

STAT 482 Final Report: Spotify Streaming Data Analysis

Jaden Oca

Spring 2024

1 Introduction

Spotify is the largest music streaming platform in the world with over 574 million users. At the end of each year, Spotify releases its annual "Wrapped," which is unique to each user and summarizes their listening activity throughout the year, including top songs, artists, and genres listened to.

The purpose of this project is to use my own listening data to extract trends beyond what Spotify's "Wrapped" gives us, including a time-series analysis of my listening trends over time, seeing how listening to certain songs and artists peak and fall over time.

More specifically, the project hopes to observe if any inflection points can be traced back to important events in the user's life, answering the question "can a shift in music taste shift be traced back to life events?"

By creating a more in-depth analysis of my own listening data, I hope to lay the groundwork for this project's formalization into a package that will allow anyone can analyze their Spotify listening trends on a deeper level.

2 Description of Dataset

2.1 Source

The data set used was pulled directly from Spotify. It tracks every single one of my Spotify streams ranging from June 24th, 2022 to January 5th, 2024.

2.2 Data

The data contains four columns. They are as follows:

- endTime - time in UTC that the song finished by the minute.

- `artistName` - name of artist of each track streamed
- `trackName` - name of each track streamed
- `msPlayed` - duration each track was listened to

The image below is an example of the head of the dataframe.

	<code>endTime</code>	<code>artistName</code>	<code>trackName</code>	<code>msPlayed</code>
0	2022-06-25 00:01	No Rome	Narcissist (feat. The 1975)	195501
1	2022-06-25 00:04	The 1975	Me & You Together Song	207223
2	2022-06-25 00:08	Valley	homebody	191086
3	2022-06-25 00:11	One Direction	Steal My Girl	228133
4	2022-06-25 00:15	HONNE	Day 1 🎧	233600

Figure 1: Head of dataframe

Evident from Figure 1 (above), there are 4 columns as mentioned before, `endTime`, `artistName`, `trackName`, and `msPlayed`.

3 Exploratory Data Analysis and Results

The exploratory data analysis for this project involved looking at top songs listened to, top tracks listened to, and observing visual trends on a time series line chart of listening behavior over the recorded period.

3.1 Artist and Track time series

To begin the analysis, observed was the length, top artists, and top tracks,

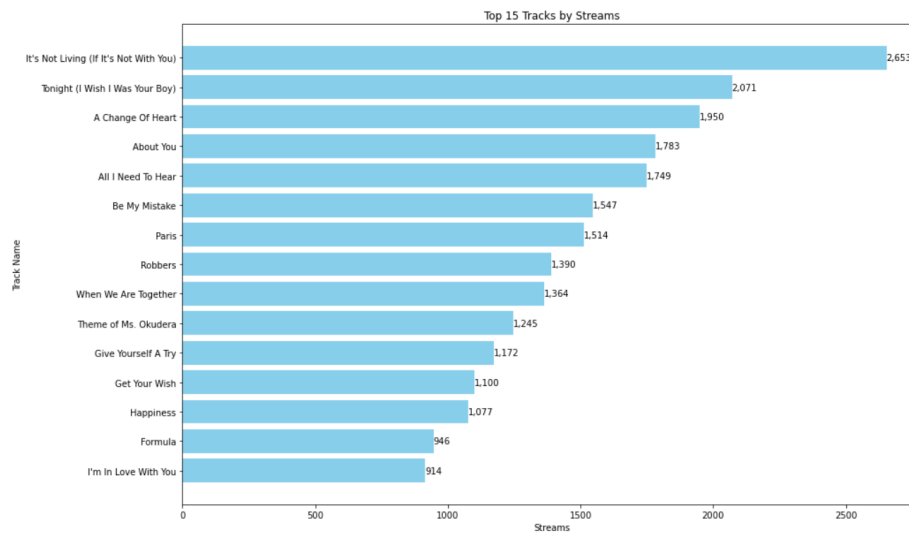


Figure 2: Stream counts of top 15 tracks

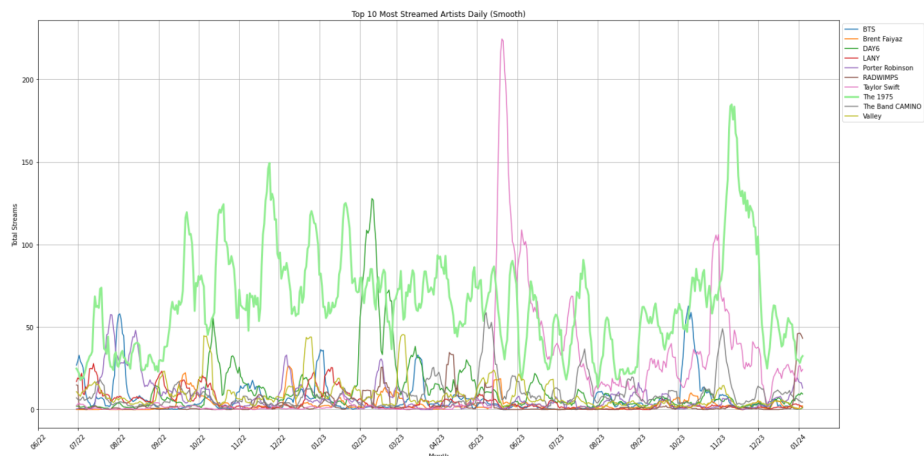


Figure 3: Top artists over time (smooth)

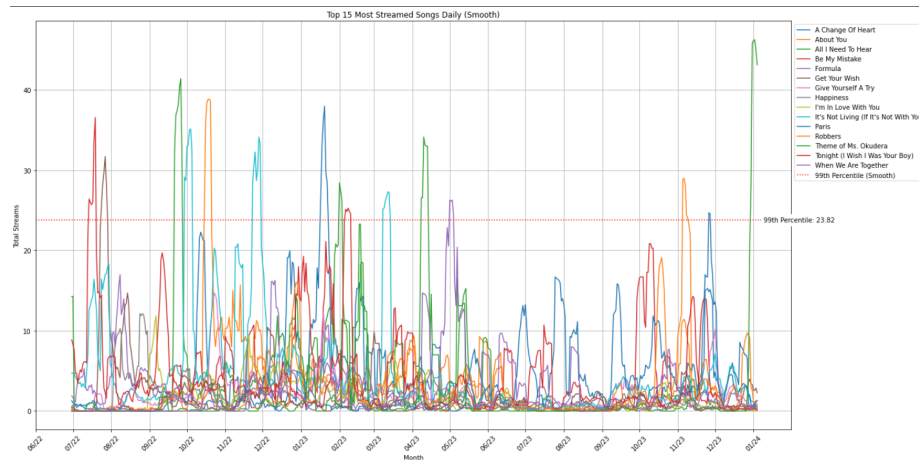


Figure 5: Top tracks over time (raw)

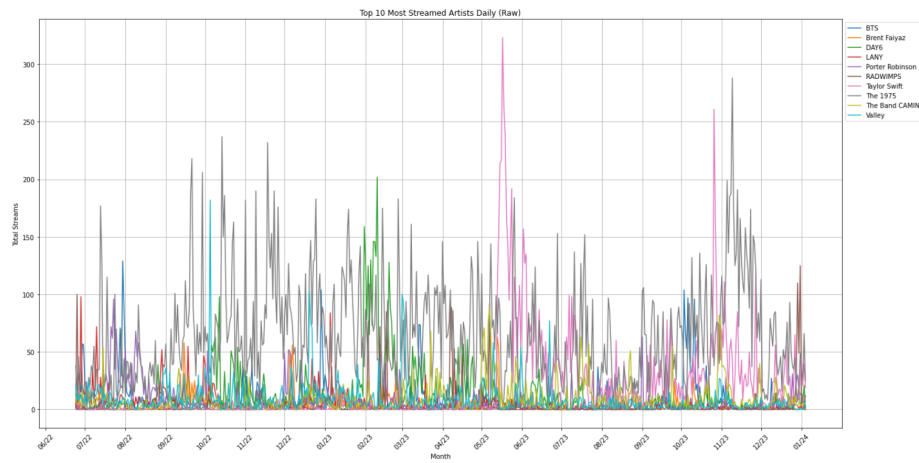


Figure 4: Top artists over time (raw)

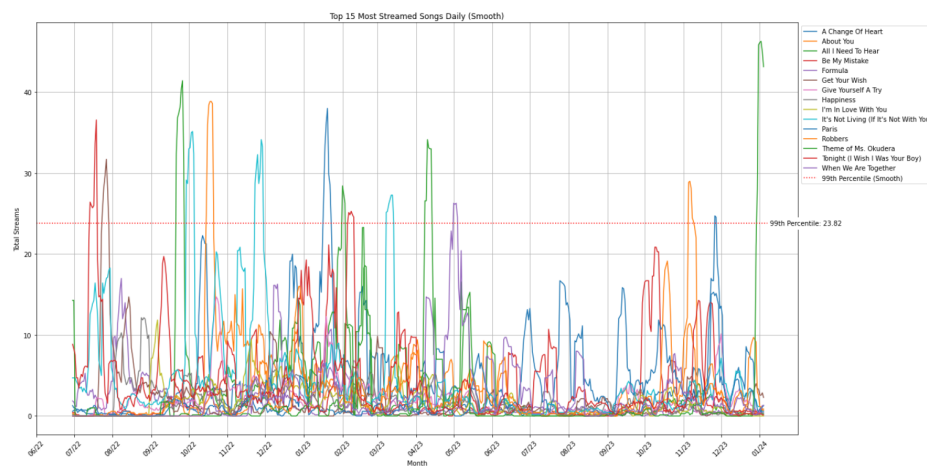


Figure 6: Top tracks over time (smooth)

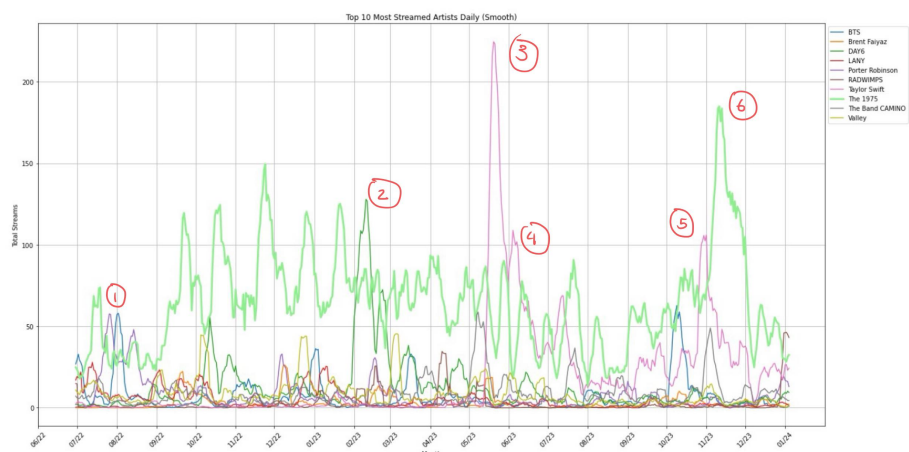


Figure 7: Labeled Top Artists over time (smooth)

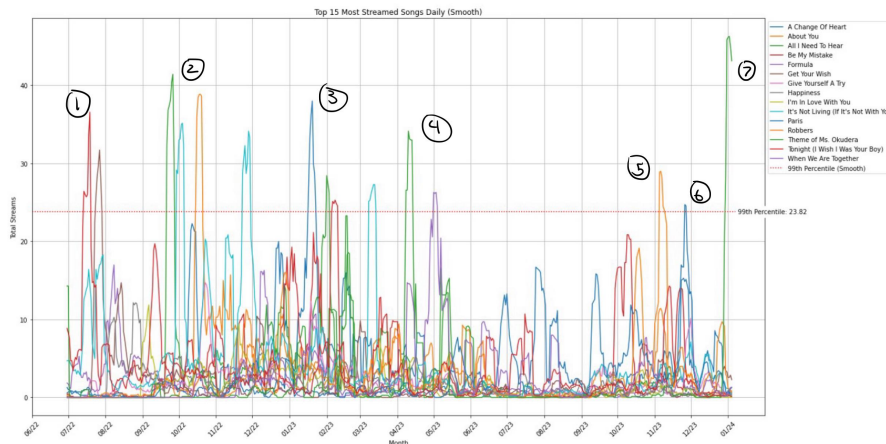


Figure 8: Labeled Top tracks over time (smooth)

Figure 3 shows the top artists listened to over time from June 2022 to January 2024. Evident from the many spikes in the graph, there is much fluctuation over time concerning which artist the user listens to with higher frequency, often changing from day to day. To make the visualization easier to process, a smoothed version is showed in figure 4 using a 7-day moving average.

Figure 4 shows a smoothed version of the top artists listened to over time from June 2022 to January 2024. Highlighted in light green is the moving average of streams per day for The 1975, the top artist listened to by the user. The 1975 is by far the user's most listened artist, so perhaps the time periods where another artist's moving average surpasses that of The 1975 represents an important period of turbulence for the user.

Figure 5 shows the top tracks listened to over time from June 2022 to January 2024. Even more so than the line chart for artists, the visualization is difficult to interpret.

Figure 6 shows a smoothed version of the top tracks listened to over time from June 2022 to January 2024. Furthermore, a line has been drawn at the 99th percentile of entries on the line chart, representing the top 1 percent of song streams per day among the top 15 tracks. One can safely assume that values above that line are very abnormal based on the user's listening history, indicating that these inflection points may be points of interest as the project progresses.

In the following section, major inflection points will be analyzed and traced

back to the user's life.

Figure 7 is similar to figure 4, but it also contains labels for inflection points of interest for artists listened to.

1. Turbulent time in Summer with lots of life changes
2. Devastating hand injury
3. Vacation in Europe
4. Moving to New York
5. Release of 1989 Taylor's Version
6. Attending concert for The 1975

Figure 8 is similar to figure 6, but it also contains labels for inflection points of interest for tracks listened to.

1. Summer 2022
2. Fall 2022 Semester
3. New Year's, start of Spring 2023 Semester
4. 2022-2023 School year ending
5. Attending The 1975 concert
6. Thanksgiving break
7. Winter break

Interestingly, many of these behavioral changes occur during periods of change, such as breaks from school or the start of new semesters. This may be an indication that when one's environment changes, their listening behavior may change as well.

Many changes can also be seen in both the track and artist visualizations. Furthermore, a large gap in inflection points can be observed between May 2023 and November 2023 on Figure 8. This may be evident of a calmer lifestyle that includes less big changes.

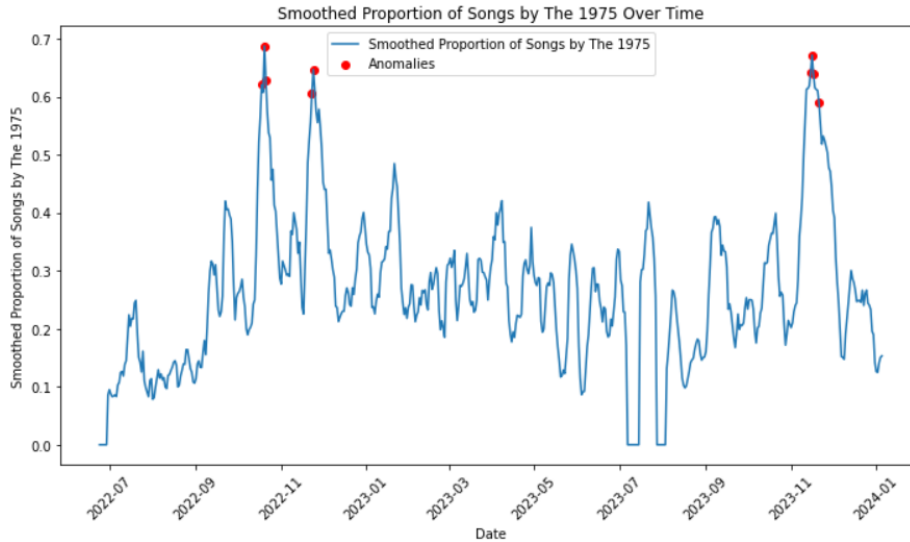


Figure 9: Abnormality Detection High

4 Isolation Forest High Anomaly Detection

This next section discusses the fit of an Isolation Forest to the data to detect anomalies in listening behavior. An isolation forest works by isolating anomalies (outliers) in a dataset by randomly partitioning the data into subsets, making it more efficient in detecting anomalies compared to traditional methods. It identifies anomalies as points that are easier to isolate due to their uniqueness, making them stand out from the majority of normal data points.

Figure 9 illustrates the output from using Isolation Forest anomaly detection to find abnormalities in our smoothed averages. The y-axis represents the proportion of The 1975 tracks of total streams over a 7-day period. The 'red' abnormal points represent statistical anomalies as detected by the random forest using a contamination rate of 0.05.

As we can see, the three main anomaly areas are October 2022, November 2022, and November 2023. This means that during these periods, I was listening to an abnormally high amount of The 1975 as a proportion of my total listens. Because the October 2022 and November 2022 period are adjacent, I believe that it's safe to group them as one "event."

I believe that these line up with life events that I was experiencing because in both November 2022 and November 2023 I watched The 1975 live in concert. This is evidence to suggest that when I see an artist live in concert, I have a tendency to increase how much I listen to them in the near future. For a more

extensive analysis, I would like to observe how each concert I goes to affects my listening to each artist, beyond just my top artist.

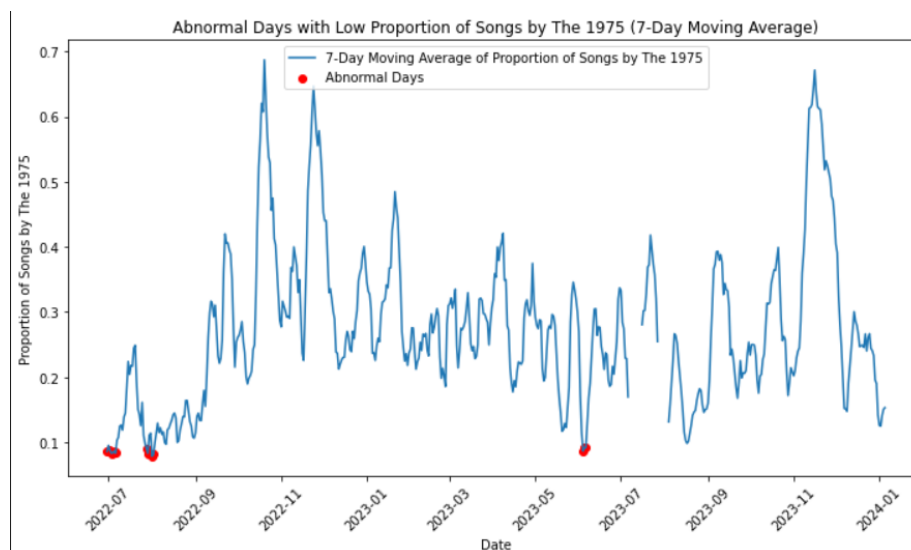


Figure 10: Caption

5 Z-Score Low Anomaly Detection

Figure 10 (featured above) expresses the output of z-score anomaly detection in Python using our time series data. The threshold for anomaly was set using a z-score of -1.5. This means that proportion values whose z-score falls under -1.5 were marked as anomalies and denoted with a red dot in our visualization.

As you can see, the abnormal lows in our data are in July 2022, August 2022, and May 2023. Since July 2022 and August 2022 are adjacent, I also believe that it suffices to group them as the same "event" in my timeline.

This means that the two outlier instances in my listening are late summer 2022 and early summer of 2023. This might suggest that summer is one of the most regularly turbulent periods of my life where I experience the most change. In summer 2022, I was doing a lot of new things and sorting out how I wanted to approach the upcoming school year. During the May 2023 dip, I was traveling to London, Paris, Italy, and eventually settling in New York following the conclusion of the May 2023 semester.

This data suggests that I experience the most change in listening when there is a lot of change in my life, whether it be mentally (Summer 2022) or physically/geographically (Summer 2023). I would love to spend more time further

exploring whether or not the same would be true in the data of other users.

6 Time Series Decomposition

Time series decomposition is a statistical method used to break down a time series into its constituent components, typically trend, seasonality, and noise, enabling the identification and analysis of underlying patterns and trends within the data.

By applying such decomposition to my listening data of The 1975, we hope to find trends of seasonality in the listening data.

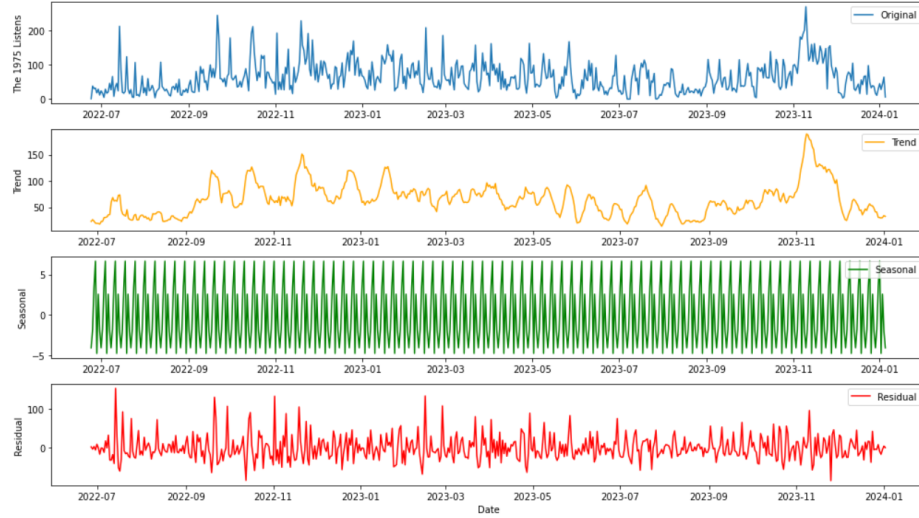


Figure 11: Time Series Decomposition

Based on Figure 11 shown above, there is no evident trend in seasonality. It ranges from -5 to 5, which is a too small when compared to the bulk of daily listens to be considered relevant. Interestingly, the trend graph roughly follows the overall graph, further expressing the lack of seasonality evident from the decomposition. Lastly, the residual graph shows heightened variance in the first half of the dataset, with "spikes" in residuals occurring every few months. These high variance periods could be indicative of times where listening behavior is shifting.

7 Conclusion

The most important findings in this data indicate that my listening trends shift the most during summer, likely to both physical and non-physical aspects of my

life changing drastically during these times. Examples of these changes include the start of new semesters and traveling to new cities in the world. In other words, shifts in my music taste reflect major changes in my own life.

Despite these findings, this project has even more potential to grow in the future. Suggestions for future analysis include:

1. Formalizing code into a package or application which will allow anyone to easily visualize their exact listening trends.
2. Using sentiment analysis to analyze song titles or lyrics. Lyrics can be scraped using Genius.com's API. Looking at common words or phrases could be more indicative of "emotion-based" listening trends at different times of day or eras.
3. Collecting more data over time. Since there has only been two summers, it's difficult to generalize trends occurring annually each summer. Collecting more data over more years would allow a more seasoned analysis of how listening trends vary throughout the year.

In conclusion, this project reveals a potential relationship between Spotify listening behavior and life events while the proposed future analyses, including formalizing code, sentiment analysis on lyrics, and longitudinal data collection, promise further insights into this dynamic relationship. These avenues offer valuable opportunities for furthering our understanding of listener behavior alongside possible standardization into something that anyone can use.

8 References

1. GeeksforGeeks. (n.d.). What is Isolation Forest? Retrieved from <https://www.geeksforgeeks.org/what-is-isolation-forest/>
2. ChatGPT. (n.d.). Retrieved from <https://chatgpt.com/>
3. Spotify. (n.d.). Data Rights and Privacy Settings. Retrieved from <https://support.spotify.com/us/article/data-rights-and-privacy-settings/>

9 Appendix

```
1 # %%
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import warnings
6 warnings.filterwarnings("ignore")
7
8 # %%
9 df_2022 = pd.read_csv("Streams_0722_0723.csv")
```

```

10
11 # %%
12 df = pd.read_json("StreamingHistory0.json")
13 df1 = pd.read_json("StreamingHistory1.json")
14 df2 = pd.read_json("StreamingHistory2.json")
15 df3 = pd.read_json("StreamingHistory3.json")
16 df4 = pd.read_json("StreamingHistory4.json")
17 df5 = pd.read_json("StreamingHistory5.json")
18 df6 = pd.read_json("StreamingHistory6.json")
19 df7 = pd.read_json("StreamingHistory7.json")
20 df8 = pd.read_json("StreamingHistory8.json")
21
22 # %%
23 dflist = [df_2022, df, df1, df2, df3, df4, df5, df6, df7, df8]
24 master_stream_df = pd.concat(dflist)
25 if 'Unnamed: 0' in master_stream_df.columns:
26     master_stream_df.drop(columns=['Unnamed: 0'], inplace=True)
27 master_stream_df.head()
28
29
30 # %% [markdown]
31 # Next, we will check the length of the df
32
33 # %%
34 print("Length of combined df including duplicates: {}".format(len(
    master_stream_df)))
35 master_stream_df.drop_duplicates(inplace=True)
36 print("Length of combined df without duplicates: {}".format(len(
    master_stream_df)))
37
38
39 # %%
40 master_stream_df.head()
41
42 # %%
43 print(master_stream_df.iloc[0,0])
44 print(master_stream_df.iloc[len(master_stream_df)-1],0)
45
46 # %%
47 from pytz import timezone
48
49 # Assuming master_stream_df contains your dataframe
50 master_stream_df['endTime'] = pd.to_datetime(master_stream_df['
    endTime'])
51
52 # Define timezones
53 utc_timezone = timezone('UTC')
54 central_timezone = timezone('US/Central')
55
56 # Localize the endTime column to UTC timezone
57 master_stream_df['endTime'] = master_stream_df['endTime'].dt.
    tz_localize(utc_timezone)
58
59 # Convert times from UTC to Central Time and format to retain
    seconds count
60 master_stream_df['endTime'] = master_stream_df['endTime'].dt.
    tz_convert(central_timezone).dt.strftime('%Y-%m-%d %H:%M:%S')

```

```

61
62 # Filter the dataframe to keep rows with endTime in 2023 and reset
    the index
63 #filtered_df = master_stream_df[master_stream_df['endTime'].str.
    startswith('2023')].reset_index(drop=True)
64
65 filtered_df = master_stream_df
66
67 # %%
68 len(filtered_df)
69
70 # %%
71 def getTopTracks(df):
72     playCounts = df['trackName'].value_counts()
73     playCounts_df = pd.DataFrame({'Track Name': playCounts.index, '
        Streams': playCounts.values})
74     return playCounts_df
75 def getTopArtists(df):
76     artistCounts = df['artistName'].value_counts()
77     artistCounts_df = pd.DataFrame({'Artist Name' : artistCounts.
        index, 'Streams': artistCounts.values})
78     return artistCounts_df
79
80 # %%
81 top_filtered = getTopTracks(filtered_df)
82 top_filtered.head(15)
83
84 # %%
85 unknownFilter = filtered_df['trackName'] != 'Unknown Track'
86 no_unknown = filtered_df[unknownFilter]
87 print(len(no_unknown))
88
89 # %%
90 filtered_unknown = getTopTracks(no_unknown)
91 filtered_unknown.head(15)
92
93 # %%
94 filtered_unknown = filtered_unknown.head(15)
95 plt.figure(figsize = (15,10))
96 bars = plt.barh(filtered_unknown['Track Name'], filtered_unknown['
    Streams'],color = 'skyblue')
97 for bar in bars:
98     plt.text(bar.get_width(), bar.get_y() + bar.get_height() / 2, f
        '{bar.get_width():.0f}', va='center')
99 plt.gca().invert_yaxis()
100 plt.xlabel('Streams')
101 plt.ylabel('Track Name')
102 plt.title('Top 15 Tracks by Streams')
103 plt.show()
104
105 # %%
106 artist_counts = getTopArtists(no_unknown)
107 artist_counts.head(20)
108
109 # %%
110 top_artists = artist_counts.head(10)
111 plt.figure(figsize = (10,10))

```

```

112 bars = plt.barh(top_artists['Artist Name'], top_artists['Streams'],
113                 color = 'skyblue')
113 for bar in bars:
114     plt.text(bar.get_width(), bar.get_y() + bar.get_height() / 2, f
115             '{bar.get_width():.0f}', va='center')
115 plt.gca().invert_yaxis()
116 plt.xlabel('Streams')
117 plt.ylabel('Artist Name')
118 plt.title('Top 10 Artists by Streams')
119 plt.show()
120
121 # %%
122 no_unknown.head(10)
123
124 # %%
125 import matplotlib.dates as mdates
126
127 def plot_top_artists_smooth(df, period, n_artists, plot_type,
128                             display_agg=True):
129     # Convert 'endTime' to datetime format
130     df.loc[:, 'endTime'] = pd.to_datetime(df['endTime'])
131
132     if period == 'Day':
133         df.loc[:, 'DaysElapsed'] = (df['endTime'] - df['endTime'].
134                                     min()).dt.days
135         grouped_df = df.groupby(['artistName', 'DaysElapsed']).size()
136         .reset_index(name='Counts')
137     elif period == 'Week':
138         df.loc[:, 'Week'] = df['endTime'].dt.isocalendar().week
139         df.loc[:, 'YearWeek'] = df['endTime'].dt.strftime('%Y-%W')
140         grouped_df = df.groupby(['artistName', 'YearWeek']).size()
141         .reset_index(name='Counts')
142     elif period == 'Month':
143         df.loc[:, 'Month'] = df['endTime'].dt.to_period('M')
144         grouped_df = df.groupby(['artistName', 'Month']).size()
145         .reset_index(name='Counts')
146
147     top_artists = grouped_df.groupby('artistName')['Counts'].sum().
148     nlargest(n_artists).index
149     top_artists_df = grouped_df[grouped_df['artistName'].isin(
150         top_artists)]
151
152     if plot_type == 'smooth':
153         pivot_df = top_artists_df.pivot_table(index=['DaysElapsed',
154             ], columns='artistName', values='Counts', fill_value=0)
155         plt.figure(figsize=(20, 10))
156         for column in pivot_df.columns:
157             pivot_df[column] = pivot_df[column].rolling(window=7).
158             mean()
159             if column == top_artists[0]: # Assuming the first
160                 artist in the top list is the top artist of the
161                 year
162                 plt.plot(pivot_df.index, pivot_df[column], label=
163                     column, linewidth=3, color='lightgreen') #
164                     Adjust linewidth and color as needed
165             else:
166                 plt.plot(pivot_df.index, pivot_df[column], label=

```

```

        column)

154
155     # Calculate aggregate streams across top artists if
        display_agg is True
156     if display_agg:
157         aggregate_streams = pivot_df.sum(axis=1)
158         if 'smooth' in plot_type:
159             aggregate_streams = aggregate_streams.rolling(
                window=7).mean()
160         plt.plot(pivot_df.index, aggregate_streams, label='
            Aggregate Streams', linestyle='--', color='blue')
            # Adjust linestyle and color as needed

161
162     elif plot_type == 'raw':
163         pivot_df = top_artists_df.pivot_table(index=['DaysElapsed'
164             ], columns='artistName', values='Counts', fill_value=0)
165         plt.figure(figsize=(20, 10))
166         for column in pivot_df.columns:
167             plt.plot(pivot_df.index, pivot_df[column], label=column
168                 )

169     # Calculate aggregate streams across top artists if
        display_agg is True
170     if display_agg:
171         aggregate_streams = pivot_df.sum(axis=1)
172         plt.plot(pivot_df.index, aggregate_streams, label='
            Aggregate Streams', linestyle='--', color='blue')
            # Adjust linestyle and color as needed

173
174     plt.xlabel('Month')
175     plt.ylabel('Total Streams')
176     plt.title('Top {} Most Streamed Artists Daily ({}).format(
        n_artists, plot_type.capitalize())

177
178     # Move the legend outside the plot area
179     plt.legend(loc='upper left', bbox_to_anchor=(1, 1))
180
181     # Add month labels
182     first_day_of_month = pd.date_range(start='2022-06-01', end='
        2024-01-01', freq='MS')
183     plt.xticks([(month_start - pd.to_datetime('2022-06-24')).days
184         for month_start in first_day_of_month], [month_start.
        strftime('%m/%y') for month_start in first_day_of_month],
        rotation=45)

185
186     plt.grid(True)
187     plt.tight_layout()
188     plt.show()
189
190     # %%
191     def plot_top_songs_smooth(df, period, n_songs, plot_type,
        display_agg=True, threshold=None):
192         # Convert 'endTime' to datetime format
193         df['endTime'] = pd.to_datetime(df['endTime'])
194
195         if period == 'Day':
196             df['DaysElapsed'] = (df['endTime'] - df['endTime'].min()).

```

```

        dt.days
195     grouped_df = df.groupby(['DaysElapsed', 'trackName']).size
        ().reset_index(name='Streams')
196 elif period == 'Week':
197     df['Week'] = df['endTime'].dt.isocalendar().week
198     df['YearWeek'] = df['endTime'].dt.strftime('%Y-%W')
199     grouped_df = df.groupby(['Week', 'trackName']).size().
        reset_index(name='Streams')
200 elif period == 'Month':
201     df['Month'] = df['endTime'].dt.to_period('M')
202     grouped_df = df.groupby(['Month', 'trackName']).size().
        reset_index(name='Streams')
203
204 top_songs = grouped_df.groupby('trackName')['Streams'].sum().
        nlargest(n_songs).index
205 top_songs_df = grouped_df[grouped_df['trackName'].isin(
        top_songs)]
206
207 if plot_type == 'smooth':
208     pivot_df = top_songs_df.pivot_table(index=['DaysElapsed'],
        columns='trackName', values='Streams', fill_value=0)
209     plt.figure(figsize=(20, 10))
210     for column in pivot_df.columns:
211         pivot_df[column] = pivot_df[column].rolling(window=7).
            mean()
212         plt.plot(pivot_df.index, pivot_df[column], label=column
            )
213
214     # Calculate the 99th percentile of song occurrences per day
        after applying rolling mean
215     all_values = np.nan_to_num(pivot_df.values.flatten())
216     percentile_99 = np.percentile(all_values, 99)
217     plt.axhline(y=percentile_99, color='r', linestyle=':',
        label='99th Percentile (Smooth)')
218
219     # Add label for the line outside the plot area
220     if np.isfinite(pivot_df.index[-1]) and np.isfinite(
        percentile_99):
221         plt.text(pivot_df.index[-1] + 5, percentile_99, f'99th
            Percentile: {percentile_99:.2f}', ha='left', va='
            center', backgroundcolor='w')
222
223 elif plot_type == 'raw':
224     pivot_df = top_songs_df.pivot_table(index=['DaysElapsed'],
        columns='trackName', values='Streams', fill_value=0)
225     plt.figure(figsize=(20, 10))
226     for column in pivot_df.columns:
227         plt.plot(pivot_df.index, pivot_df[column], label=column
            )
228
229     # Calculate the 99th percentile of song occurrences per day
230     all_values = np.nan_to_num(top_songs_df.groupby('
        DaysElapsed')['Streams'].sum().values)
231     percentile_99 = np.percentile(all_values, 99)
232     plt.axhline(y=percentile_99, color='r', linestyle=':',
        label='99th Percentile (Raw)')
233

```



```

234     # Add label for the line outside the plot area
235     if np.isfinite(pivot_df.index[-1]) and np.isfinite(
        percentile_99):
236         plt.text(pivot_df.index[-1] + 5, percentile_99, f'99th
            Percentile: {percentile_99:.2f}', ha='left', va='
            center', backgroundcolor='w')
237
238     plt.xlabel('Month')
239     plt.ylabel('Total Streams')
240     plt.title('Top {} Most Streamed Songs Daily ({}).format(
        n_songs, plot_type.capitalize())
241
242     # Move the legend outside the plot area
243     plt.legend(loc='upper left', bbox_to_anchor=(1, 1))
244
245     # Add month labels
246     first_day_of_month = pd.date_range(start='2022-06-01', end='
        2024-01-01', freq='MS')
247     plt.xticks([(month_start - pd.to_datetime('2022-06-24')).days
        for month_start in first_day_of_month], [month_start.
        strftime('%m/%y') for month_start in first_day_of_month],
        rotation=45)
248
249     plt.grid(True)
250
251     # Draw a horizontal line at the specified threshold if provided
252     if threshold is not None:
253         plt.axhline(y=threshold, color='g', linestyle='--', label='
            Threshold')
254
255     plt.tight_layout()
256     plt.show()
257
258     # %%
259     plot_top_artists_smooth(no_unknown, 'Day', 10, 'smooth', False)
260
261     # %%
262     plot_top_artists_smooth(no_unknown, 'Day', 10, 'smooth', True)
263
264     # %%
265     plot_top_songs_smooth(no_unknown, 'Day', 15, 'smooth', True)
266
267     # %%
268     plot_top_songs_smooth(no_unknown, 'Day', 15, 'smooth', True)
269
270     # %%
271     plot_top_songs_smooth(no_unknown, 'Day', 15, 'raw', False)
272
273     # %%
274     from sklearn.preprocessing import LabelEncoder
275     from sklearn.ensemble import IsolationForest
276     # Convert 'endTime' to datetime format
277     no_unknown['endTime'] = pd.to_datetime(no_unknown['endTime'])
278     # Encode 'artistName' using Label Encoding
279     label_encoder = LabelEncoder()
280     no_unknown['artistCode'] = label_encoder.fit_transform(no_unknown['
        artistName'])

```

```

281 # Calculate the daily counts of all songs and "The 1975" songs
282 daily_counts = no_unknown.groupby(no_unknown['endTime'].dt.date)['trackName'].count()
283 the_1975_songs = no_unknown[no_unknown['artistName'] == 'The 1975']
284 the_1975_daily_counts = the_1975_songs.groupby(the_1975_songs['endTime'].dt.date)['trackName'].count()
285 # Calculate the proportion of songs by "The 1975" each day
286 proportion_by_day = the_1975_daily_counts / daily_counts
287 # Initialize Isolation Forest model
288 model = IsolationForest(contamination=0.05) # Adjust contamination
        based on expected anomaly rate
289 # Fit the model to your data (using the raw proportion)
290 proportion_by_day = proportion_by_day.fillna(0) # Replace NaN
        values with 0 for the initial period
291 model.fit(proportion_by_day.values.reshape(-1, 1))
292 # Predict anomalies (-1 indicates anomalies, 1 indicates normal
        data points)
293 anomaly_predictions = model.predict(proportion_by_day.values.
        reshape(-1, 1))
294 # Filter the dataframe to keep only anomalies
295 anomalies = proportion_by_day[anomaly_predictions == -1]
296 # Visualize anomalies with raw proportion of songs by "The 1975" as
        y-axis
297 plt.figure(figsize=(10, 6))
298 plt.plot(proportion_by_day.index, proportion_by_day.values, label='Proportion of Songs by The 1975')
299 plt.scatter(anomalies.index, anomalies.values, color='red', label='Anomalies')
300 plt.xlabel('Date')
301 plt.ylabel('Proportion of Songs by The 1975')
302 plt.title('Proportion of Songs by The 1975 Over Time')
303 plt.legend()
304 plt.xticks(rotation=45) # Rotate x-axis labels for better
        readability
305 plt.tight_layout()
306 plt.show()
307 # Print or further analyze the anomalies dataframe
308 print("Number of anomalies detected:", len(anomalies))
309 from sklearn.preprocessing import StandardScaler
310 # Calculate daily counts of all songs and "The 1975" songs
311 daily_counts = no_unknown.groupby(no_unknown['endTime'].dt.date)['trackName'].count()
312 the_1975_songs = no_unknown[no_unknown['artistName'] == 'The 1975']
313 the_1975_daily_counts = the_1975_songs.groupby(the_1975_songs['endTime'].dt.date)['trackName'].count()
314 # Calculate the proportion of songs by "The 1975" for each day
315 proportion_by_day = the_1975_daily_counts / daily_counts
316 # Calculate the 7-day moving average of the proportion
317 moving_average = proportion_by_day.rolling(window=7).mean()
318 # Calculate z-scores for each day's proportion
319 scaler = StandardScaler()
320 z_scores = scaler.fit_transform(moving_average.values.reshape(-1,
        1))
321 # Set thresholds for abnormality detection (e.g., z-score below
        -1.5 and above 1.5)
322 lower_threshold = -1.5
323 upper_threshold = 1.5

```

```

324 # Find abnormal days based on thresholds
325 abnormal_low_days = moving_average[z_scores.flatten() <
    lower_threshold]
326 abnormal_high_days = moving_average[z_scores.flatten() >
    upper_threshold]
327 plt.figure(figsize=(10, 6))
328 plt.plot(moving_average.index, moving_average.values, label='7-Day
    Moving Average of Proportion of Songs by The 1975')
329 plt.scatter(abnormal_low_days.index, abnormal_low_days.values,
    color='red', label='Abnormal Low Days')
330 plt.scatter(abnormal_high_days.index, abnormal_high_days.values,
    color='red', label='Abnormal High Days')
331 plt.xlabel('Date')
332 plt.ylabel('Proportion of Songs by The 1975')
333 plt.title('Abnormal Days with Low and High Proportion of Songs by
    The 1975 (7-Day Moving Average)')
334 plt.legend()
335 plt.xticks(rotation=45)
336 plt.tight_layout()
337 plt.show()
338 from statsmodels.tsa.seasonal import seasonal_decompose
339 no_unknown['endTime'] = pd.to_datetime(no_unknown['endTime'])
340 # Filter rows where 'artistName' is 'The 1975'
341 the_1975_data = no_unknown[no_unknown['artistName'] == 'The 1975']
342 # Group by date and count occurrences, setting the frequency to 'D'
    (daily)
343 the_1975_daily_counts = the_1975_data.groupby(pd.Grouper(key='
    endTime', freq='D')).size()
344 # Perform seasonal decomposition
345 decomposition = seasonal_decompose(the_1975_daily_counts, model='
    additive')
346 # Plotting the decomposition with a wider figure size
347 plt.figure(figsize=(14, 8)) # Adjust the width and height as needed
348 # Original time series
349 plt.subplot(4, 1, 1)
350 plt.plot(the_1975_daily_counts, label='Original')
351 plt.legend()
352 plt.ylabel('The 1975 Listens')
353 # Trend component
354 plt.subplot(4, 1, 2)
355 plt.plot(decomposition.trend, label='Trend', color='orange')
356 plt.legend()
357 plt.ylabel('Trend')
358 # Seasonal component
359 plt.subplot(4, 1, 3)
360 plt.plot(decomposition.seasonal, label='Seasonal', color='green')
361 plt.legend()
362 plt.ylabel('Seasonal')
363 # Residual component
364 plt.subplot(4, 1, 4)
365 plt.plot(decomposition.resid, label='Residual', color='red')
366 plt.legend()
367 plt.ylabel('Residual')
368 plt.xlabel('Date')
369 plt.tight_layout()
370 plt.show()
371 no_unknown['endTime'] = pd.to_datetime(no_unknown['endTime'])

```

```

372 # Encode 'artistName' using Label Encoding
373 label_encoder = LabelEncoder()
374 no_unknown['artistCode'] = label_encoder.fit_transform(no_unknown['
    artistName'])
375 # Calculate the daily counts of all songs and "The 1975" songs
376 daily_counts = no_unknown.groupby(no_unknown['endTime'].dt.date)['
    trackName'].count()
377 the_1975_songs = no_unknown[no_unknown['artistName'] == 'The 1975']
378 the_1975_daily_counts = the_1975_songs.groupby(the_1975_songs['
    endTime'].dt.date)['trackName'].count()
379 # Calculate the proportion of songs by "The 1975" each day
380 proportion_by_day = the_1975_daily_counts / daily_counts
381 # Smooth the proportion using a 3-day moving average
382 smoothed_proportion = proportion_by_day.rolling(window=5).mean()
383 # Initialize Isolation Forest model
384 model = IsolationForest(contamination=0.05) # Adjust contamination
    based on expected anomaly rate
385 # Fit the model to your data (using the smoothed proportion)
386 smoothed_proportion = smoothed_proportion.fillna(0) # Replace NaN
    values with 0 for the initial period
387 model.fit(smoothed_proportion.values.reshape(-1, 1))
388 # Predict anomalies (-1 indicates anomalies, 1 indicates normal
    data points)
389 anomaly_predictions = model.predict(smoothed_proportion.values.
    reshape(-1, 1))
390 # Filter the dataframe to keep only anomalies
391 anomalies = smoothed_proportion[anomaly_predictions == -1]
392 # Visualize anomalies with smoothed proportion of songs by "The
    1975" as y-axis
393 plt.figure(figsize=(10, 6))
394 plt.plot(smoothed_proportion.index, smoothed_proportion.values,
    label='Smoothed Proportion of Songs by The 1975')
395 plt.scatter(anomalies.index, anomalies.values, color='red', label='
    Anomalies')
396 plt.xlabel('Date')
397 plt.ylabel('Smoothed Proportion of Songs by The 1975')
398 plt.title('Smoothed Proportion of Songs by The 1975 Over Time')
399 plt.legend()
400 plt.xticks(rotation=45) # Rotate x-axis labels for better
    readability
401 plt.tight_layout()
402 plt.show()
403 # Print or further analyze the anomalies dataframe
404 print("Number of anomalies detected:", len(anomalies))

```

482FinalCode.py