

Spotify Playlist Success Case Study Presentation

February 25, 2022

1 Executive Summary

1. Across all available metrics of success, Spotify-created playlists are much more successful than nearly all user-created playlists.
2. All considered measures of success are highly correlated.
3. The number of tracks is the most important predictor of short-term and long-term playlist success, with longer playlists tending to be more successful.
4. For Spotify-created playlists, having greater musical diversity (as measured by tracks per album and tracks per artist) is associated with increased success. This effect is present for the user-created playlists, but is less pronounced.
5. Greater musical diversity is more important for a playlist's long-term staying power than for short-term success.
6. For Spotify-created playlists, the primary genres of 'dance & house' and 'pop' are associated with both short-term and long-term success.
7. The most success-inducing genres for Spotify-created and user-created vary considerably. For user-created playlists, "Children's", "Latin" and "Traditional" increase both short- and long-term success.

```
[1]: %%capture
# 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?

```

```

[2]: import numpy as np
import pandas as pd
import seaborn as sns
import sklearn
import matplotlib.pyplot as plt
import lightgbm as lgb
from sklearn.decomposition import FactorAnalysis, PCA
from sklearn.preprocessing import StandardScaler, OneHotEncoder
from sklearn.linear_model import MultiTaskLassoCV
from sklearn.model_selection import KFold
import shap
import statsmodels.api as sm
import statsmodels.formula.api as smf
from patsy import dmatrices
pd.set_option('display.max_columns', None)
import warnings
import plotly.express as px
warnings.filterwarnings('ignore')
%matplotlib inline

```

```

[3]: outcome_vars_dict = {'streams': 'streams_today_any',
                          'stream30s': 'streams_today_30s',
                          'monthly_stream30s': 'streams_month_30s',
                          'dau': 'users_today_30s',

```

```

        'wau': 'users_week_30s',
        'mau': 'users_month_30s',
        'users': 'users_month_any',
        'mau_previous_month': 'users_previousmonth_30s',
        'mau_both_months': 'users_bothmonths_30s',
    }

outcome_vars = list(outcome_vars_dict.values())

raw_data = pd.read_table('data/playlist_summary_external.txt')
raw_data = raw_data.rename(columns = outcome_vars_dict)
spotify_data = raw_data.query('owner == "spotify"')
nonspotify_data = raw_data.query('owner != "spotify"')

predictor_vars = ['n_tracks',
                  'n_local_tracks',
                  'n_artists',
                  'n_albums',
                  'genre_1',
                  'genre_2',
                  'genre_3',
                  'mood_1',
                  'mood_2',
                  'mood_3'
                  ]

```

2 Outline

1. Explore massive differences in playlist success between the Spotify-created and user-created playlists.
2. Establish the substantial correlation between each of the measures of success across both Spotify-created and user-created playlists will next be discussed. A principal components analysis will be performed.
3. Briefly present the distribution of the predictors of success, including both the numeric (number of tracks, tracks per artist and tracks per album) as well as the categorical (mood and genre).
4. Fit gradient boosting models to predict measures of both short-term and long-term success, stratifying by Spotify-created and user-created playlists. Use SHAP values to assess the inferential implications of the models.
5. The Appendix contains a number of other subanalyses and assumption checks.

3 Brief Data Introduction

- The data consists of 403,366 distinct playlists, with 314,899 distinct playlist creators.
- Of the creators with more than one playlist, Spotify itself has the most, with 399.
- The data is composed of only playlists from US creators, and thus extrapolating any of these analyses to other countries is likely unwarranted or should be done with great caution.
- The data consists of a number of potential metrics of playlist success, including measures of total users and streams today, over the past week, month, and over the past two months.
- Each playlist is categorized by its top three genres and top three moods. There are 26 genres and 27 moods under consideration.

4 Comparing Spotify-Created and User-Created Playlists

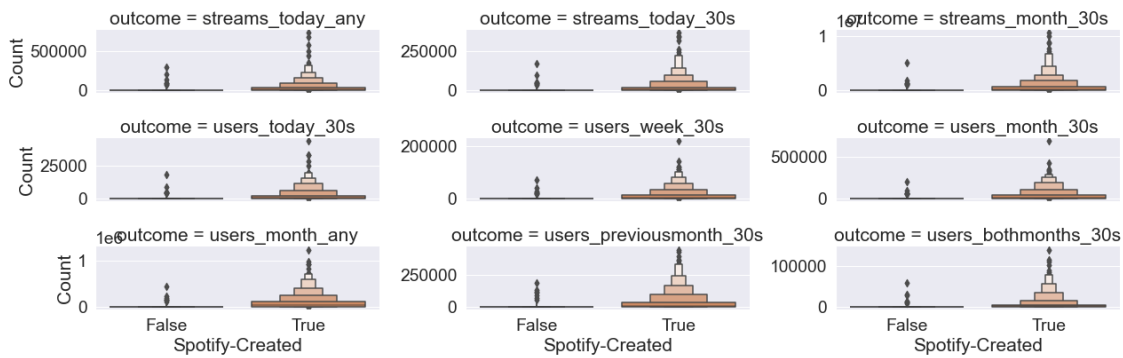
- There are two Spotify-created playlists that constitute extreme outliers across each of the potential outcome variables, which have more than four times as many > 30 second streams in the past month as the nearest competitor. For the purposes of plotting, these playlists will be removed.
- Even after removing the two most successful Spotify-created playlists, there is still a wide gulf between the Spotify-created and user-created playlists.

```
[4]: raw_data['spotify_created'] = raw_data['owner'] == 'spotify'
comparison_data = (raw_data.
                    drop(index = [163726, 152032]).
                    loc[:, ['spotify_created'] + outcome_vars].
                    melt(id_vars = 'spotify_created',
                        var_name = 'outcome',
                        value_name = 'value')
                    )
```

```
[5]: ###capture
sns.set(font_scale = 1.75)
#plot_data['value'] = plot_data['value'] / 1000

p = sns.catplot(data = comparison_data,
                kind = 'boxen',
                col = 'outcome',
                col_wrap = 3,
                sharey = False,
                y = 'value',
                x = 'spotify_created',
                height = 2,
                aspect = 3)

p.set_ylabels('Count')
p.set_xlabels('Spotify-Created');
```

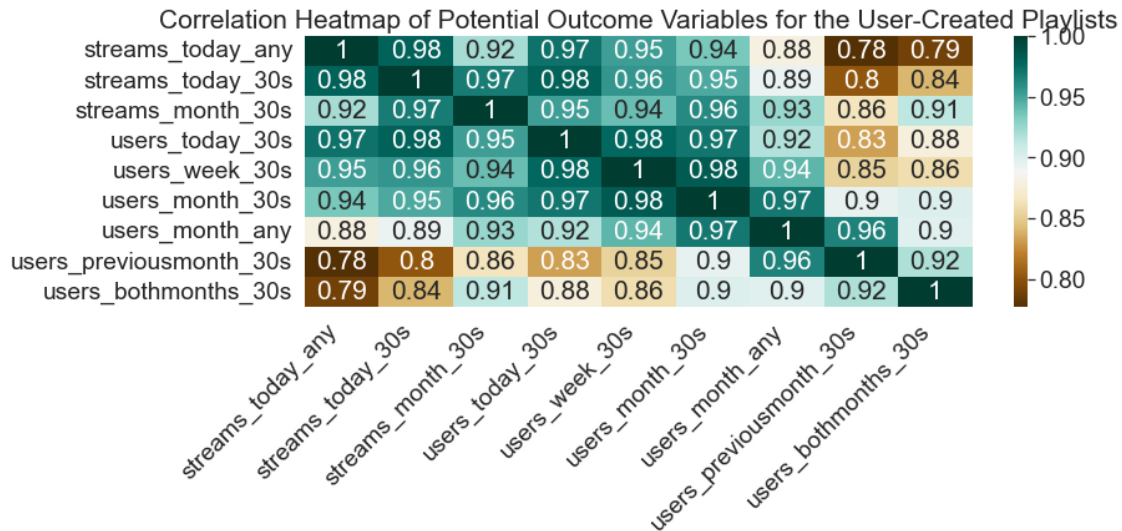


5 Exploration of Potential Outcomes

- Amongst the user-created playlists, we see that the minimum correlation between any of the outcomes is 0.78.
- Weakest correlation is between measures of short- and long-term success.

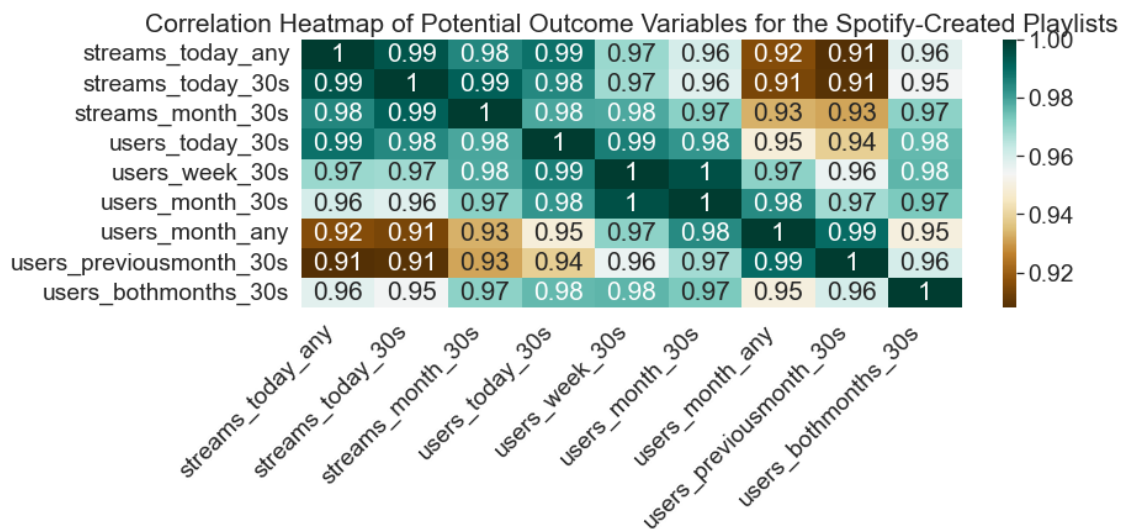
```
[6]: def plot_outcome_corr(data):
    plot_data = data.loc[:, outcome_vars]
    plot_data_corr = plot_data.corr()
    #sns.set(rc = {'figure.figsize': (10, 5)})
    plt.figure(figsize = (12, 4))
    return plot_data, sns.heatmap(plot_data_corr, annot = True, cmap = 'BrBG')
```

```
[7]: #%%capture
nonspotify_outcome_data, nonspotify_heatmap = plot_outcome_corr(nonspotify_data)
nonspotify_heatmap.set_xticklabels(nonspotify_heatmap.get_xticklabels(),
    ↪rotation = 45, horizontalalignment = 'right')
nonspotify_heatmap.set_title('Correlation Heatmap of Potential Outcome_
    ↪Variables for the User-Created Playlists');
```



- Amongst the Spotify-created playlists, the potential outcomes are even more highly correlated, with the smallest correlation being 0.91.
- The Spotify-created playlists exhibit a similar general pattern as the user-created in that the weakest correlation is between the long-term and short-term measures of success.

```
[8]: # NOTE: A FRAGMENT JUST MEANS THAT YOU HAVE TO CLICK AGAIN TO BUILD OUT THE
      ↪SLIDE
spotify_outcome_data, spotify_heatmap = plot_outcome_corr(spotify_data)
spotify_heatmap.set_xticklabels(spotify_heatmap.get_xticklabels(), rotation =
      ↪45, horizontalalignment = 'right')
spotify_heatmap.set_title('Correlation Heatmap of Potential Outcome Variables
      ↪for the Spotify-Created Playlists');
```



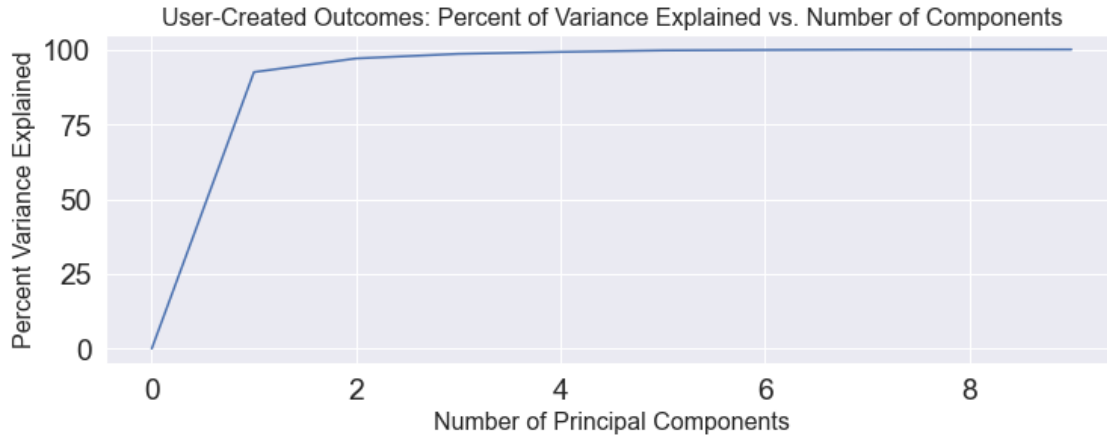
5.1 Principal Components Analysis of Potential Outcomes

- We use a Principal Components Analysis (PCA) to address the degree to which the information in the nine potential outcomes can be reduced to one or a few factors.
- We will show results for the User-Created playlists. The Spotify-created playlist results are similar and can be found in the Appendix.
- 92% of the variation is explained by the first component and an additional 5% is explained by the second component.

```
[9]: nonspotify_outcome_scaled = StandardScaler().  
      ↪fit_transform(nonspotify_outcome_data)  
      spotify_outcome_scaled = StandardScaler().fit_transform(spotify_outcome_data)
```

```
[10]: def plot_percent_variation_explained(data, title):  
        pca = PCA(n_components = data.shape[1])  
        pca.fit(data)  
        percent_explained = 100 * pca.explained_variance_ratio_.cumsum()  
        percent_explained = np.insert(percent_explained, obj = 0, values = 0.0,   
        ↪axis = 0)  
        plt.figure(figsize = (12, 4))  
        plt.plot(percent_explained)  
        plt.title(title,  
                  fontdict = {'fontsize': 16})  
        plt.xlabel('Number of Principal Components',  
                  fontdict = {'fontsize': 16})  
        plt.ylabel('Percent Variance Explained',  
                  fontdict = {'fontsize': 16})  
        return p
```

```
[11]: ###capture  
p = plot_percent_variation_explained(  
    nonspotify_outcome_scaled,  
    title = 'User-Created Outcomes: Percent of Variance Explained vs. Number   
    ↪of Components'  
)
```



- The plot below presents the loadings for each of the first two principal components.
- The first component loads nearly equally across each of the 9 potential outcomes.
- The second component loads negatively on the long-term measures of success and positively on total streams today.

```
[12]: def plot_pca_loadings(data_scaled, data, title, n_components = 2):
    pca = PCA(n_components = n_components)
    pca.fit(data_scaled)
    pca_loadings = pca.components_

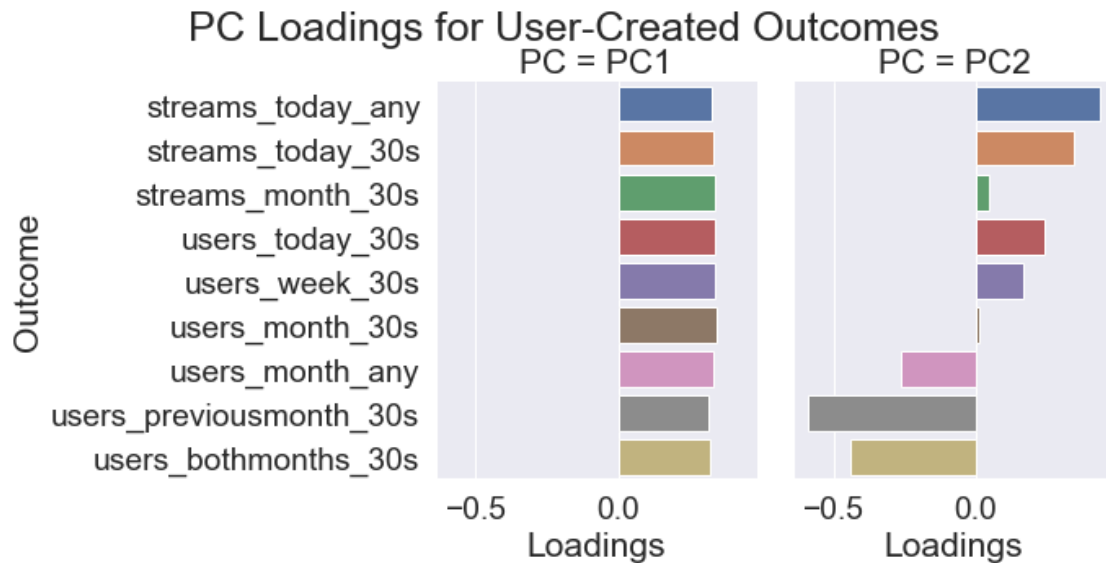
    pca_df = (pd.DataFrame(pca_loadings.T,
                           columns = ['PC1', 'PC2'],
                           index = data.columns).
               reset_index().
               rename(columns = {'index': 'outcome'}).
               melt(id_vars = 'outcome',
                    value_vars = ['PC1', 'PC2'],
                    var_name = 'PC',
                    value_name = 'loading')
               )

    load_plot = sns.catplot(data = pca_df,
                             kind = 'bar',
                             y = 'outcome',
                             x = 'loading',
                             col = 'PC')

    load_plot.set_ylabels('Outcome')
    load_plot.set_xlabels('Loadings')
    load_plot.fig.suptitle(title, y = 1)
```



```
[13]: ###capture
plot_pca_loadings(data_scaled = nonspotify_outcome_scaled,
                  data = nonspotify_outcome_data,
                  title = 'PC Loadings for User-Created Outcomes')
```



5.2 Exploration of Predictors of Interest

- We will briefly describe the distribution of the numeric predictors (number of tracks, artists and albums) and the categorical predictors (genre and mood).
- We create inverse measures of “musical diversity” as measured by number of tracks per artist and number of tracks per album.

```
[14]: predictor_data = raw_data.loc[:, ['spotify_created'] + predictor_vars]
```

```
[15]: np.random.seed(2)
continuous_predictors = (predictor_data.
                        loc[:, ['spotify_created',
                                'n_tracks',
                                #'n_local_tracks',
                                'n_artists',
                                'n_albums',]],
                        groupby('spotify_created').
                        # melt(var_name = 'predictor',
                        #     value_name = 'value').
                        sample(399)
                        )

# sns.pairplot(continuous_predictors,
```

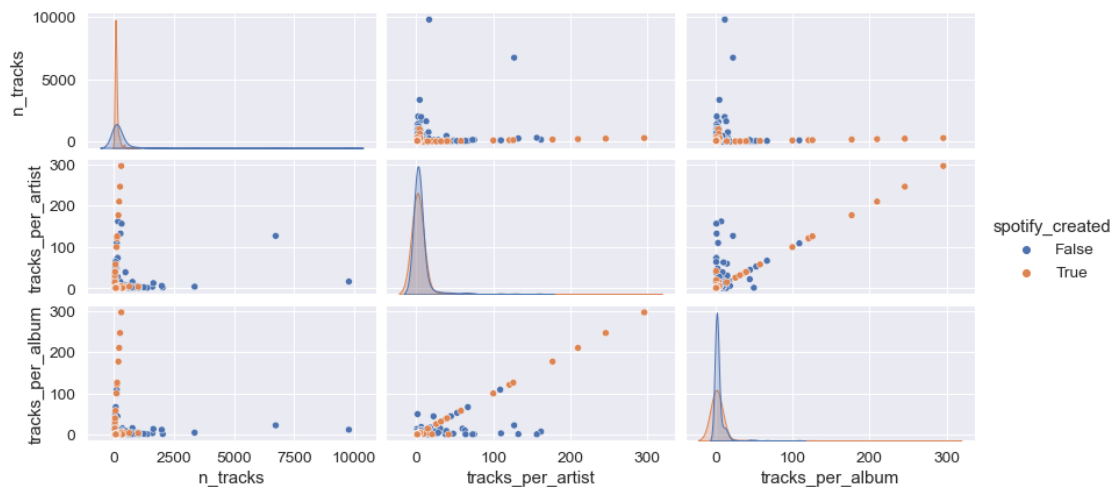
```
#             hue = 'spotify_created')
```

```
[16]: predictor_data['tracks_per_artist'] = predictor_data.n_tracks / predictor_data.  
      ↪ n_artists  
      predictor_data['tracks_per_album'] = predictor_data.n_tracks / predictor_data.  
      ↪ n_albums
```

- We drew a random sample of user-created playlists for the plots below in order for them not to visually swamp the Spotify-created playlists.
- Spotify-created playlists exhibit a bit more artist and album diversity than the user-created playlists, but the difference is miniscule compared to the outcome variables.

```
[18]: #%%capture  
np.random.seed(2)  
scaled_predictors = (predictor_data.  
                      loc[:, ['spotify_created',  
                              'n_tracks',  
                              # 'n_local_tracks',  
                              'tracks_per_artist',  
                              'tracks_per_album'],].  
                      groupby('spotify_created').  
                      # melt(var_name = 'predictor',  
                      #      value_name = 'value').  
                      sample(399)  
                      )
```

```
sns.set(font_scale = 1.25)  
sns.pairplot(scaled_predictors,  
            hue = 'spotify_created',  
            height = 2,  
            aspect = 2);
```



- We now turn our attention to describing the categorical predictors: genre and mood.
- In results now shown, the distribution of genre and mood are similar for user-created and Spotify-created playlists, so we combine them in the plots below.

```
[19]: cat_vars = ['genre_1',
                  'genre_2',
                  'genre_3',
                  'mood_1',
                  'mood_2',
                  'mood_3'
                  ]

cat_preds = (predictor_data.
             loc[:, ['spotify_created'] + cat_vars]
             )
```

```
[20]: def get_cat_order(pred):

    cat_order = (cat_preds[pred].
                  value_counts().
                  reset_index().
                  rename(columns = {pred: 'count',
                                    'index': 'level'
                                    }).
                  sort_values('count', ascending = False)
                  )

    cat_order['variable'] = pred

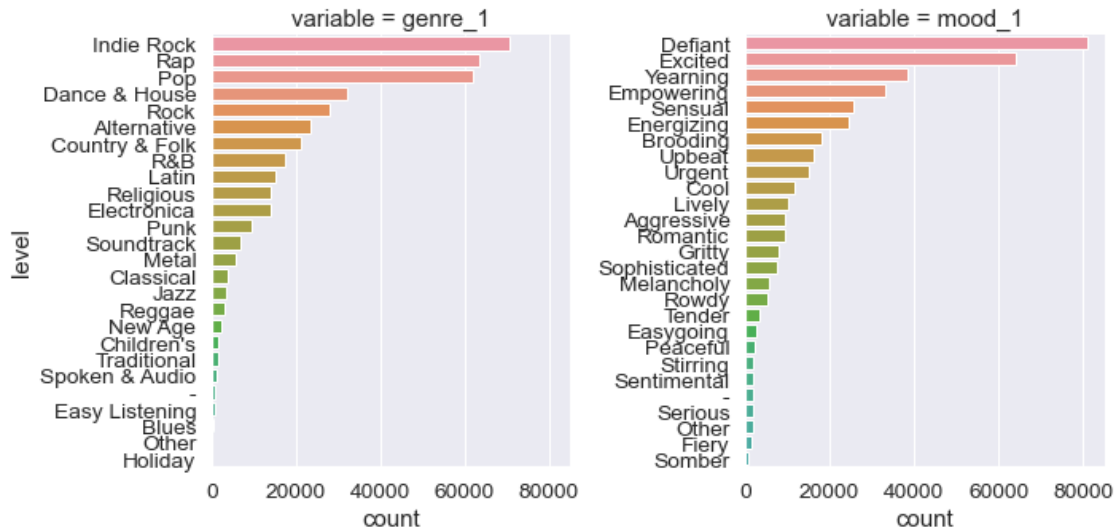
    return cat_order

genre_cat_order = get_cat_order('genre_1')
mood_cat_order = get_cat_order('mood_1')

new_cat_order = genre_cat_order.append(mood_cat_order)
```

```
[21]: my_plot = sns.catplot(
    data = new_cat_order,
    kind = 'bar',
    y = 'level',
    x = 'count',
    col = 'variable',
    sharey = False
```

)



```
[22]: def plot_cat_var(pred):  
  
    cat_order = (cat_preds[pred].  
                  value_counts().  
                  reset_index().  
                  rename(columns = {pred: 'count',  
                                     'index': pred  
                                     })).  
    sort_values('count', ascending = False)  
    )  
  
    my_plot = sns.barplot(  
        data = cat_order,  
        y = pred,  
        x = 'count'  
    )  
  
    my_plot.set_title(f'Distribution of {pred}')  
  
    return my_plot
```

```
[23]: ###capture  
# sns.set(font_scale = 1)  
# plot_cat_var('genre_1');
```

```
[24]: ###capture  
# plot_cat_var('mood_1');
```

6 Predicting Playlist Success

- We now seek to predict playlist success as measured by the number of users with stream greater than 30 seconds today.
- Separate models are fit for the user-created and Spotify-created playlists.
- For sake of brevity, we leave models predicting number of users with a stream in the current and previous months for the Appendix.
- All modeling results presented here should be taken as hypothesis-generating rather than hypothesis-confirming. This is fundamentally a data-mining exercise rather than a testing of well-defined hypotheses defined prior to the construction of models.

6.1 The Modeling Process

- Train a lightgbm gradient boosting model using tenfold cross validation (CV).
- Use early stopping to avoid overfitting and prevent wasted time and compute resources.
- Use the Huber loss function to reduce the effect of the extreme outliers present in our data.
- Based on best number of trees from CV, refit model using all available data.
- We use SHAP values to draw inference from the models.
 1. Create plots of the most important predictors ranked by mean absolute SHAP values.
 2. For genre and mood, create plots that highlight which levels of the variable are related to higher and lower success.
- Deriving measures of inferential uncertainty from gradient boosting models can be challenging. Thus, we also fit a robust M-estimation linear regression with a Huber loss function (results in Appendix).

```
[25]: def fit_lgb_mod(X, y):

    lgb_data = lgb.Dataset(data = X,
                           label = y,
                           params = {'verbose': -1}
                           )

    lgb_params = {'learning_rate': 1.0,
                  'objective': 'huber' # deal with outliers
                  }

    lgb_cv = lgb.cv(params = lgb_params,
                    train_set = lgb_data,
                    num_boost_round = 10000,
                    nfold = 10,
                    verbose_eval = -1,
                    stratified = False,
```

```

        early_stopping_rounds = 25,
        return_cvbooster = True)

lgb_booster = lgb_cv.get('cvbooster')
n_trees = lgb_booster.best_iteration

lgb_mod = lgb.train(params = lgb_params,
                    train_set = lgb_data,
                    num_boost_round = n_trees,
                    verbose_eval = -1)

return lgb_mod

```

```

[26]: X = predictor_data.drop(columns = ['n_artists', 'n_albums'])
cat_columns = ['genre_1', 'genre_2', 'genre_3', 'mood_1', 'mood_2', 'mood_3']
for column in cat_columns:
    X[column] = X[column].astype('category')

```

```

[27]: # Get the modeling input data for the Spotify-only models and no-Spotify models
X_spotify = X.query('spotify_created == True').drop(columns = 'spotify_created')
users_today_spotify = raw_data.query('spotify_created == True').loc[:,
    ↪ 'users_today_30s']
users_bothmonths_spotify = raw_data.query('spotify_created == True').loc[:,
    ↪ 'users_bothmonths_30s']

# Get the modeling input data for the Spotify-only models and no-Spotify models
X_nonspotify = X.query('spotify_created == False').drop(columns =
    ↪ 'spotify_created')
users_today_nonspotify = raw_data.query('spotify_created == False').loc[:,
    ↪ 'users_today_30s']
users_bothmonths_nonspotify = raw_data.query('spotify_created == False').loc[:,
    ↪ 'users_bothmonths_30s']

```

```

[28]: %%capture
users_today_spotify_mod = fit_lgb_mod(X = X_spotify,
                                     y = users_today_spotify)

```

```

[29]: class shap_model_explanation():

    def __init__(self, mod, X):
        self.mod = mod
        self.X = X

    def create_shap_values(self):
        explainer = shap.TreeExplainer(self.mod)
        self.shap_values = explainer.shap_values(self.X)

```

```

def plot_summary(self):
    shap_summary = shap.summary_plot(self.shap_values, self.X)
    plt.show(shap_summary)

def plot_dependence(self):
    fig, (ax1, ax2) = plt.subplots(1,2, figsize = (15, 6))
    shap.dependence_plot('genre_1',
                          self.shap_values,
                          self.X,
                          show = False,
                          interaction_index = None,
                          ax = ax1)
    shap.dependence_plot('mood_1',
                          self.shap_values,
                          self.X,
                          show = False,
                          interaction_index = None,
                          ax = ax2)

    plt.show()

```

```

[30]: spotify_today_exp = shap_model_explanation(mod = users_today_spotify_mod,
                                              X = X_spotify)

spotify_today_exp.create_shap_values()

```

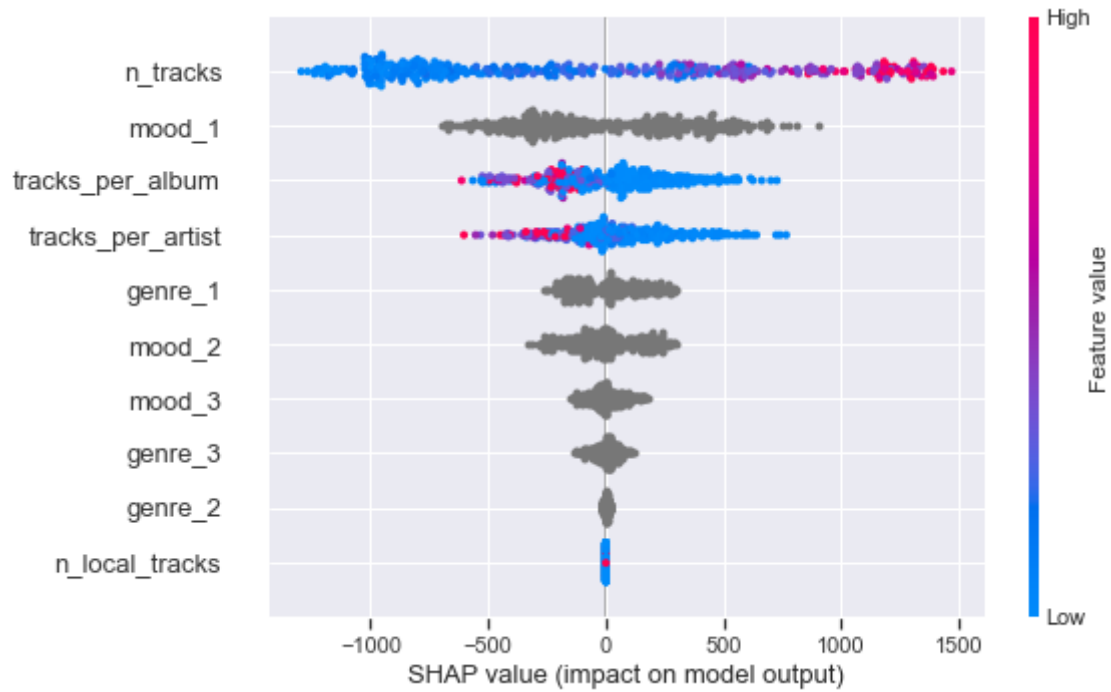
6.2 Spotify-Created Playlists: Modeling Results for Number of Users Today

- Amongst the Spotify-created playlists, the most important variable predicting success is the number of tracks, followed by the primary mood, tracks per album and tracks per artist.
- Playlists with more tracks tend to have more users today.
- Increasing the musical diversity of a playlist tends to increase playlist success as measured by number of users today.

```

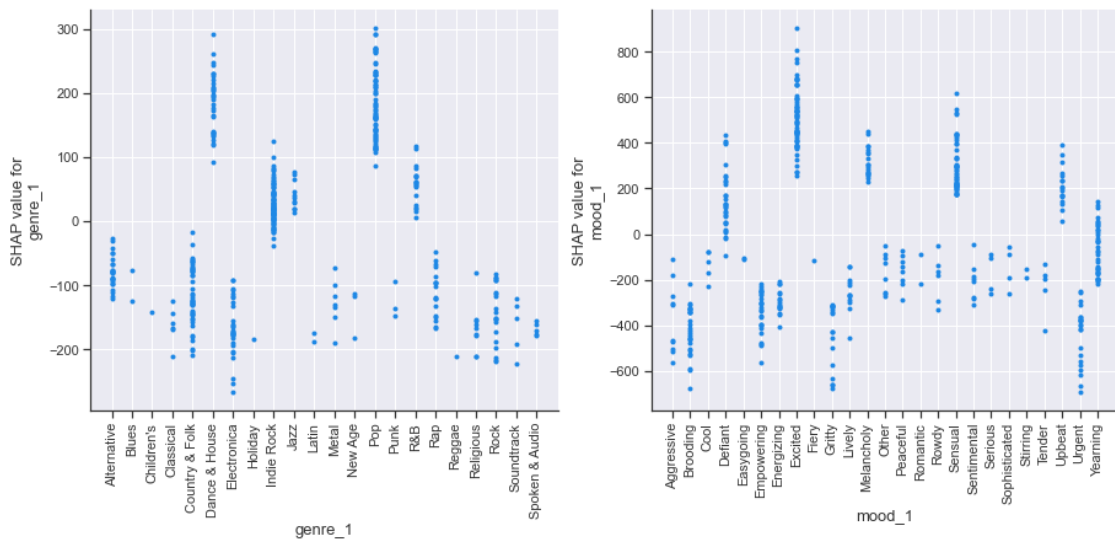
[31]: spotify_today_exp.plot_summary()

```



6.2.1 Spotify-Created Playlists: Genre and Mood Effects

[32]: `spotify_today_exp.plot_dependence()`



6.3 User-Created Playlists: Modeling Results for Number of Users Today

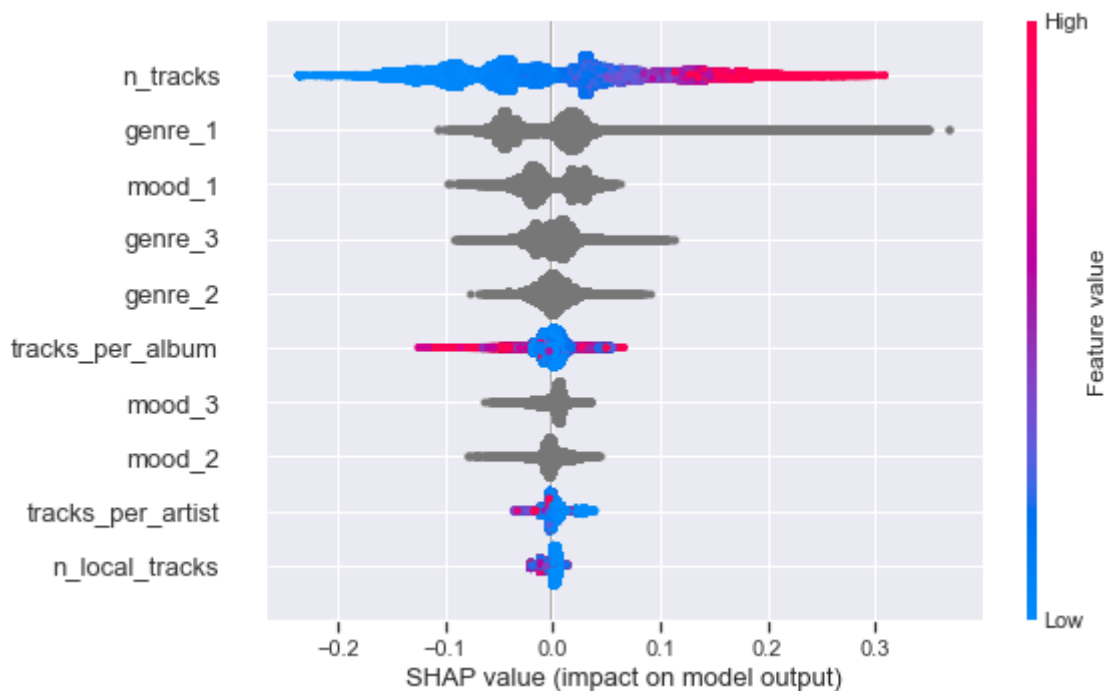
- Increasing number of tracks is also associated with increased number of users today.
- Genre appears relatively more important than in the Spotify playlists, with the most successful genres being “Children’s”, “Latin” and “Traditional”.
- In a robust linear regression model, each of these genres are found to be significantly more successful than genre = “-”.

```
[33]: %%capture
users_today_nonspotify_mod = fit_lgb_mod(X = X_nonspotify,
                                         y = users_today_nonspotify)
```

```
[34]: nonspotify_today_exp = shap_model_explanation(mod = users_today_nonspotify_mod,
                                                  X = X_nonspotify)

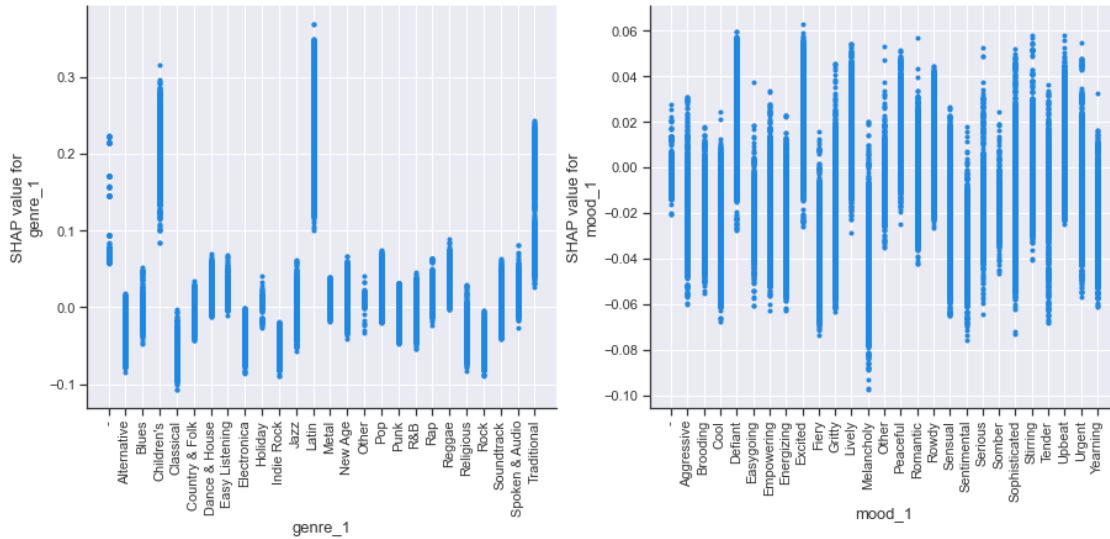
nonspotify_today_exp.create_shap_values()
```

```
[35]: nonspotify_today_exp.plot_summary()
```



6.3.1 User-Created Playlists: Genre and Mood Effects

```
[36]: nonspotify_today_exp.plot_dependence()
```



7 Appendix

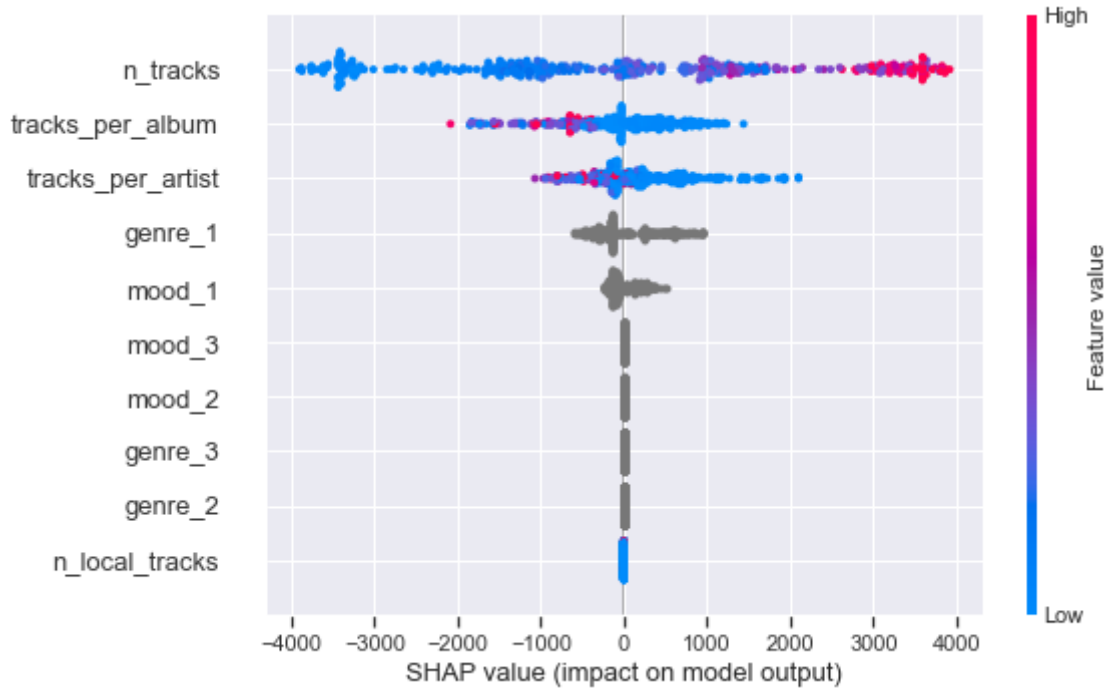
7.1 Spotify-Created Playlists: Modeling Results for Number of Users in Current and Previous Month

```
[37]: %%capture
users_bothmonths_spotify_mod = fit_lgb_mod(X = X_spotify,
                                           y = users_bothmonths_spotify)
```

```
[38]: spotify_bothmonths_exp = shap_model_explanation(mod = ↳
↳users_bothmonths_spotify_mod,
                                           X = X_spotify)

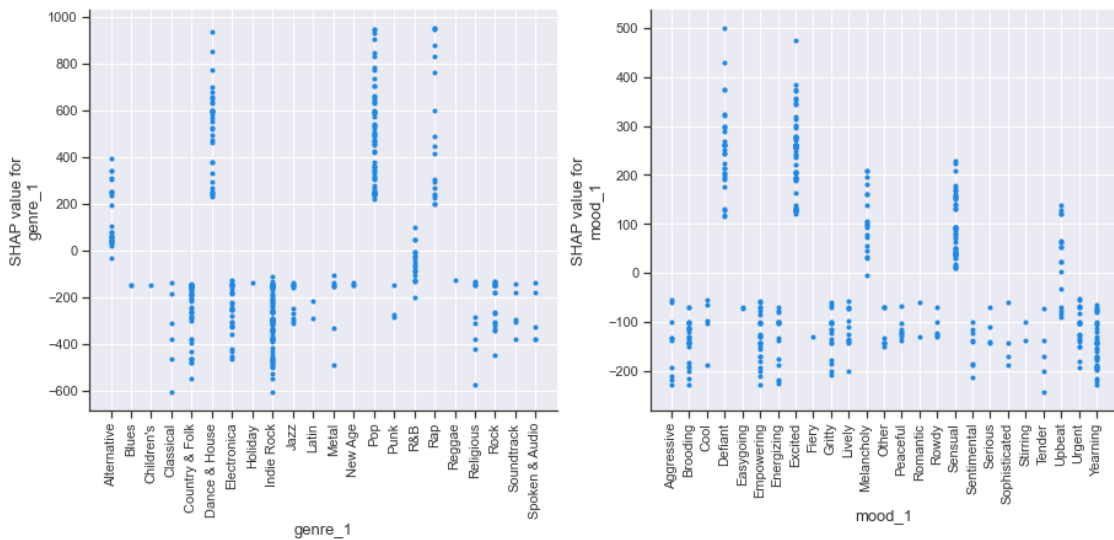
spotify_bothmonths_exp.create_shap_values()
```

```
[39]: spotify_bothmonths_exp.plot_summary()
```



7.2 Spotify-Created Playlists: Modeling Results for Number of Users in Current and Previous Month

[40]: `spotify_bothmonths_exp.plot_dependence()`



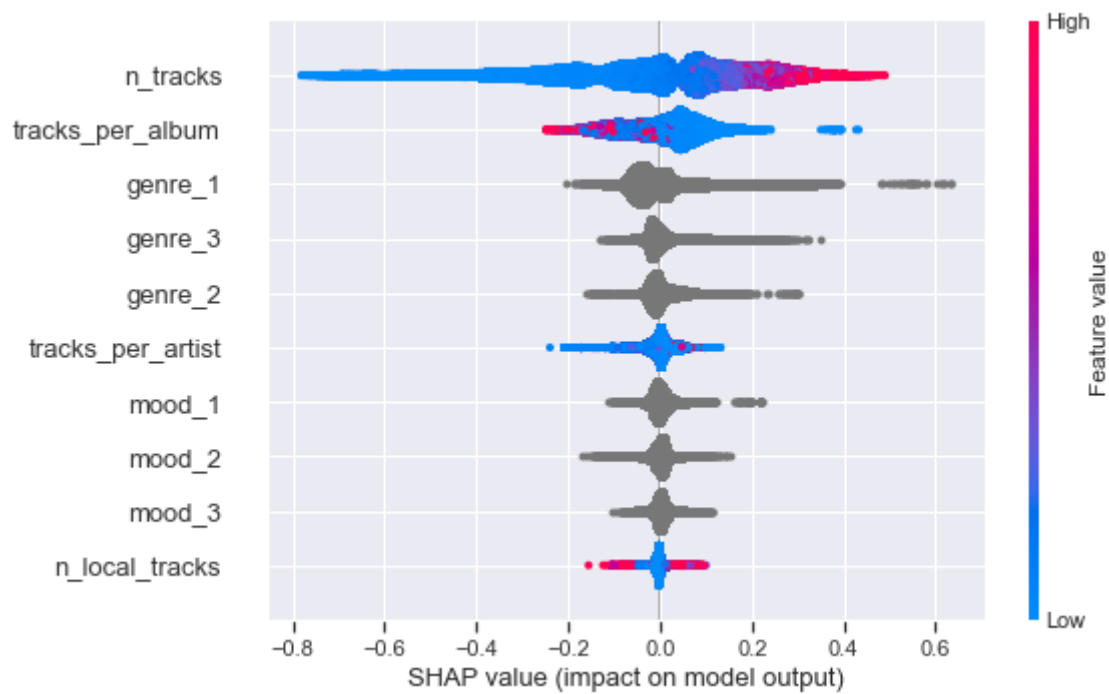
7.3 User-Created Playlists: Modeling Results for Number of Users in Current and Previous Month

```
[41]: %%capture
users_bothmonths_nonspotify_mod = fit_lgb_mod(X = X_nonspotify,
                                              y = users_bothmonths_nonspotify)
```

```
[42]: nonspotify_bothmonths_exp = shap_model_explanation(mod = users_bothmonths_nonspotify_mod,
                                                       X = X_nonspotify)

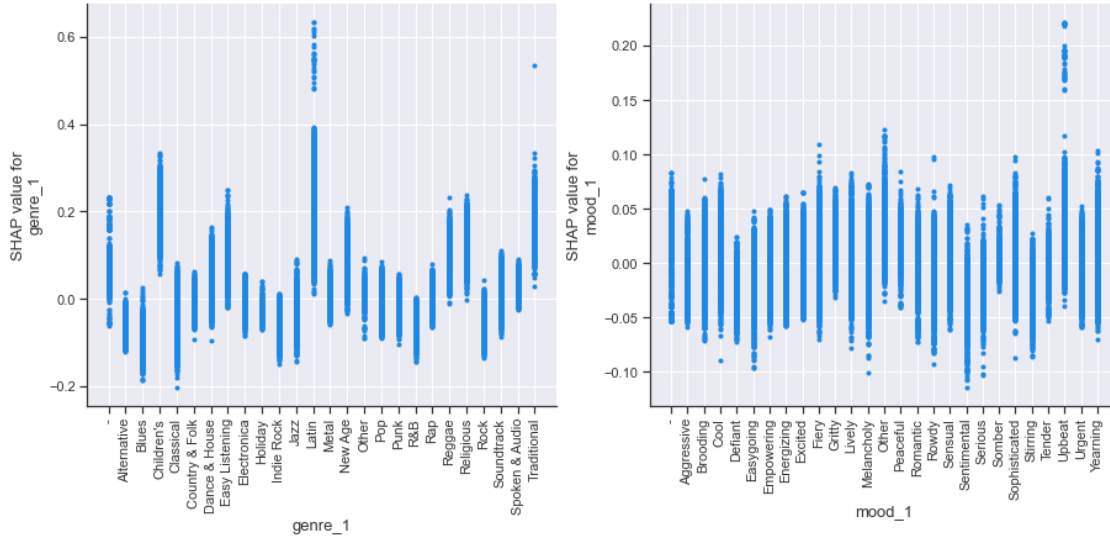
nonspotify_bothmonths_exp.create_shap_values()
```

```
[43]: nonspotify_bothmonths_exp.plot_summary()
```



7.4 User-Created Playlists: Modeling Results for Number of Users in Current and Previous Month

```
[44]: nonspotify_bothmonths_exp.plot_dependence()
```

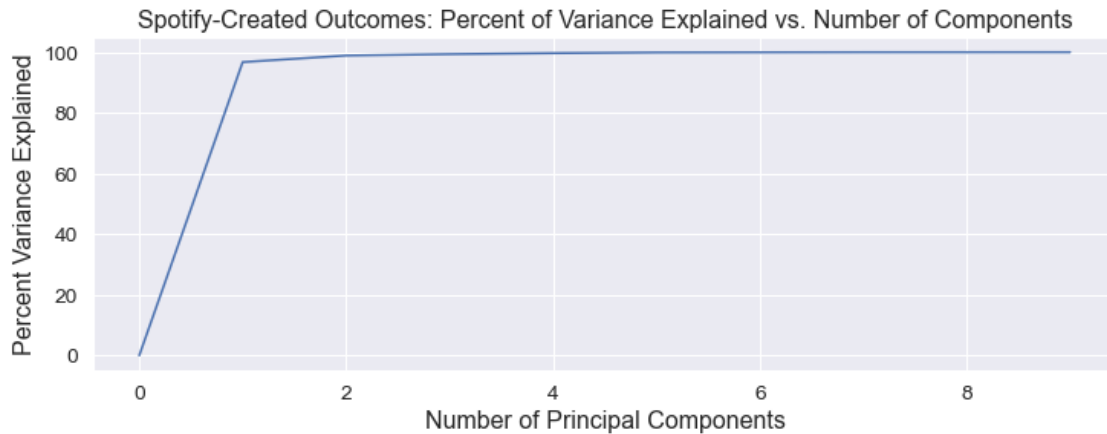


7.5 Assumptions:

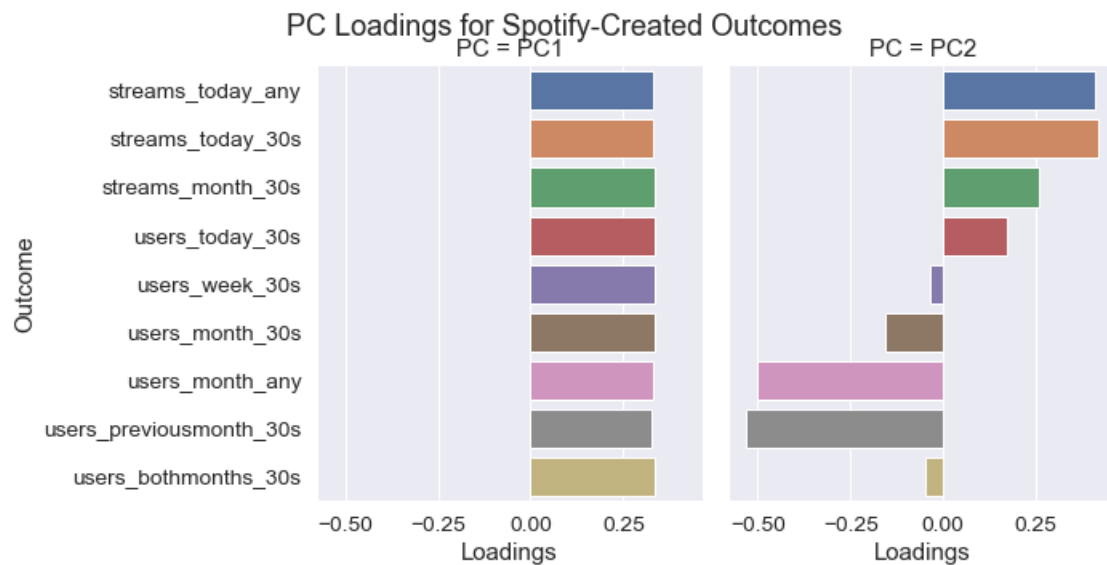
1. This is a random sample of playlists. If this is a biased sample, then any generalizations that we make from the data is likely to be meaningfully inaccurate.
2. Spotify treats each user-created playlist equally in terms of promotion. For examples, if the Spotify algorithms were promoting some genres 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.
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.

7.6 PCA Analysis of Outcomes for Spotify-Created Playlists

```
[45]: plot_percent_variation_explained(  
    spotify_outcome_scaled,  
    title = 'Spotify-Created Outcomes: Percent of Variance Explained vs. ↵  
    ↪Number of Components'  
);
```

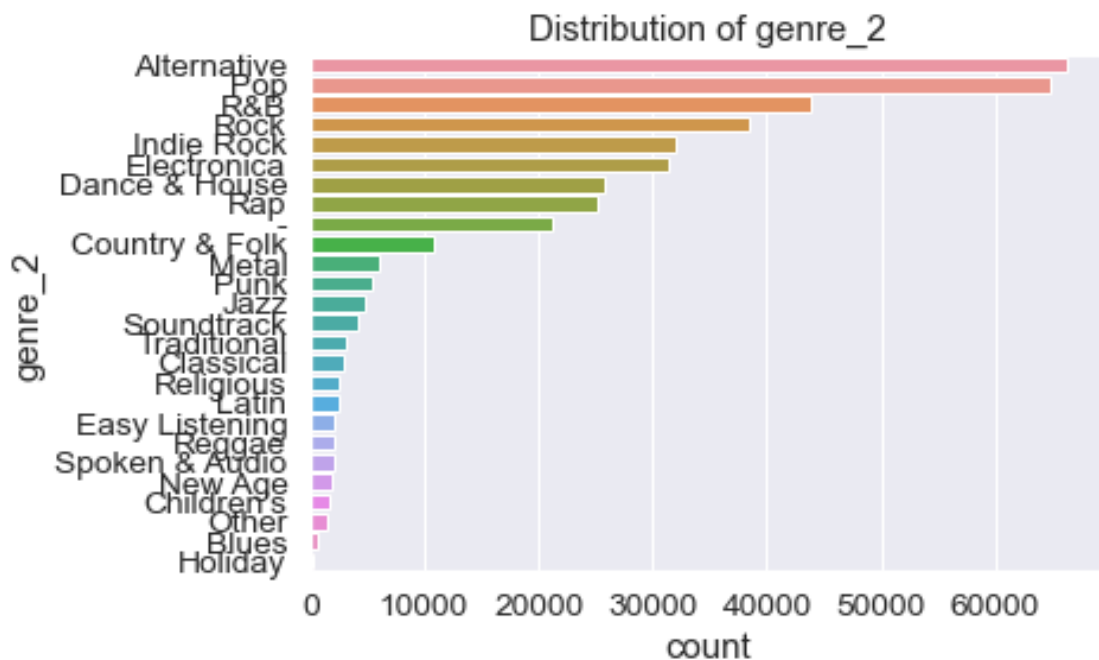


```
[46]: plot_pca_loadings(data_scaled = spotify_outcome_scaled,  
    data = spotify_outcome_data,  
    title = 'PC Loadings for Spotify-Created Outcomes')
```

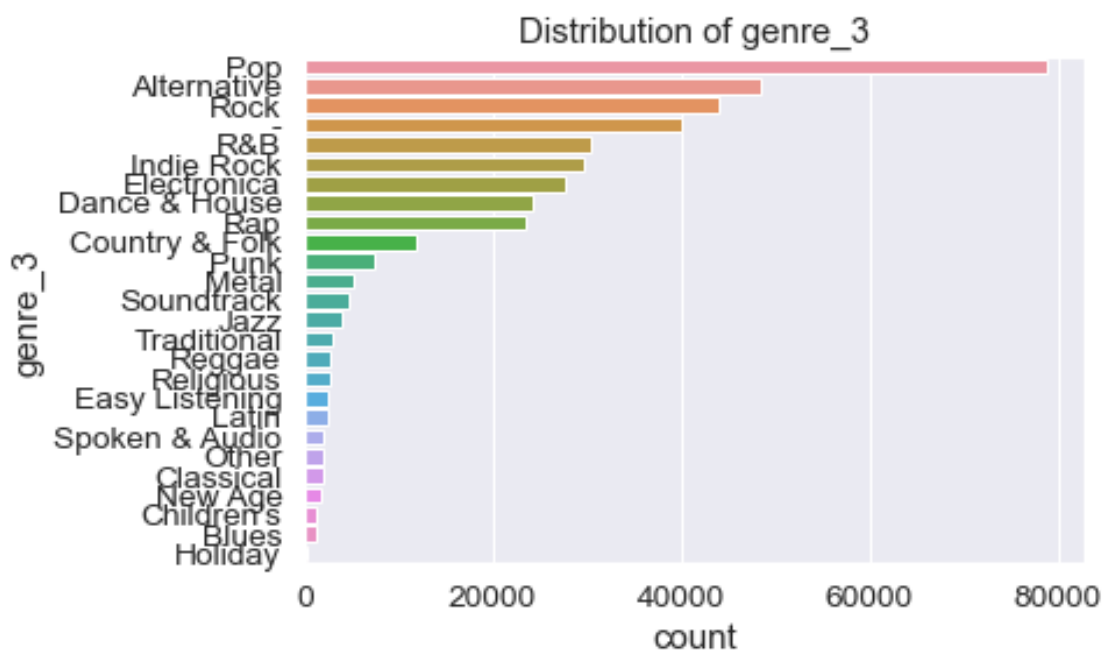


7.7 Plots of Distribution of Secondary and Tertiary Genre and Mood

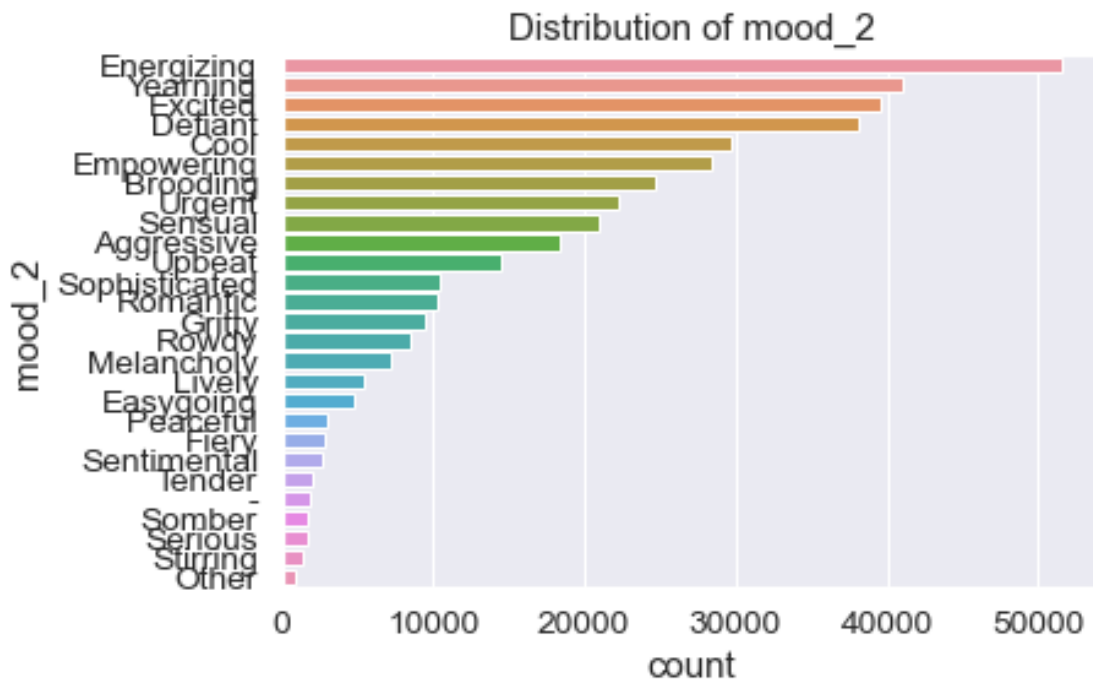
```
[47]: plot_cat_var('genre_2');
```



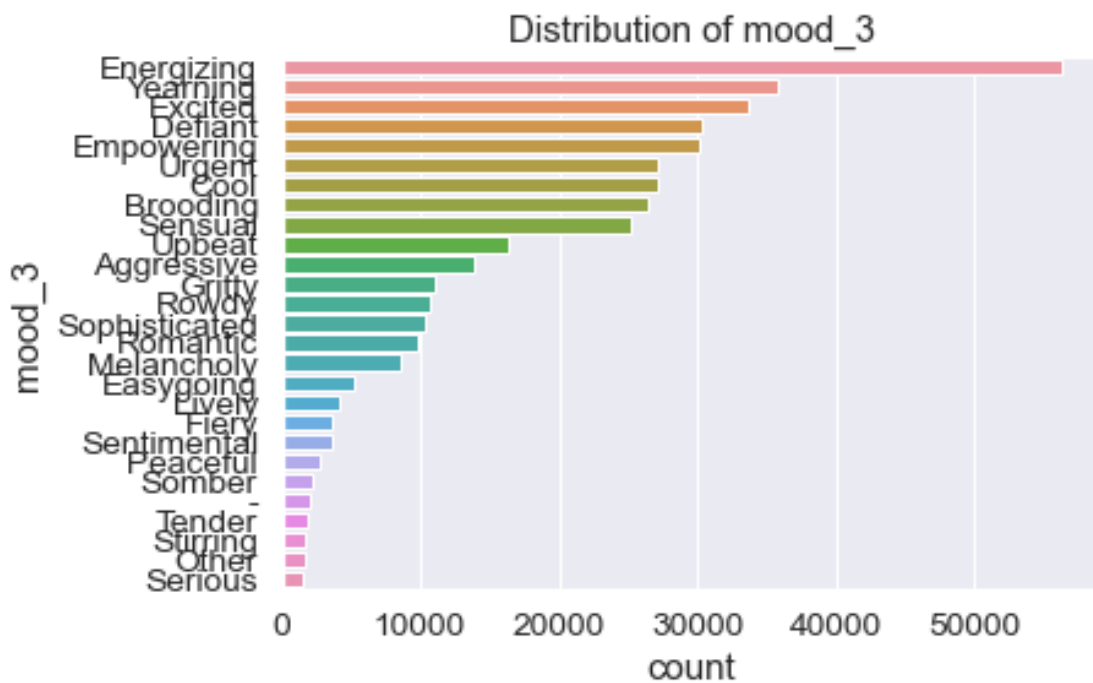
```
[48]: plot_cat_var('genre_3');
```



```
[49]: plot_cat_var('mood_2');
```



```
[50]: plot_cat_var('mood_3');
```



7.8 Robust Linear Regression Model Results

7.8.1 RLM Model for Spotify-Created Playlists, Outcome = Users Today

```
[51]: user_data = pd.concat([users_today_nonspotify, users_bothmonths_nonspotify,
    ↪X_nonspotify], axis = 1)
spotify_data = pd.concat([users_today_spotify, users_bothmonths_spotify,
    ↪X_spotify], axis = 1)
```

```
[52]: def fit_huber_model(data, outcome):

    y, X = dmatrices(f'{outcome} ~ n_tracks + tracks_per_album + genre_1 +
    ↪mood_1',
                    data = data,
                    return_type = 'dataframe')

    huber_t = sm.RLM(y, X, M = sm.robust.norms.HuberT()).fit()

    print(huber_t.summary2())

    return huber_t
```

```
[53]: fit_huber_model(data = spotify_data,
    outcome = 'users_today_30s')
```

Results: Robust linear model

```
=====
==
Model:                                RLM                                Df Residuals:          351
Dependent Variable:                  users_today_30s                    Norm:
HuberT
Date:                                2022-02-25 12:28                    Scale Est.:          mad
No. Observations:                    399                                Cov. Type:            H1
Df Model:                            47                                Scale:
1289.6
-----
--
                                Coef.      Std.Err.    z    P>|z|    [0.025
0.975]
-----
--
Intercept                          592.5850   210.9870   2.8086  0.0050   179.0581
1006.1119
genre_1[T.Alternative]              321.4248   387.0639   0.8304  0.4063  -437.2066
1080.0561
```

genre_1[T.Blues] 4013.7515	1730.9971	1164.6920	1.4862	0.1372	-551.7574
genre_1[T.Children's] 2697.7050	-445.0210	1603.4611	-0.2775	0.7814	-3587.7470
genre_1[T.Classical] 3281.6336	1301.9412	1010.0657	1.2890	0.1974	-677.7513
genre_1[T.Country & Folk] 546.3991	-169.5201	365.2716	-0.4641	0.6426	-885.4392
genre_1[T.Dance & House] 1400.2630	660.8945	377.2358	1.7519	0.0798	-78.4740
genre_1[T.Easy Listening] 0.0000	-0.0000	0.0000	-0.5453	0.5855	-0.0000
genre_1[T.Electronica] 191.2900	-603.0586	405.2874	-1.4880	0.1368	-1397.4072
genre_1[T.Holiday] 2961.9536	-199.2876	1612.9078	-0.1236	0.9017	-3360.5288
genre_1[T.Indie Rock] 823.0166	266.7984	283.7900	0.9401	0.3472	-289.4198
genre_1[T.Jazz] 730.0144	-360.3224	556.3045	-0.6477	0.5172	-1450.6592
genre_1[T.Latin] 2261.4439	-5.5204	1156.6357	-0.0048	0.9962	-2272.4847
genre_1[T.Metal] 1252.7475	-147.9796	714.6698	-0.2071	0.8360	-1548.7068
genre_1[T.New Age] 1390.6299	-570.1787	1000.4309	-0.5699	0.5687	-2530.9872
genre_1[T.Other] 0.0000	0.0000	0.0000	1.7542	0.0794	-0.0000
genre_1[T.Pop] 1276.5127	695.8092	296.2827	2.3485	0.0189	115.1058
genre_1[T.Punk] 1901.7465	17.1923	961.5249	0.0179	0.9857	-1867.3619
genre_1[T.R&B] 1252.4876	365.5038	452.5511	0.8077	0.4193	-521.4800
genre_1[T.Rap] 1279.5104	296.8161	501.3839	0.5920	0.5539	-685.8783
genre_1[T.Reggae] 3737.3008	198.2037	1805.6950	0.1098	0.9126	-3340.8934
genre_1[T.Religious] 269.8714	-931.6740	613.0447	-1.5197	0.1286	-2133.2195
genre_1[T.Rock] 303.8178	-516.6896	418.6339	-1.2342	0.2171	-1337.1970
genre_1[T.Soundtrack] 826.4040	-1068.7898	966.9534	-1.1053	0.2690	-2963.9836
genre_1[T.Spoken & Audio] 2118.3000	-244.9543	1205.7641	-0.2032	0.8390	-2608.2086
genre_1[T.Traditional] 0.0000	-0.0000	0.0000	-0.3063	0.7594	-0.0000

mood_1[T.Aggressive] 984.3709	-254.2457	631.9588	-0.4023	0.6875	-1492.8622
mood_1[T.Brooding] 273.7997	-437.9987	363.1691	-1.2060	0.2278	-1149.7970
mood_1[T.Cool] 1529.7971	-221.0991	893.3308	-0.2475	0.8045	-1971.9953
mood_1[T.Defiant] 888.2322	140.3623	381.5733	0.3679	0.7130	-607.5076
mood_1[T.Easygoing] 1317.0227	-968.6157	1166.1635	-0.8306	0.4062	-3254.2542
mood_1[T.Empowering] 367.1239	-348.6438	365.1943	-0.9547	0.3397	-1064.4115
mood_1[T.Energizing] 731.9462	-157.2315	453.6704	-0.3466	0.7289	-1046.4091
mood_1[T.Excited] 1285.2594	700.3149	298.4466	2.3465	0.0189	115.3704
mood_1[T.Fiery] 1970.0611	-1175.7843	1605.0526	-0.7326	0.4638	-4321.6296
mood_1[T.Gritty] 118.7260	-759.8694	448.2712	-1.6951	0.0901	-1638.4648
mood_1[T.Lively] 822.8642	-206.9745	525.4376	-0.3939	0.6936	-1236.8133
mood_1[T.Melancholy] 1166.5315	346.0206	418.6357	0.8265	0.4085	-474.4904
mood_1[T.Other] 1164.9759	-790.2208	997.5676	-0.7921	0.4283	-2745.4174
mood_1[T.Peaceful] 1800.7184	594.4401	615.4594	0.9658	0.3341	-611.8382
mood_1[T.Romantic] 3052.4426	167.0885	1472.1465	0.1135	0.9096	-2718.2655
mood_1[T.Rowdy] 1211.1502	-126.0988	682.2825	-0.1848	0.8534	-1463.3478
mood_1[T.Sensual] 1049.7519	372.3687	345.6100	1.0774	0.2813	-305.0144
mood_1[T.Sentimental] 2154.9162	951.1252	614.1904	1.5486	0.1215	-252.6658
mood_1[T.Serious] 1323.0182	-504.0052	932.1720	-0.5407	0.5887	-2331.0287
mood_1[T.Somber] 0.0000	-0.0000	0.0000	-0.8718	0.3833	-0.0000
mood_1[T.Sophisticated] 4734.9158	3097.9696	835.1919	3.7093	0.0002	1461.0235
mood_1[T.Stirring] 1722.1606	-778.8670	1276.0579	-0.6104	0.5416	-3279.8946
mood_1[T.Tender] 2847.4945	989.2884	948.0818	1.0435	0.2967	-868.9177
mood_1[T.Upbeat] 1624.3361	752.1248	445.0139	1.6901	0.0910	-120.0865

```

mood_1[T.Urgent]          -625.8362   399.0035  -1.5685  0.1168  -1407.8687
156.1962
mood_1[T.Yearning]        -163.0275   328.3333  -0.4965  0.6195   -806.5490
480.4940
n_tracks                   4.3197     1.0171   4.2470  0.0000     2.3261
6.3132
tracks_per_album           23.7652     3.3748   7.0419  0.0000    17.1507
30.3798
=====
==

```

[53]: <statsmodels.robust.robust_linear_model.RLMResultsWrapper at 0x1d1de7290a0>

7.8.2 RLM Model of User-Created Playlists, Outcome = Users Today

Note that for the following model outputs, the baseline level for both the genre and mood categorical predictors is “-”, which I am taking to mean there is no easily discernible genre or mood detectable for the given playlist.

```
[54]: fit_huber_model(data = user_data,
                      outcome = 'users_today_30s')
```

```

Results: Robust linear model
=====
Model:                RLM                Df Residuals:      402913
Dependent Variable:    users_today_30s      Norm:              HuberT
Date:                 2022-02-25 12:29      Scale Est.:         mad
No. Observations:     402967              Cov. Type:          H1
Df Model:              53                 Scale:              0.73545
-----
              Coef.  Std.Err.   z    P>|z|    [0.025  0.975]
-----
Intercept          0.4775    0.0191 25.0304 0.0000   0.4401   0.5149
genre_1[T.Alternative] -0.0373   0.0280 -1.3342 0.1821  -0.0922   0.0175
genre_1[T.Blues]     -0.0499   0.0392 -1.2741 0.2026  -0.1268   0.0269
genre_1[T.Children's]  0.1435   0.0312  4.6023 0.0000   0.0824   0.2046
genre_1[T.Classical] -0.0907   0.0293 -3.0895 0.0020  -0.1482  -0.0331
genre_1[T.Country & Folk]  0.0057   0.0281  0.2024 0.8396  -0.0495   0.0609
genre_1[T.Dance & House]  0.0254   0.0280  0.9064 0.3648  -0.0295   0.0802
genre_1[T.Easy Listening]  0.0205   0.0338  0.6068 0.5440  -0.0457   0.0867
genre_1[T.Electronica] -0.0404   0.0282 -1.4323 0.1520  -0.0956   0.0149
genre_1[T.Holiday]     0.0271   0.0542  0.5004 0.6168  -0.0792   0.1335
genre_1[T.Indie Rock]  -0.0407   0.0279 -1.4625 0.1436  -0.0953   0.0139
genre_1[T.Jazz]        -0.0234   0.0296 -0.7878 0.4308  -0.0815   0.0347
genre_1[T.Latin]        0.2287   0.0282  8.1150 0.0000   0.1734   0.2839
genre_1[T.Metal]        0.0101   0.0289  0.3495 0.7267  -0.0465   0.0667
genre_1[T.New Age]     -0.0051   0.0302 -0.1702 0.8649  -0.0643   0.0540

```

genre_1[T.Other]	-0.0759	0.0519	-1.4627	0.1436	-0.1775	0.0258
genre_1[T.Pop]	0.0018	0.0278	0.0654	0.9478	-0.0527	0.0564
genre_1[T.Punk]	-0.0318	0.0284	-1.1182	0.2635	-0.0875	0.0239
genre_1[T.R&B]	-0.0094	0.0281	-0.3351	0.7375	-0.0646	0.0457
genre_1[T.Rap]	0.0399	0.0279	1.4292	0.1529	-0.0148	0.0946
genre_1[T.Reggae]	0.0639	0.0297	2.1478	0.0317	0.0056	0.1221
genre_1[T.Religious]	-0.0631	0.0282	-2.2376	0.0253	-0.1183	-0.0078
genre_1[T.Rock]	-0.0452	0.0280	-1.6176	0.1058	-0.1000	0.0096
genre_1[T.Soundtrack]	-0.0039	0.0286	-0.1350	0.8926	-0.0599	0.0522
genre_1[T.Spoken & Audio]	-0.0047	0.0332	-0.1404	0.8883	-0.0698	0.0605
genre_1[T.Traditional]	0.0993	0.0314	3.1655	0.0015	0.0378	0.1609
mood_1[T.Aggressive]	-0.0045	0.0212	-0.2115	0.8325	-0.0460	0.0370
mood_1[T.Brooding]	-0.0226	0.0207	-1.0891	0.2761	-0.0632	0.0180
mood_1[T.Cool]	-0.0174	0.0210	-0.8280	0.4077	-0.0586	0.0238
mood_1[T.Defiant]	0.0532	0.0204	2.6097	0.0091	0.0133	0.0932
mood_1[T.Easygoing]	-0.0553	0.0233	-2.3711	0.0177	-0.1009	-0.0096
mood_1[T.Empowering]	0.0055	0.0205	0.2668	0.7896	-0.0347	0.0456
mood_1[T.Energizing]	-0.0176	0.0206	-0.8570	0.3914	-0.0579	0.0227
mood_1[T.Excited]	0.0387	0.0204	1.8972	0.0578	-0.0013	0.0786
mood_1[T.Fiery]	-0.0221	0.0256	-0.8632	0.3880	-0.0723	0.0281
mood_1[T.Gritty]	-0.0355	0.0213	-1.6634	0.0962	-0.0772	0.0063
mood_1[T.Lively]	0.0377	0.0211	1.7841	0.0744	-0.0037	0.0790
mood_1[T.Melancholy]	-0.0545	0.0217	-2.5092	0.0121	-0.0970	-0.0119
mood_1[T.Other]	0.0040	0.0254	0.1591	0.8736	-0.0457	0.0538
mood_1[T.Peaceful]	-0.0040	0.0236	-0.1716	0.8638	-0.0503	0.0422
mood_1[T.Romantic]	-0.0045	0.0211	-0.2133	0.8311	-0.0458	0.0368
mood_1[T.Rowdy]	0.0042	0.0218	0.1908	0.8487	-0.0386	0.0469
mood_1[T.Sensual]	-0.0135	0.0206	-0.6539	0.5132	-0.0538	0.0269
mood_1[T.Sentimental]	-0.0833	0.0246	-3.3923	0.0007	-0.1315	-0.0352
mood_1[T.Serious]	-0.0143	0.0250	-0.5703	0.5685	-0.0633	0.0348
mood_1[T.Somber]	-0.0623	0.0280	-2.2258	0.0260	-0.1171	-0.0074
mood_1[T.Sophisticated]	0.0033	0.0213	0.1543	0.8774	-0.0384	0.0450
mood_1[T.Stirring]	0.0200	0.0244	0.8224	0.4109	-0.0277	0.0678
mood_1[T.Tender]	0.0062	0.0229	0.2732	0.7847	-0.0386	0.0510
mood_1[T.Upbeat]	0.0372	0.0209	1.7755	0.0758	-0.0039	0.0783
mood_1[T.Urgent]	0.0178	0.0208	0.8549	0.3926	-0.0230	0.0586
mood_1[T.Yearning]	-0.0143	0.0205	-0.6987	0.4847	-0.0544	0.0258
n_tracks	0.0001	0.0000	53.0255	0.0000	0.0001	0.0001
tracks_per_album	-0.0000	0.0000	-1.6380	0.1014	-0.0000	0.0000
=====						

[54]: <statsmodels.robust.robust_linear_model.RLMResultsWrapper at 0x1d1b4cba070>

7.8.3 OLS Model of User Created Playlists, Outcome = Users Today

This is a model of number of users today for the non-Spotify playlists. Take note of the incredibly large standard errors on our estimates when we do not employ a loss function that is robust to

outliers.

```
[55]: ols_user = smf.ols('users_today_30s ~ n_tracks + tracks_per_artist + genre_1 + mood_1',  
                        data = user_data).fit()  
  
print(ols_user.summary2())
```

Results: Ordinary least squares

```
=====
Model:                OLS                Adj. R-squared:    0.000
Dependent Variable:   users_today_30s    AIC:                4140974.1152
Date:                2022-02-25 12:29    BIC:                4141563.0721
No. Observations:    402967              Log-Likelihood:     -2.0704e+06
Df Model:             53                  F-statistic:        1.180
Df Residuals:         402913              Prob (F-statistic): 0.173
R-squared:           0.000                Scale:           1699.7
=====
```

```
-----
                Coef.  Std.Err.  t    P>|t|    [0.025 0.975]
-----
Intercept                1.2035    1.3420   0.8968  0.3698  -1.4268  3.8338
genre_1[T.Alternative]  -0.5632    1.9695  -0.2860  0.7749  -4.4233  3.2969
genre_1[T.Blues]        -0.4999    2.7576  -0.1813  0.8562  -5.9047  4.9049
genre_1[T.Children's]    1.1479    2.1936   0.5233  0.6008  -3.1515  5.4473
genre_1[T.Classical]    -1.1414    2.0647  -0.5528  0.5804  -5.1882  2.9054
genre_1[T.Country & Folk] -0.6222    1.9803  -0.3142  0.7534  -4.5036  3.2591
genre_1[T.Dance & House] -0.0923    1.9686  -0.0469  0.9626  -3.9507  3.7661
genre_1[T.Easy Listening] -0.6473    2.3767  -0.2724  0.7854  -5.3056  4.0110
genre_1[T.Electronica]  -0.8544    1.9821  -0.4310  0.6664  -4.7392  3.0305
genre_1[T.Holiday]      -0.5885    3.8160  -0.1542  0.8774  -8.0678  6.8909
genre_1[T.Indie Rock]   -0.7986    1.9597  -0.4075  0.6836  -4.6396  3.0424
genre_1[T.Jazz]         -0.4959    2.0855  -0.2378  0.8120  -4.5834  3.5916
genre_1[T.Latin]         0.5322    1.9823   0.2685  0.7883  -3.3531  4.4176
genre_1[T.Metal]        -0.5730    2.0306  -0.2822  0.7778  -4.5529  3.4069
genre_1[T.New Age]      -0.9145    2.1236  -0.4306  0.6667  -5.0767  3.2477
genre_1[T.Other]        -0.8709    3.6496  -0.2386  0.8114  -8.0240  6.2821
genre_1[T.Pop]          -0.2399    1.9576  -0.1225  0.9025  -4.0766  3.5969
genre_1[T.Punk]         -0.8336    2.0002  -0.4168  0.6769  -4.7540  3.0868
genre_1[T.R&B]          -0.7183    1.9798  -0.3628  0.7167  -4.5986  3.1620
genre_1[T.Rap]          -0.3503    1.9636  -0.1784  0.8584  -4.1988  3.4983
genre_1[T.Reggae]       -0.3277    2.0920  -0.1566  0.8755  -4.4279  3.7725
genre_1[T.Religious]    -0.7106    1.9828  -0.3584  0.7201  -4.5968  3.1756
genre_1[T.Rock]         -0.0104    1.9670  -0.0053  0.9958  -3.8656  3.8448
genre_1[T.Soundtrack]    0.3482    2.0116   0.1731  0.8626  -3.5945  4.2910
genre_1[T.Spoken & Audio] -0.5224    2.3385  -0.2234  0.8232  -5.1059  4.0610
genre_1[T.Traditional]   0.6116    2.2080   0.2770  0.7818  -3.7159  4.9392
mood_1[T.Aggressive]     0.5344    1.4906   0.3585  0.7199  -2.3870  3.4559
mood_1[T.Brooding]      0.0589    1.4573   0.0404  0.9677  -2.7973  2.9151
-----
```

mood_1[T.Cool]	-0.0634	1.4776	-0.0429	0.9658	-2.9594	2.8326
mood_1[T.Defiant]	0.3676	1.4347	0.2562	0.7978	-2.4445	3.1796
mood_1[T.Easygoing]	-0.2866	1.6398	-0.1748	0.8612	-3.5005	2.9273
mood_1[T.Empowering]	0.2961	1.4410	0.2055	0.8372	-2.5281	3.1203
mood_1[T.Energizing]	0.4055	1.4465	0.2803	0.7792	-2.4297	3.2407
mood_1[T.Excited]	0.7369	1.4337	0.5140	0.6072	-2.0730	3.5468
mood_1[T.Fiery]	0.0122	1.8012	0.0068	0.9946	-3.5181	3.5426
mood_1[T.Gritty]	0.0013	1.4995	0.0009	0.9993	-2.9376	2.9402
mood_1[T.Lively]	0.2418	1.4850	0.1628	0.8706	-2.6686	3.1523
mood_1[T.Melancholy]	0.1516	1.5269	0.0993	0.9209	-2.8410	3.1442
mood_1[T.Other]	0.0283	1.7856	0.0158	0.9874	-3.4714	3.5279
mood_1[T.Peaceful]	0.1250	1.6594	0.0753	0.9399	-3.1273	3.3773
mood_1[T.Romantic]	0.9962	1.4824	0.6720	0.5016	-1.9093	3.9017
mood_1[T.Rowdy]	0.4852	1.5336	0.3164	0.7517	-2.5206	3.4909
mood_1[T.Sensual]	0.5650	1.4482	0.3901	0.6965	-2.2735	3.4034
mood_1[T.Sentimental]	-0.1200	1.7282	-0.0694	0.9446	-3.5073	3.2673
mood_1[T.Serious]	1.8377	1.7595	1.0444	0.2963	-1.6109	5.2863
mood_1[T.Somber]	0.0115	1.9681	0.0058	0.9953	-3.8459	3.8689
mood_1[T.Sophisticated]	0.1403	1.4974	0.0937	0.9253	-2.7945	3.0751
mood_1[T.Stirring]	0.1790	1.7150	0.1044	0.9169	-3.1823	3.5403
mood_1[T.Tender]	0.3622	1.6079	0.2252	0.8218	-2.7892	3.5136
mood_1[T.Upbeat]	0.7530	1.4738	0.5109	0.6094	-2.1356	3.6417
mood_1[T.Urgent]	0.2116	1.4655	0.1444	0.8852	-2.6608	3.0839
mood_1[T.Yearning]	0.1184	1.4399	0.0822	0.9345	-2.7037	2.9404
n_tracks	0.0003	0.0001	2.5696	0.0102	0.0001	0.0005
tracks_per_artist	-0.0005	0.0007	-0.7858	0.4320	-0.0018	0.0008

Omnibus:	2142164.888	Durbin-Watson:	1.999
Prob(Omnibus):	0.000	Jarque-Bera (JB):	180443014585193.719
Skew:	268.776	Prob(JB):	0.000
Kurtosis:	103668.698	Condition No.:	108308

=====

* The condition number is large (1e+05). This might indicate strong multicollinearity or other numerical problems.