# Agenda

- Problem
- Data
- Models
- Interface
- Up Next

# Problem

# Scope and Background

- The goal: Predict success of a song defined by the number of Spotify plays
  - Help small artists improve the likelihood of success based on metrics
- Success is defined this way as music is because the goal of most artists is to share their music with as many people as possible!

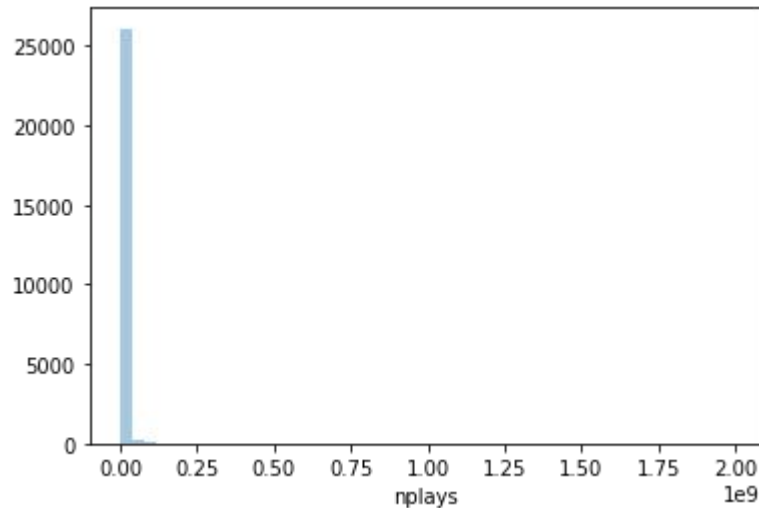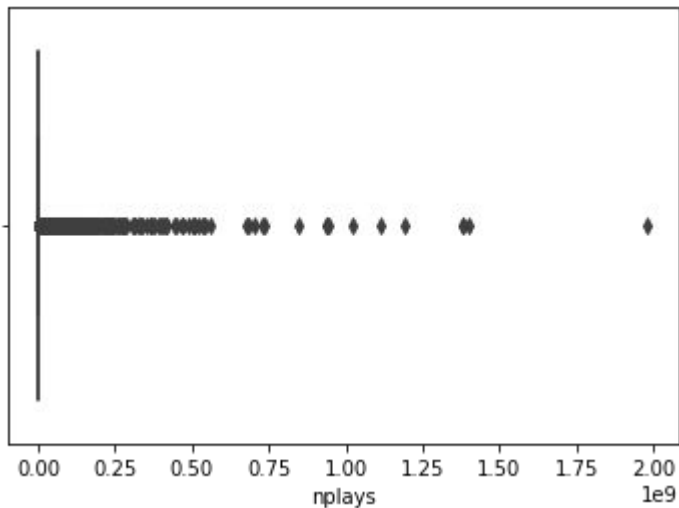Problem          Data          Models          Interface          Up Next
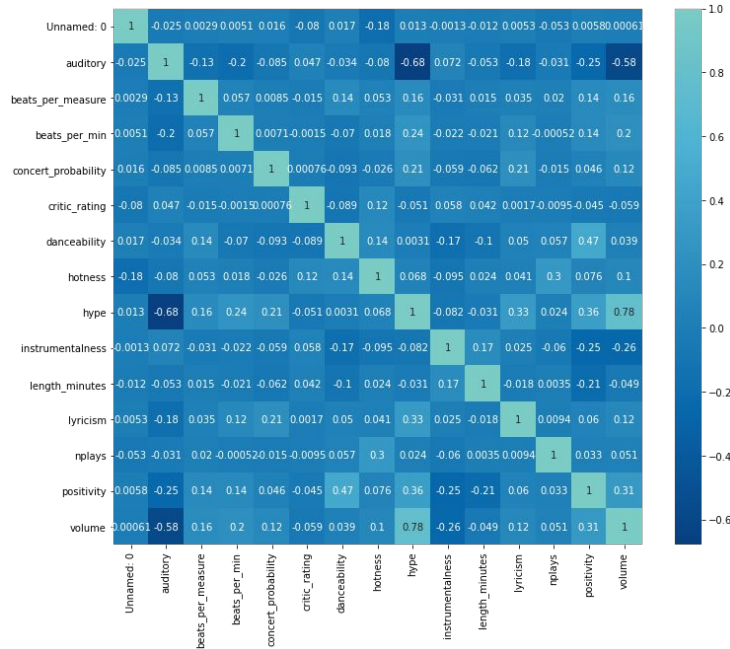
Data

# Data Scope - nplays

- Nplays has a huge spread
  - Eliminate outliers outside 3 standard deviations
  - Log scale the values

# Data Scope - Correlations

- Upon looking at correlations decided to remove all about 0.5

# Data Preprocessing

- **Eliminate metrics that can not be predicted before release(lookahead bias)**
  - Hotness
  - Critic information
- Translated information directly related to individual artist
  - Artist name → length, is uppercase, positive sentiment, negative sentiment
  - Album name → length, is uppercase, positive sentiment, negative sentiment
  - Song name → length, is uppercase, positive sentiment, negative sentiment
- Log scaled nplays
- Replace all NaNs in features with the most common class
- Removed variables with correlation above 0.5
- Min Max scaled the numerical features to [0, 1]

Problem          Data          Models          Interface          Up Next

Models

# Models

- XGBoost, Catboost
  - Performed similarly overall
  - Some hyper parameter tuning revealed that we didn't need to deep of a tree
  - Since the styles feature ended up mattering a lot, catboost automatic categorical feature handling could have helped it outperform xgb
  - Achieved final $R^2$ score of 0.49 with Catboost
- Recurrent neural networks
  - Used textgenrnn to train many different models to generate reviews
  - User input used to select which model to use for review generation (currently based on genre)

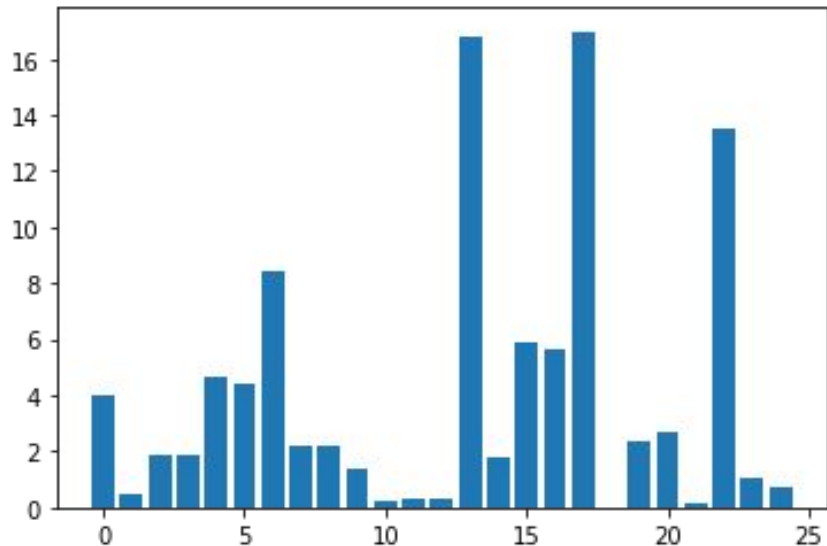Problem          Data          Models          Interface          Up Next

# Insights

- Top 3 features according to our tuned Catboost model
    - Length of album name
    - Length of artist name
    - Style
- Could be indicative that the visual aspect of the song could influence whether it is streamed or not
    - Also possible that these metrics still too closely associate the song with artist
- Some styles are just more popular than others
    - Styes were also consistently an important feature through our model exploration

Interface

# Practical Application

- Flask and React App built on our model, so a user may input info about their song and be returned the expected number of Spotify plays
- The user would not see the complex model underneath



Artist Name*    Song Name*    Album Name*

### Song Properties

Beats per Measure    Beats per Minute    Length in Minutes    Auditory

Lyricism    Volume    Danceability    Positivity

Hype    Instrumentalness    Style

### Tone

Major    Minor
○        ●

Key*

### Vulgarity

Yes    No
○      ●

### Misc.

Concert Probability*

CALCULATE

**Projected Streams**

-

**Sample Review**

-

Problem    Data    Models    Interface    Up Next

Up Next

# Future Plans

- Upload a finished song and our algorithm would calculate the song statistics
- Displaying what metrics had the most positive and negative influence on the number of plays so artists can adjust this in real time
- Implement grammarly API to help generate reviews that make sense
- Categorize music by genre to improve metrics and success prediction

Problem          Data          Models          Interface          Up Next

Questions