

## Table of contents

## PART 1: DATA CLEANING

Set Up Dataframe

Convert Data Types

Check Missing Values

Feature Transformation

Merge Release Day, Month, and Year to Release Date

Create Average Monthly Stream Value

Convert Artist Count to Categorical Value

## PART 2: DESCRIPTIVE ANALYTICS, DATA VISUALIZATION AND TAKEAWAYS

## Total Streams &amp; Average Monthly Streams

Log Transformation of Total Streams

Remove Outliers

## Song Features

Numerical Song Features

Artist Count

Key &amp; Mode

Pivot Table of Artistic Features

## Playlists

Spotify Playlist

Apple Playlist

## Release Month

## Release Year

## PART 3: DESCRIPTIVE ANALYTICS, DATA VISUALIZATION AND TAKEAWAYS

## ANOVA

Artist Count vs Streams

ANOVA Post-Hoc Test

Release Month vs Streams

ANOVA Post-Hoc Test

Key vs Streams

NING

ie

pes

'values

motion

ise Day, Month, and

ase Date

ige Monthly Stream

st Count to

Value

## DATA VISUALIZATION AND TAKEAWAYS

## Average Monthly

Information of Total

liers

ong Features

## of Artistic Features

ist

it

pes

'values

motion

ise Day Month, and

ase Date

ige Monthly Stream

st Count to

Value

## DATA VISUALIZATION AND TAKEAWAYS

vs Streams

Post-Hoc Test

th vs Streams

Post-Hoc Test

ms

NING

ie

pes

'values

motion

ise Day Month, and

ase Date

ige Monthly Stream

st Count to

Value

## DATA VISUALIZATION AND TAKEAWAYS

## Average Monthly

Information of Total

liers

## + Code + Text

Group Number: 35

Student ID's: 1006031820, 1003556831, 1009692110

Student Names: Berke Derin Berkay, Paris Wang, Darenne Christabel

Instructor: Maher Elshakankiri

Course code: INF1340

Course name: Programming for Data Science

University of Toronto

Final Project

Please refer to the README file &amp; API documentation for information regarding the program.

## PART 1: DATA CLEANING

## Set Up Dataframe

```
[1] import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
import statsmodels.api as sm
from sklearn.preprocessing import LabelEncoder
from statsmodels.formula.api import ols
```

```
[2] from google.colab import drive
drive.mount('/content/drive') # reading the data from google drive
df = pd.read_csv('/content/drive/MyDrive/spotify-2023.csv', encoding='iso-8859-1')
df.head()
```

Mounted at /content/drive

|   | track_name          | artist(s)_name  | artist_count     | released_year | released_month | released_day | in_spotify_playlists | in_spotify_charts | streams | in_apple_playlists | ... bpm | key | mode | danceability_% | valence_% | energy_% | acousticness_% |
|---|---------------------|-----------------|------------------|---------------|----------------|--------------|----------------------|-------------------|---------|--------------------|---------|-----|------|----------------|-----------|----------|----------------|
| 0 | Sever (feat. Latto) | (Explicit Ver.) | Latto, Jung Kook | 2             | 2023           | 7            | 14                   | 553               | 147     | 141381703          | 43      | ... | D    | Major          | 80        | 89       | 83             |
| 1 |                     | LALA            | Myke Towers      | 1             | 2023           | 3            | 23                   | 1474              | 48      | 133716286          | 48      | ... | C#   | Major          | 71        | 61       | 74             |
| 2 | vampire             |                 | Olivia Rodrigo   | 1             | 2023           | 6            | 30                   | 1397              | 113     | 140003974          | 94      | ... | F    | Major          | 51        | 32       | 53             |
| 3 | Cruel Summer        |                 | Taylor Swift     | 1             | 2019           | 8            | 23                   | 7858              | 100     | 800840817          | 116     | ... | A    | Major          | 55        | 58       | 72             |
| 4 | WHERE SHE GOES      |                 | Bad Bunny        | 1             | 2023           | 5            | 18                   | 3133              | 50      | 303236322          | 84      | ... | A    | Minor          | 65        | 23       | 80             |

5 rows × 24 columns

[3] df.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 953 entries, 0 to 952
Data columns (total 24 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   track_name      953 non-null    object 
 1   artist(s)_name  953 non-null    object 
 2   artist_count     953 non-null    int64  
 3   released_year    953 non-null    int64  
 4   released_month   953 non-null    int64  
 5   released_day     953 non-null    int64  
 6   in_spotify_playlists  953 non-null    int64  
 7   in_spotify_charts  953 non-null    int64  
 8   streams          953 non-null    float64
 9   in_apple_playlists  953 non-null    int64  
 10  in_apple_charts   953 non-null    int64  
 11  in_deezer_playlists  953 non-null    int64  
 12  in_deezer_charts  953 non-null    int64  
 13  in_shazam_charts  953 non-null    float64
 14  bpm              953 non-null    int64  
 15  key              953 non-null    object 
 16  mode              953 non-null    object 
 17  danceability_%   953 non-null    int64  
 18  valence_%        953 non-null    int64  
 19  energy_%         953 non-null    int64  
 20  loudness_%       953 non-null    int64  
 21  instrumentalness_% 953 non-null    int64  
 22  liveness_%       953 non-null    int64  
 23  speechiness_%   953 non-null    int64  
dtypes: float64(2), int64(18), object(4)
memory usage: 178.8+ KB
```

## Check Missing Values

```
[6] # we are checking for na values to see if there are anomalies
empty_cells = df.isna().sum()
print(empty_cells)
```

```
track_name          0
artist(s)_name      0
artist_count         0
released_year        0
released_month       0
released_day         0
in_spotify_playlists 0
in_spotify_charts    0
streams              1
in_apple_playlists   0
in_apple_charts      0
in_deezer_playlists  0
in_deezer_charts     0
in_shazam_charts     50
bpm                  0
key                  95
mode                 0
danceability_%       0
valence_%            0
energy_%             0
acousticness_%       0
instrumentalness_%  0
liveness_%           0
speechiness_%        0
dtype: int64
```

As we can see, streams column has one missing value.

```
[7] # check the row with missing streams
missing_streams_rows = df[df['streams'].isnull()]
missing_streams_rows
```

|     | track_name                          | artist(s)_name    | artist_count | released_year | released_month | released_day | in_spotify_playlists | in_spotify_charts | streams | in_apple_playlists | ... bpm | key | mode  | danceability_% | valence_% | energy_% |
|-----|-------------------------------------|-------------------|--------------|---------------|----------------|--------------|----------------------|-------------------|---------|--------------------|---------|-----|-------|----------------|-----------|----------|
| 574 | Love Grows (Where My Rosemary Goes) | Edison Lighthouse | 1            | 1970          | 1              | 1            | 2877                 | 0                 | Nan     | 16                 | ...     | I   | Major | 53             | 75        |          |

1 rows × 24 columns

```
50% 1.000000 2020.000000 6.500000 14.000000 10236.500000 2.500000 1.199093e+09 168.000000 82.000000 521.500000 0.000000 NaN 122.500000 65.500000
75% 1.000000 2021.000000 9.000000 21.000000 16200.500000 17.000000 1.726432e+09 257.250000 122.000000 960.000000 6.000000 NaN 143.250000 74.000000
max 7.000000 2022.000000 12.000000 31.000000 43899.000000 130.000000 3.703895e+09 672.000000 199.000000 6508.000000 46.000000 NaN 206.000000 86.000000
```

[10] # plot histogram

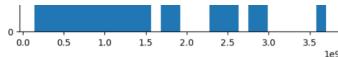
missing\_in\_shazam\_charts\_rows['streams'].plot(kind='hist', bins=30)

plt.title("Histogram of streams for data points missing shazam charts information")

Text(0.5, 1.0, 'Histogram of streams for data points missing shazam charts information')

Histogram of streams for data points missing shazam charts information

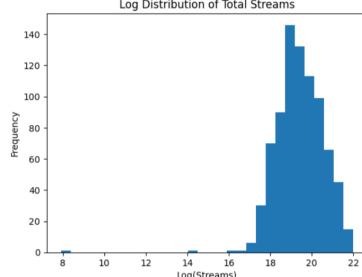




```
plt.title("Log Distribution of Total Streams")
```

```
# Show the plot
```

```
plt.show()
```



We can see from the frequency distribution that the shape of streams is nearly normal.

Next, let's find the outliers using the lognormal distribution of streams

```
Q1 = np.log1p(df['streams']).quantile(0.25)
Q3 = np.log1p(df['streams']).quantile(0.75)
IQR = Q3 - Q1
```

```
# Define the upper and lower bounds to identify outliers
```

```
lower_bound = Q1 - 1.5 * IQR
```

```
upper_bound = Q3 + 1.5 * IQR
```

```
# Identify rows with outliers
```

```
outliers = df[(np.log1p(df['streams']) < lower_bound) | (np.log1p(df['streams']) > upper_bound)]
```

```
outliers
```

| track_name | artist(s)_name              | artist_count                                      | released_year | released_month | released_day | in_spotify_playlists | in_spotify_charts | streams | in_apple_playlists | ... mode     | danceability_% | valence_% | energy_% | acousticness_% | instrumentalness_% |
|------------|-----------------------------|---|---------------|----------------|--------------|----------------------|-------------------|---------|--------------------|--------------|----------------|-----------|----------|----------------|--------------------|
| 123        | Que Vuelvas                 | Carin Leon, Grupo Frontera                        | 2             | 2022           | 12           | 9                    | 783               | 26      | 2762.0             | 21 ... Major | 49             | 78        |          |                |                    |
| 142        | Gol Bolinha, Gol Quadrado 2 | Mc Pedrinho, DJ 900                               | 2             | 2023           | 6            | 1                    | 293               | 8       | 11956641.0         | 5 ... Minor  | 93             | 68        |          |                |                    |
| 393        | Jhoomie Jo Pathaan          | Arijit Singh, Vishal Dadlani, Sukriti Kakar, V... | 4             | 2022           | 12           | 22                   | 138               | 4       | 1365184.0          | 13 ... Major | 82             | 62        |          |                |                    |

3 rows × 26 columns

#### Remove Outliers

```
df = df.drop(outliers.index)
```

We also graphed the relationship between total streams and these song features using a scatterplot.

```
sns.set(style="whitegrid")
```

```
# Set up subplots
```

```
fig, axes = plt.subplots(nrows=1, ncols=len(song_features), figsize=(32, 4))
```

```
# Create side-by-side scatterplots
```

```
for i, feature in enumerate(song_features):
    sns.scatterplot(x=feature, y='streams', data=df, ax=axes[i])
    axes[i].set_title(f"Scatterplot of {feature}")
    axes[i].set_xlabel(feature)
    axes[i].set_ylabel('streams')
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([100, 150, 200])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([20, 40, 60, 80])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

```
axes[i].set_yticks([1e9, 3.0, 1.5, 0.5, 0.0])
```

```
axes[i].set_xticks([0, 25, 50, 75, 100])
```

```
axes[i].set_yticks([0.0, 0.5, 1.0, 1.5, 2.0, 2.5, 3.0, 3.5])
```

```
axes[i].set_xscale('log')
```

```
axes[i].set_yscale('log')
```

</div

```

if Artistic Features
    df = pd.read_csv('data.csv')
    plt.figure(figsize=(10, 6))
    sns.boxplot(x='artist_count', y='streams', data=df)
    plt.title('Box Plot of Song Streams by Number of Artists')
    plt.xlabel('Number of Artists')
    plt.ylabel('Streams')
    plt.show()

```

DATA  
) TAKEAWAYS

vs Streams

Post-Hoc Test

tth vs Streams

Post-Hoc Test

ms

VING

ie

pes

'values

mation

use Day, Month, and

use Date

age Monthly Stream

st Count to

Value

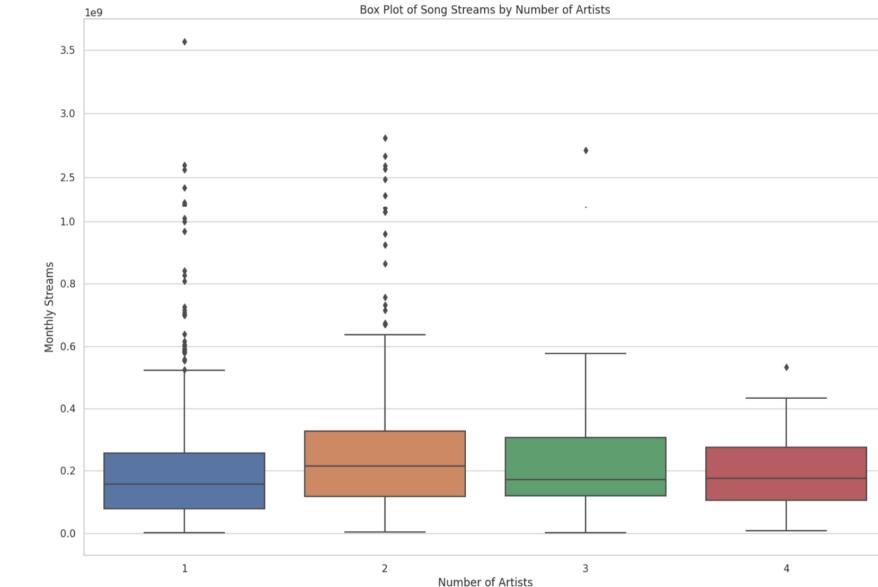
DATA  
) TAKEAWAYS

Average Monthly

mation of Total

liers

ong Features



We can see from the box plots that although average streams for songs look similar across artist counts, the spread looks very different. There appears to be a pattern of the standard deviations of streams as artist counts grew; however, this is most likely due to most songs have less than 3 artists.

Again, there is no obvious patterns for average monthly streams across artist counts, though the average monthly streams for different artist counts look very different. In part 3, we will perform ANOVA to determine if the average are really different across groups given the difference in spread (aka standard deviation).

#### ✓ Key & Mode

```

[32] # we can also visualize the most commonly and rarely used song key
key_counts = df['key'].value_counts()
key_counts

```

| Key | Count |
|-----|-------|
| C#  | 115   |
| G   | 90    |
| F   | 87    |
| G#  | 85    |
| D   | 78    |
| B   | 76    |
| A   | 70    |
| F#  | 69    |
| E   | 59    |
| A#  | 54    |
| D#  | 30    |

DATA  
) TAKEAWAYS

ms

VING

ie

pes

'values

mation

use Day, Month, and

use Date

age Monthly Stream

st Count to

Value

DATA  
) TAKEAWAYS

Average Monthly

mation of Total

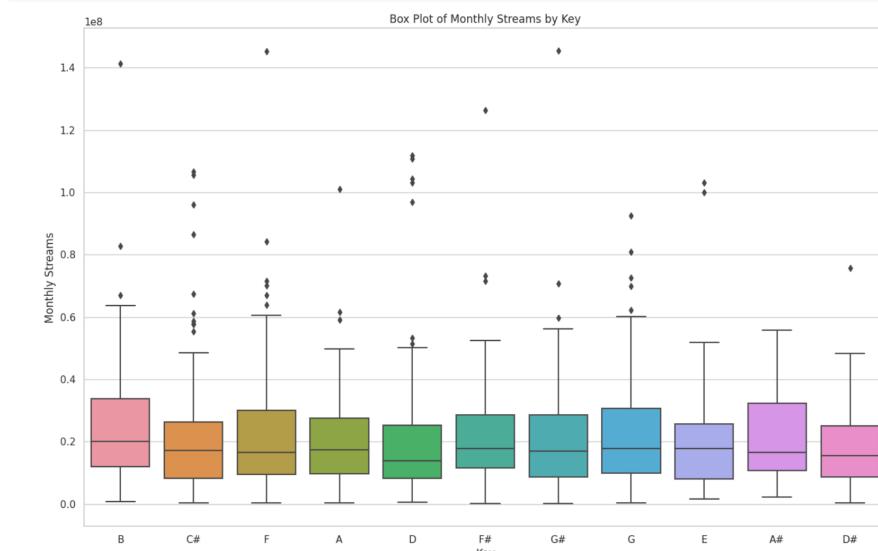
liers

ong Features

if Artistic Features

ist

it



```

[35] # we can do the same for mode
mode_counts = df['mode'].value_counts()
mode_counts

```

| Mode  | Count |
|-------|-------|
| Major | 449   |
| Minor | 364   |

DATA  
) TAKEAWAYS

vs Streams

Post-Hoc Test

tth vs Streams

Post-Hoc Test

ms

VING

ie

pes

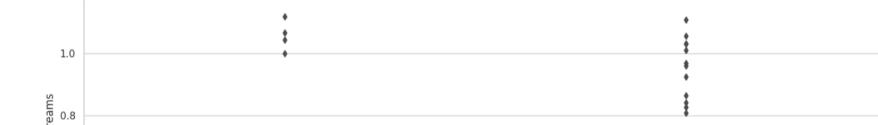
'values

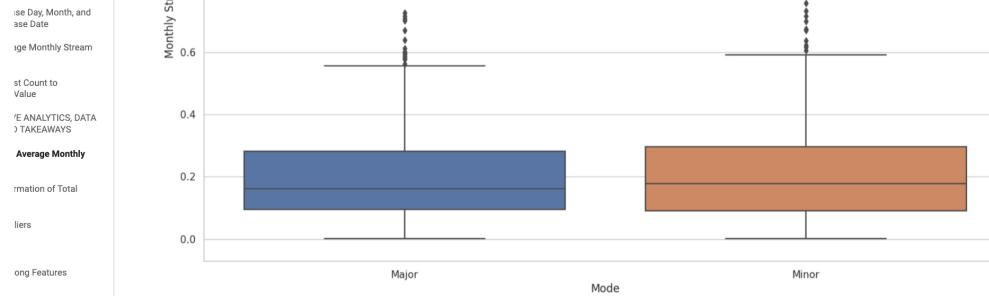
mation

```

[36] plt.figure(figsize=(10, 6))
sns.barplot(x=mode_counts.index, y=mode_counts)
plt.title('Count of Songs by Mode')
plt.xlabel('Song Mode')
plt.ylabel('Count')
plt.xticks(rotation=70)
plt.show()

```





#### ▼ Pivot Table of Artistic Features

After doing some research, we realized that the mode of a song plays a significant role in shaping the mood and vibe of the music.

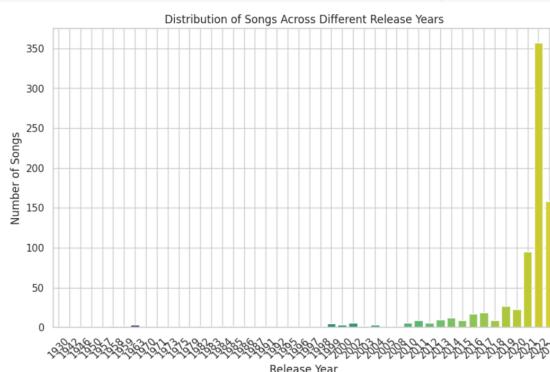
```
[38] pivot_table_result = pd.pivot_table(df,
                                         values=['danceability_%',
                                                  'energy_%', 'acousticness_%',
                                                  'liveness_%', 'speechiness_%'],
                                         index=['mode'],
                                         aggfunc='mean')
pivot_table_result.round(1)
```

| mode  | acousticness_% | danceability_% | energy_% | liveness_% | speechiness_% | valence_% |
|-------|----------------|----------------|----------|------------|---------------|-----------|
| Major | 27.9           | 65.5           | 63.3     | 18.3       | 9.7           | 49.5      |
| Minor | 24.4           | 69.7           | 65.6     | 18.0       | 11.5          | 53.1      |

Let's visualize the relationships between streams and other numerical variables that are not related to the intrinsic artistic characteristics of a song.

```
[42] #Distribution of songs across different release years
release_year_distribution = df['released_year'].value_counts().sort_index()

plt.figure(figsize=(10, 6))
sns.barplot(x=release_year_distribution.index, y=release_year_distribution.values, palette="viridis")
plt.title("Distribution of Songs Across Different Release Years")
plt.xlabel("Release Year")
plt.ylabel("Number of Songs")
plt.xticks(rotation=45)
plt.grid(True)
plt.show()
```



#### ▼ PART 3: DIAGNOSTIC ANALYTICS, DATA VISUALIZATION AND TAKEAWAYS

##### ▼ ANOVA

After visualizing data distribution by categorical variables (artist\_count, key and mode), we wanted to test if the streams and avg monthly streams are truly different across different levels of the categorical variables

##### ▼ Artist Count vs Streams

As we can see from the box plot in part 2, there is a lot of noise for songs with artist count > 4. For this reason, we decided to do an ANOVA

```
Collecting scikit-posthocs
  Downloading scikit_posthocs-0.8.0-py3-none-any.whl (32 kB)
Requirement already satisfied: numpy in /usr/local/lib/python3.10/dist-packages (from scikit-posthocs) (1.23.5)
Requirement already satisfied: scipy in /usr/local/lib/python3.10/dist-packages (from scikit-posthocs) (1.11.4)
Requirement already satisfied: statsmodels in /usr/local/lib/python3.10/dist-packages (from scikit-posthocs) (0.14.0)
Requirement already satisfied: pandas>=0.20.0 in /usr/local/lib/python3.10/dist-packages (from scikit-posthocs) (1.1.5.0)
Requirement already satisfied: matplotlib>=3.0.0 in /usr/local/lib/python3.10/dist-packages (from scikit-posthocs) (3.2.1)
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (from scikit-posthocs) (3.7.1)
Requirement already satisfied: python-dateutil<2.8.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.20.0->scikit-posthocs) (2.8.2)
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas>=0.20.0->scikit-posthocs) (2023.3.post1)
Requirement already satisfied: cycler>=0.10.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib>=3.0.0->scikit-posthocs) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in /usr/local/lib/python3.10/dist-packages (from matplotlib>=3.0.0->scikit-posthocs) (4.46.0)
Requirement already satisfied: kiwisolver>1.6.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib>=3.0.0->scikit-posthocs) (1.4.5)
Requirement already satisfied: packaging>=20.4 in /usr/local/lib/python3.10/dist-packages (from matplotlib>=3.0.0->scikit-posthocs) (25.2)
Requirement already satisfied: patsy>=0.5.2 in /usr/local/lib/python3.10/dist-packages (from statsmodels>=0.13.2->scikit-posthocs) (0.5.4)
Requirement already satisfied: parsings>2.3.1 in /usr/local/lib/python3.10/dist-packages (from matplotlib>=3.0.0->scikit-posthocs) (3.1.1)
Requirement already satisfied: patsy>=0.5.2 in /usr/local/lib/python3.10/dist-packages (from statsmodels>=0.13.2->scikit-posthocs) (0.5.4)
Requirement already satisfied: six in /usr/local/lib/python3.10/dist-packages (from patsy>=0.5.2->statsmodels>=scikit-posthocs) (1.16.0)
Installing collected packages: scikit-posthocs
Successfully installed scikit-posthocs-0.8.0
```

```
[46] import scikit_posthocs as sp
p_values1 = sp.posthoc_dunn(df, val_col='avg_streams_per_month', group_col='artist_count')
```

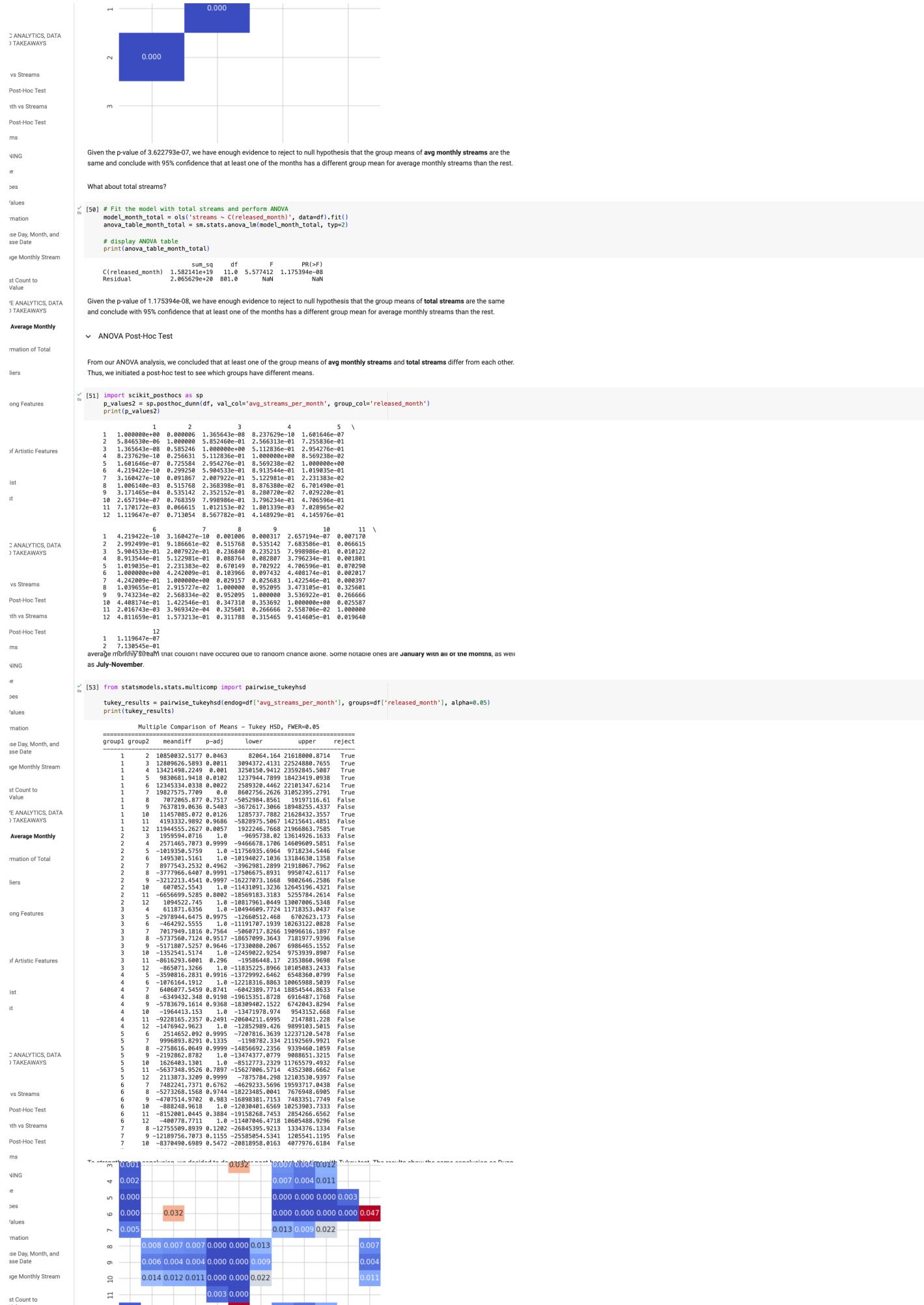
|   | 1        | 2        | 3        | 4        |
|---|----------|----------|----------|----------|
| 1 | 1.000000 | 0.000002 | 0.125544 | 0.576665 |
| 2 | 0.000002 | 1.000000 | 0.150235 | 0.230016 |
| 3 | 0.125544 | 0.150235 | 1.000000 | 0.766585 |
| 4 | 0.576665 | 0.230016 | 0.766585 | 1.000000 |

```
[47] alpha = 0.05
mask = p_values1 < alpha

plt.figure(figsize=(8, 6))
sns.heatmap(p_values1, annot=True, cmap='coolwarm', fmt=".3f", mask=mask, cbar=False)

plt.title("Dunn Test p-values of artist count and average streams per month")
plt.xlabel("Artist Count")
plt.show()
```

Dunn Test p-values of artist count and average streams per month







```

liers
    key_G#          -8.416e+06  1.66e+07   -9.588   0.612   -4.1e+07  2.41e+07
    mode_Minor     6.875e+06  1.26e+07   0.545   0.586   -1.79e+07  3.17e+07
=====
Omnibus:           218.824 Durbin-Watson:            2.035
Prob(Omnibus):      0.000 Jarque-Bera (JB):       1645.999
Skew:                 1.499 Prob(JB):                  0.00
Kurtosis:                18.774 Cond. No.:             5.93
=====

Notes:
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

of Artistic Features
    ✕ Cross-Validation

ist
    ✕ [90] from sklearn.model_selection import KFold
    from sklearn.model_selection import GridSearchCV
    from sklearn.pipeline import make_pipeline
    from sklearn.model_selection import cross_val_score

    ✕ # Perform hyperparameter tuning to find the best hyperparameters for the linear regression model.

    # k-fold Cross-validation
    # we can't implement cross validation using statsmodels
    # so to get the validation score, we switch to sklearn's linear regression model
    # fit_intercept = True by default so no need to add constant term
    lm = LinearRegression()

vs Streams
    # fit sklearn linear regression on training set
    lm.fit(X_train, y_train)

Post-Hoc Test
    # used trained sklearn linear regression model "lm" to predict mpg for testing data
    predictions = lm.predict(X_test)

Post-Hoc Test
    R2 = cross_val_score(lm, X_train, y_train, scoring='r2', cv=10) #default 5-fold cross validation

ms
    - - - - -
    'danceability_%', 'valence_%', 'energy_%', 'liveness_%', 'speechiness_%',
    'artist_count',
    'key_A#', 'key_B#', 'key_C#', 'key_D#', 'key_E#', 'key_F#', 'key_G#', 'key_H#',
    'mode_Minor',
    'released_day']

NING
ie
pes
values
imation
use Day, Month, and
use Date
ige Monthly Stream
st Count to
Value
IE ANALYTICS, DATA
) TAKEAWAYS
Average Monthly
formation of Total
liers
ong Features
    ✕ [91] # Set up the environment to ignore warnings.
    import warnings
    warnings.filterwarnings("ignore")

    ✕ [92] # add a constant term for intercept
    X_train_const_2 = sm.add_constant(X_train_2)

    # Create a linear regression model and train it with the training data.
    model2 = sm.OLS(y_train, X_train_const_2)
    result = model2.fit()

    ✕ [93] print(result.summary(xname=['const',
                                         'released_year', 'released_month',
                                         'in_spotify_playlists', 'in_spotify_charts', 'in_apple_playlists', 'in_apple_charts',
                                         'in_deezer_charts', 'in_shazam_charts', 'bpm', 'acousticness_%',
                                         'instrumentalness_%'],
                                         type='object'))
=====
OLS Regression Results
=====
Dep. Variable: streams R-squared:      0.716
Model: OLS Adj. R-squared:      0.710
Method: Least Squares F-statistic:   127.6
Date: Wed, 13 Dec 2023 Prob (F-statistic): 2.90e-144
Time: 22:15:52 Log-Likelihood: -11869.
No. Observations: 59 AIC: 2.376e+04
Df Residuals: 557 BIC: 2.381e+04
Df Model: 11
Covariance Type: nonrobust
=====
coef std err t P>|t| [0.025  0.975]
=====
const      4.519e+08  1.17e+07  38.498  0.000  4.29e+08  4.75e+08
released_year  4.566e+07  1.33e+07  3.433  0.001  1.95e+07  7.18e+07
released_month  4.456e+07  1.19e+07  3.722  0.000  0.715  -1.91e+07  2.77e+07
in_spotify_playlists  2.96e+08  1.83e+07  16.192  0.000  2.6e+08  3.32e+08
in_spotify_charts  9.433e+07  1.09e+07  8.582  0.000  7.4e+07  1.79e+08
gridSearch = GridSearchCV(Ridge, param_grid, scoring="r2", cv=5)
gridSearch.fit(X_train, y_train)

ie
    # Get the best hyperparameters
    best_alpha = gridSearch.best_params_['alpha']

pes
    # Train the model with the best hyperparameters
    best_ridge_model = Ridge(alpha=best_alpha)
    best_ridge_model.fit(X_train, y_train)

values
    # Evaluate the model on the validation set
    y_val_pred = best_ridge_model.predict(X_test)
    mse = mean_squared_error(y_test, predictions)
    r2 = r2_score(y_test, predictions)

    # calculate adj r2
    n = len(X_train)
    p = len(X.columns)
    adj_r2 = 1 - ((1-r2) * (n-1)/(n-p))

    print("Best alpha: (best_alpha)")
    print("R-squared on validation set: (r2)")
    print("Adjusted R-squared on validation set: (adj_r2)")

    Best alpha: 10
    R-squared on validation set: 0.685836267799053
    Adjusted R-squared on validation set: 0.6683178440703759

liers
    After attempting linear regression model with only selective statistically significant variables, the R-squared of the model which in this case evaluates how much variation is accounted in the model did not change much. This tells us that linear regression is not a good model for predicting streams.

ong Features
    ✕ Logistic Regression

of Artistic Features
    Since we have a lot of discrete variables in our data set, it makes sense that linear regression is not the ideal model. Instead we decided to try logistic regression for classification.

First, we need to define the binary outcome variable. In this case, we wanted to predict what makes a song a "top hit", which is the top 25%/first quartile of number of streams on Spotify.

    ✕ Define Outcome Variable: Top Hits

ist
    ✕ [102] # let's use the top 25% streams as "meg hits"
    streamsQ1 = df1['streams'].quantile(0.75)
    streamsQ1
    601863821.0

Post-Hoc Test
    ✕ [103] # Create a new data frame for logistic regression
    df2 = df1.copy()

    # Create a new binary variable named "TopHits" where the value is 1 if streams is above top 25% streams.
    TopHits = np.where(df2['streams'] > streamsQ1, 1, 0)

    # Add column to dataframe.

```





```

'values
mation
ise Day, Month, and
ase Date
ige Monthly Stream
st Count to
Value
'E ANALYTICS, DATA
) TAKEAWAYS
Average Monthly
rmation of Total
liers
ong Features
of Artistic Features
ist
st
C ANALYTICS, DATA
) TAKEAWAYS
vs Streams
Post-Hoc Test
rth vs Streams
Post-Hoc Test
ms
NING
ie
pes
'values
mation
ise Day, Month, and
ase Date
NING
ie
pes
'values
mation
ise Day, Month, and
ase Date
ige Monthly Stream
st Count to
Value
'E ANALYTICS, DATA
) TAKEAWAYS
Average Monthly
rmation of Total
liers
ong Features
of Artistic Features
ist
st
C ANALYTICS, DATA
) TAKEAWAYS
vs Streams
Post-Hoc Test
rth vs Streams
Post-Hoc Test
ms

```

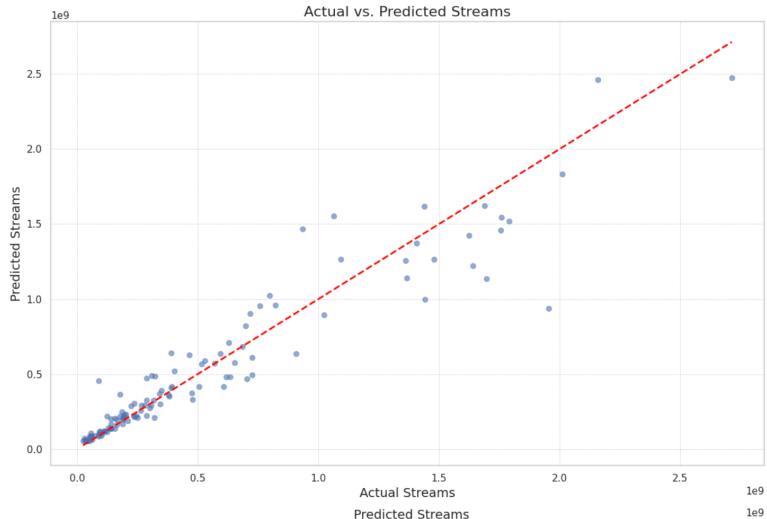
[120] #Now, let us plot the Actual vs. Predicted Streams to see how our model has performed

```

plt.figure(figsize=(12, 8))
sns.scatterplot(x=y_test, y=y_pred, alpha=0.6, edgecolor=None)
plt.plot([min(y_test), max(y_test)], [min(y_test), max(y_test)], color='red', linestyle='--', lw=2)
plt.xlabel('Actual Streams', fontsize=14)
plt.ylabel('Predicted Streams', fontsize=14)
plt.title('Actual vs. Predicted Streams', fontsize=16)
plt.grid(True, which='both', linestyle='--', linewidth=0.5)
plt.tight_layout()
plt.show()

```

#Here, especially for the less listened to songs, we see that we have reached a relatively decent accuracy based off of looking at this graphs.  
#Note the numbers at the axes are in billions



Here, we have reached a relatively decent accuracy based off of looking at this graphs.  
However, we do have some outliers on the residual plot. The shape of the residual plot is also fan-shaped. This pattern may indicate that there is relationships that are not captured by the model.

**\*\*Note the numbers at the x-axis are in billions and the numbers on the y-axis are in millions.**

- ↳ **DATASET 5: DESCRIPTIVE ANALYTICS, DATA VISUALIZATION AND TAKEAWAYS**
- We could experiment with other machine learning models like Gradient Boosting, XGBoost, or neural networks to compare performance.
- We could incorporate sentiment analysis of social media buzz around songs or artists to see if it correlates with streaming numbers.
- We could explore differences in streaming behavior across different geographic regions or demographic groups.
- We could utilize grid search or randomized search for hyperparameter tuning of the random forest model to further optimize its performance.

To strengthen the credibility of our analysis results (and in extension, our recommendations), we explored previous research done on music popularity. In addition, we wanted to know where our current results stand amongst what is currently known. To our surprise, Random Forest model is consistently the highest performing (most accurate) machine learning algorithm for music popularity predictions, compared to other approaches such as: Linear Regression, Logistic Regression, Neural Networks, Lasso Regression, Deep Learning, Boosting Tree, etc. (Arora, Rani, & Saxena, 2022; Colley et al., 2022; Essa et al., 2022; Gao, 2021; Khan et al., 2022). Several researchers were even successful in creating a Random Forest model with a 95.37% accuracy rate (Gulmatico et al., 2022). Based on previous research, we believe our Random Forest model has potential, and that its superior performance against our Linear Regression model isn't unusual.

In reviewing further research, it is evident that many external factors have profound effects on music streaming. Research shows that there is a strong correlation between numbers of followers and royalties from music streaming specifically. Set in the context of the high level of income precarity being experienced amongst independent artists, this finding suggests that a positive co-evolution is occurring between social media platforms and music streaming with potentially significant increases in artist royalty incomes resulting from increases in followers (Watson, A., Watson, J. B., & Tompkins, L., 2023). This combined with the importance of being in not only Spotify playlists but also playlists on other music streaming platforms that we've found through our model suggests that → to achieve streaming success, network and marketing is a key factor. Based on our data analysis and industry research, we propose the following recommendations for musicians:

#### Leverage Social Media Platforms for Increased Followership:

Engage actively on social media platforms to increase your followership. As evidenced by previous research, a positive correlation exists between the number of followers and artist royalty incomes from music streaming. Utilize platforms such as Instagram, Twitter, and Facebook to connect with your audience, share updates, and promote your music. Optimize Presence on Spotify Playlists:

Our model highlights the importance of being featured in Spotify playlists for increased music streaming. Work towards getting your songs included in relevant playlists to enhance discoverability. Collaborate with curators, promote your music to influencers, and ensure your songs align with the themes of popular playlists.

#### Diversify Across Multiple Music Streaming Platforms:

Extend your reach beyond Spotify and explore other music streaming platforms. Our analysis suggests that being featured in playlists on various platforms contributes to overall streaming success. Platforms like Apple Music, Deezer, and YouTube Music offer additional opportunities to reach diverse audiences. Collaborate with Other Artists and Curators:

Network within the music industry by collaborating with other artists and connecting with playlist curators. Cross-promotions and collaborative projects can broaden your audience and increase your visibility. Building relationships with curators enhances the chances of getting featured in curated playlists. Strategic Marketing Campaigns:

Develop and execute strategic marketing campaigns to promote your music. Utilize online advertising, influencer partnerships, and content marketing to reach a wider audience. Tailor your campaigns to target specific demographics and genres to maximize impact.

#### Engage with Fans and Collect Feedback:

Foster a strong connection with your fansbase by actively engaging with them on social media and other platforms. Encourage feedback, conduct polls, and consider fan preferences in your music production. Building a loyal fan community contributes to sustained streaming success. Monitor Streaming Analytics and Iterate Strategies:

Regularly analyze streaming analytics to assess the impact of your efforts. Monitor which strategies contribute most to streaming numbers and adjust your approach accordingly. Iterative improvements based on data insights are crucial for long-term success.

In conclusion, success in music streaming is not solely determined by the quality of the music but also by strategic networking and marketing

↳ **DATA 5: DESCRIPTIVE ANALYTICS, DATA VISUALIZATION AND TAKEAWAYS**