thus define tracks per artist and tracks per album as scaled measures of the number of artists and albums that account for the overall size of the playlist. Also, note that the distribution of number of tracks, artists and albums does not vary between Spotify-created and user-created to the degree that the outcome variables do. However, we still observe some difference between user-created and Spotify-created playlists, namely that user-created playlists tend to have

moer tracks, artists and albums, with some user-created playlists having many more of each.

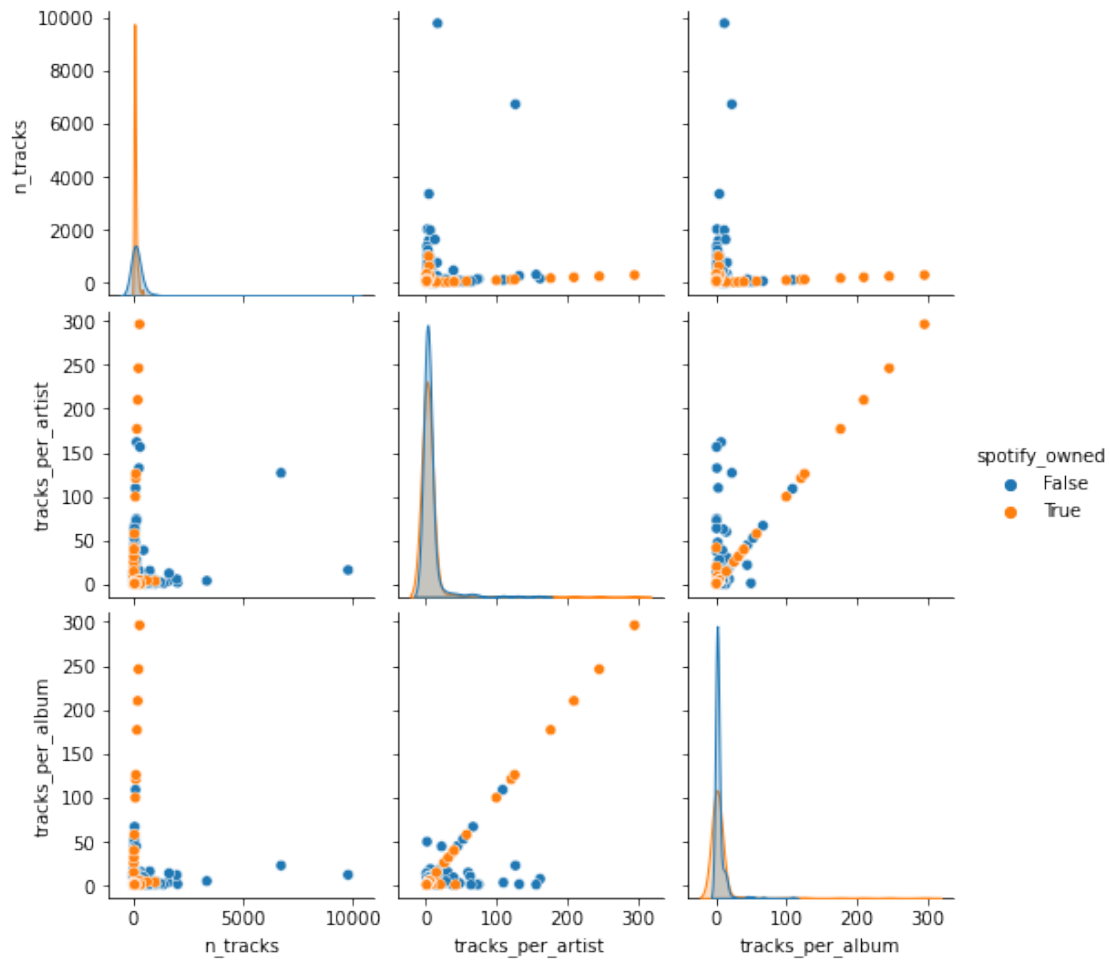
```
<seaborn.axisgrid.PairGrid at 0x1ab10df40d0>
```



Let's consider the correlation of the new scaled measures of artists and albums as well as whether we observe differences between the user-created and Spotify-created playlists. Notice that the distribution of tracks per artist and tracks per album appear quite similar for the user-created and Spotify-created playlists. Spotify-created playlist appear to exhibit a little bit more artist and album diversity than the user-created playlists, but the difference is miniscule compared to the outcome variables.

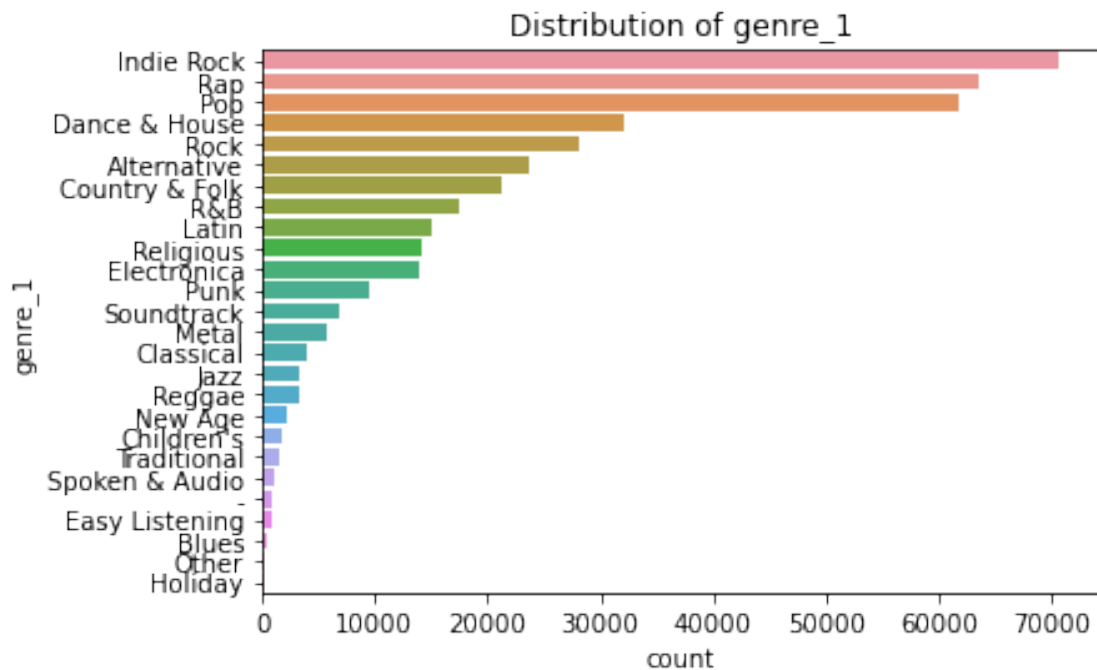
```
<seaborn.axisgrid.PairGrid at 0x1ab2e43e400>
```



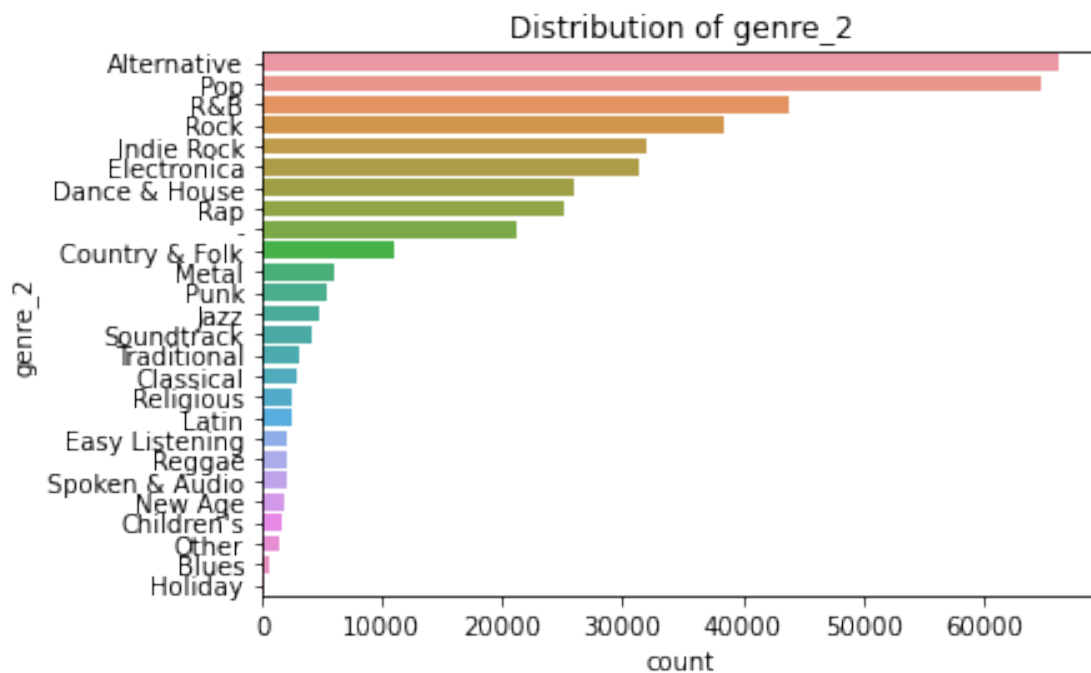


We now turn our attention to describing the categorical predictors: genre and mood.

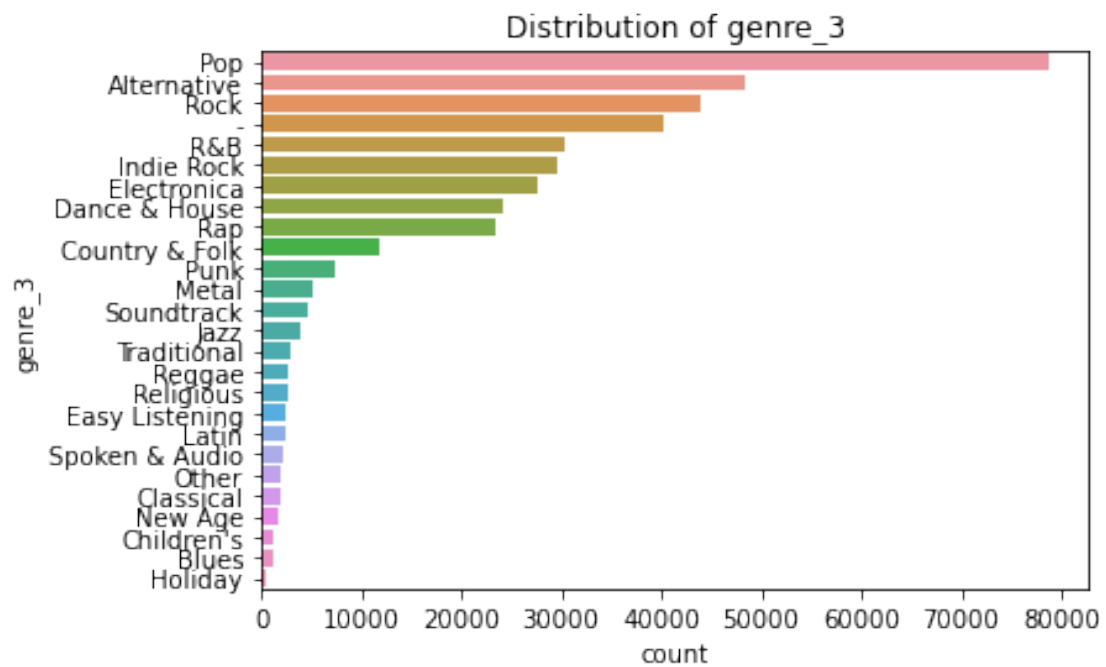
```
<AxesSubplot:title={'center':'Distribution of genre_1'}, xlabel='count',
ylabel='genre_1'>
```



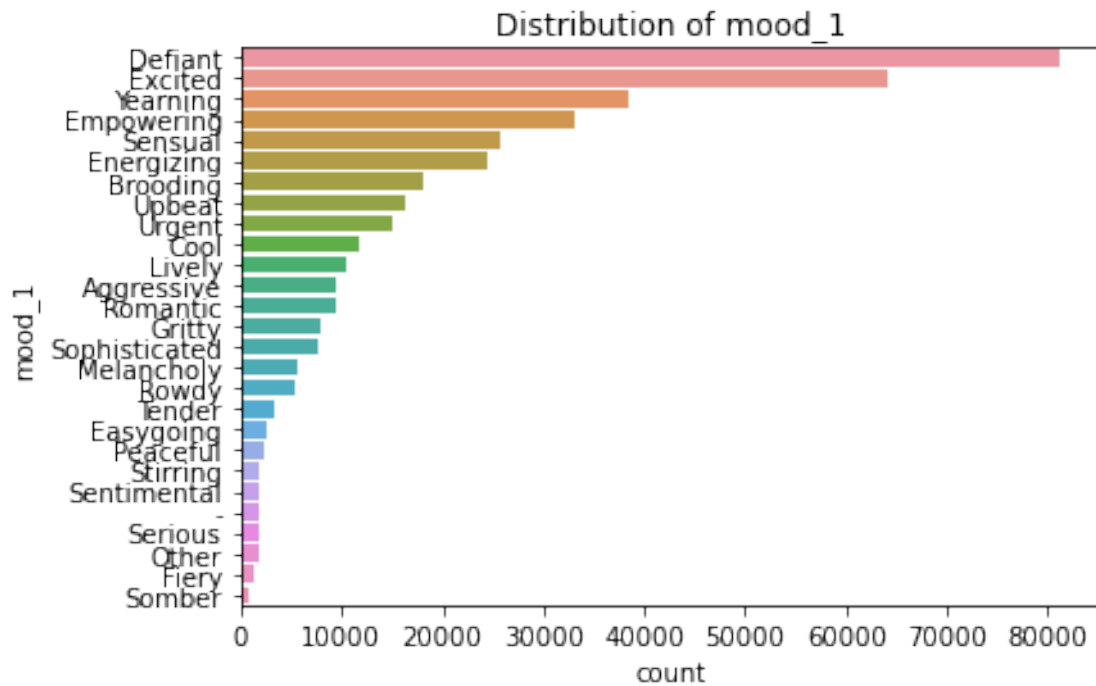
```
<AxesSubplot:title={'center':'Distribution of genre_2'}, xlabel='count',
ylabel='genre_2'>
```



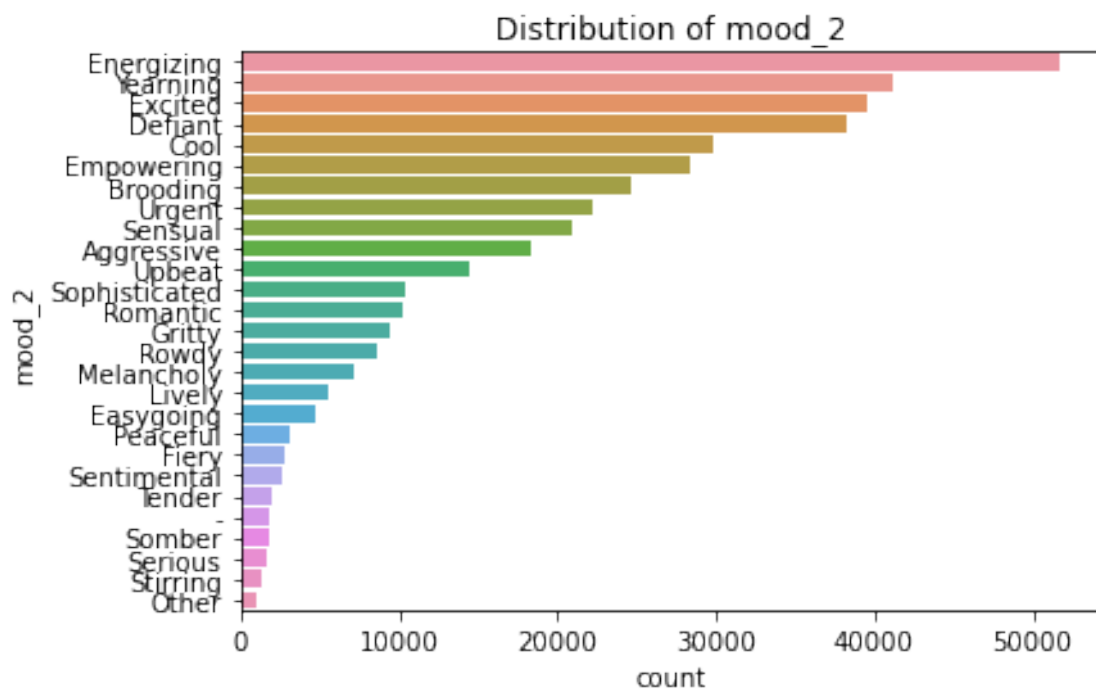
```
<AxesSubplot:title={'center':'Distribution of genre_3'}, xlabel='count',
ylabel='genre_3'>
```



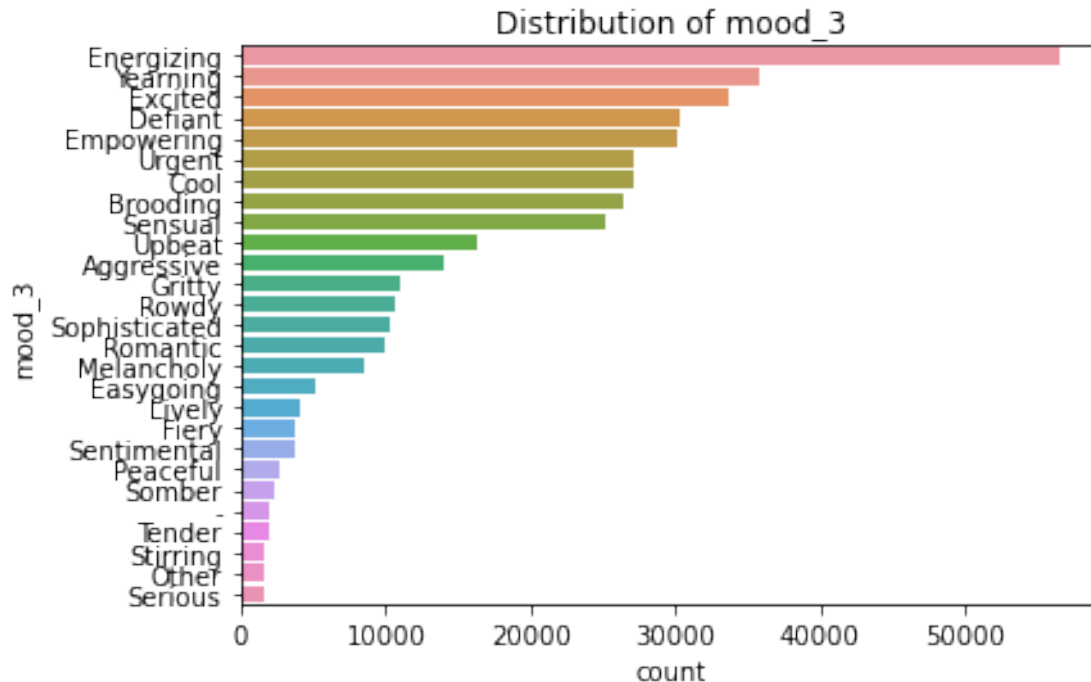
```
<AxesSubplot:title={'center':'Distribution of mood_1'}, xlabel='count',
ylabel='mood_1'>
```



```
<AxesSubplot:title={'center':'Distribution of mood_2'}, xlabel='count',
ylabel='mood_2'>
```



```
<AxesSubplot:title={'center':'Distribution of mood_3'}, xlabel='count',
ylabel='mood_3'>
```

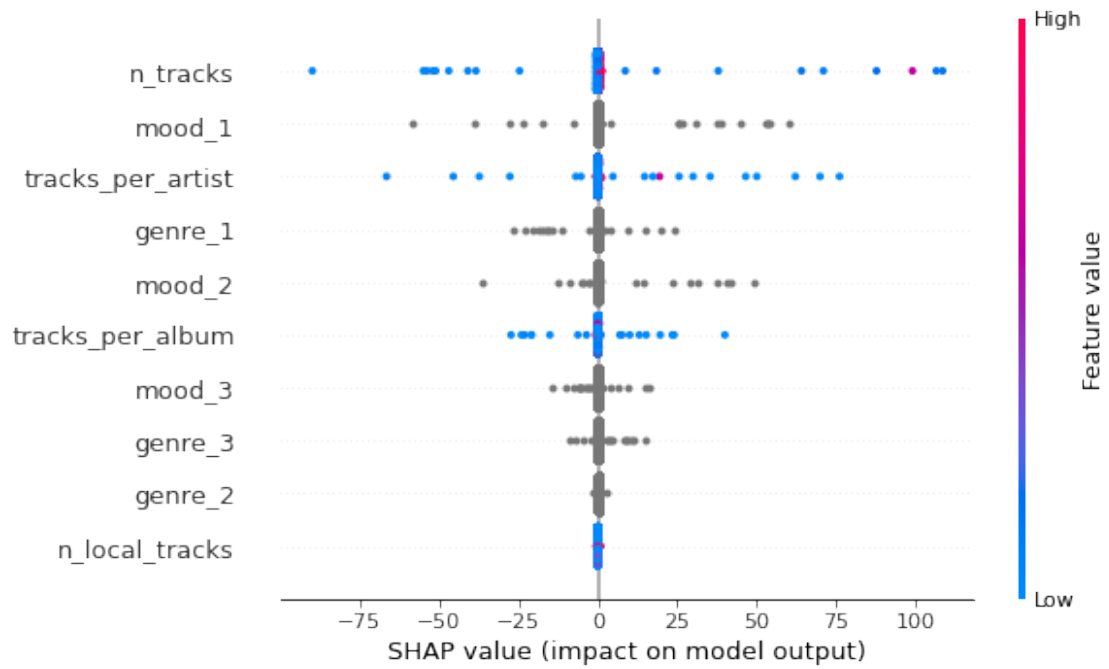


We now need to get a sklearn pipeline built for modeling the various outcomes. We will try to predict the outcomes ‘users\_today\_30s’ and ‘users\_bothmonths\_30s’. We choose ‘users\_today\_30s’ because amongst the user-created playlists, it very highly correlated with streams\_today\_any, streams\_today\_30s, streams\_month\_30s, users\_week\_30s, and users\_month\_30s (with the minimum correlation being 0.95). It is also less correlated with users\_previousmonth\_30s (0.83) and users\_bothmonths\_30s (0.88). We choose ‘users\_bothmonths\_30s’ because (1) it is less correlated than many of the other outcomes and (2) it is a good measure of a playlist’s longer-term staying power.

```
<IPython.core.display.HTML object>
```

```
<shap.plots._force.AdditiveForceVisualizer at 0x1ab3cf995b0>
```

Note that I have deleted the spotify\_owned from these plots because its effect makes all other variables effect hard to see.



11

