# Spotify Sequential Skip Prediction Challenge
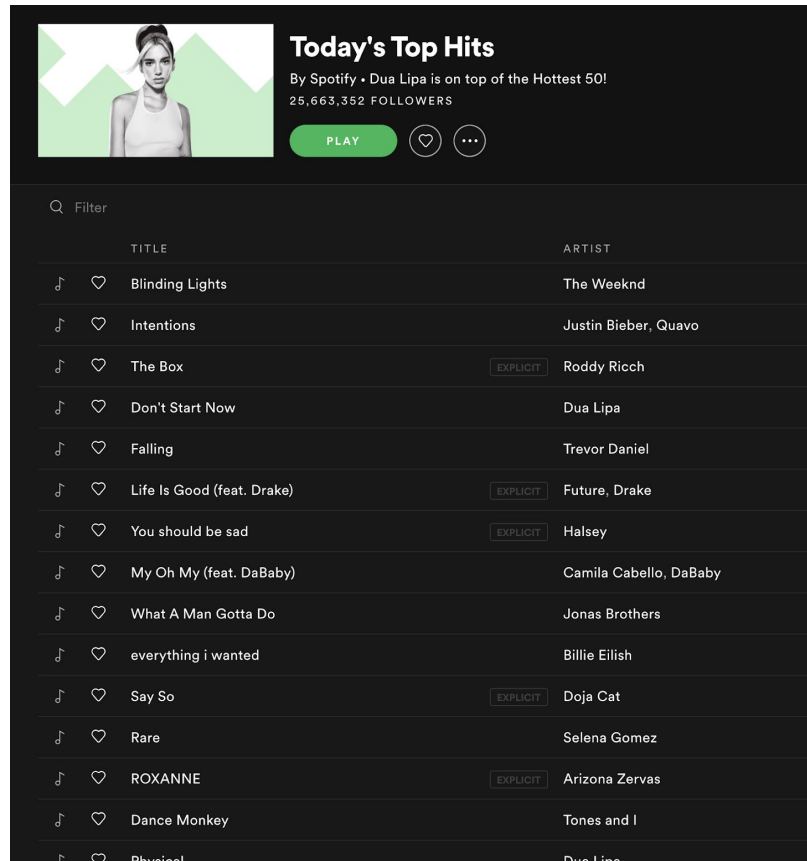
Austin Poor

# Introduction

"Spotify has over **190 million** active users interacting with over **40 million** tracks"

___

# The Challenge

**Goal:** *Predict the likelihood of a user skipping any given song during a listening session*
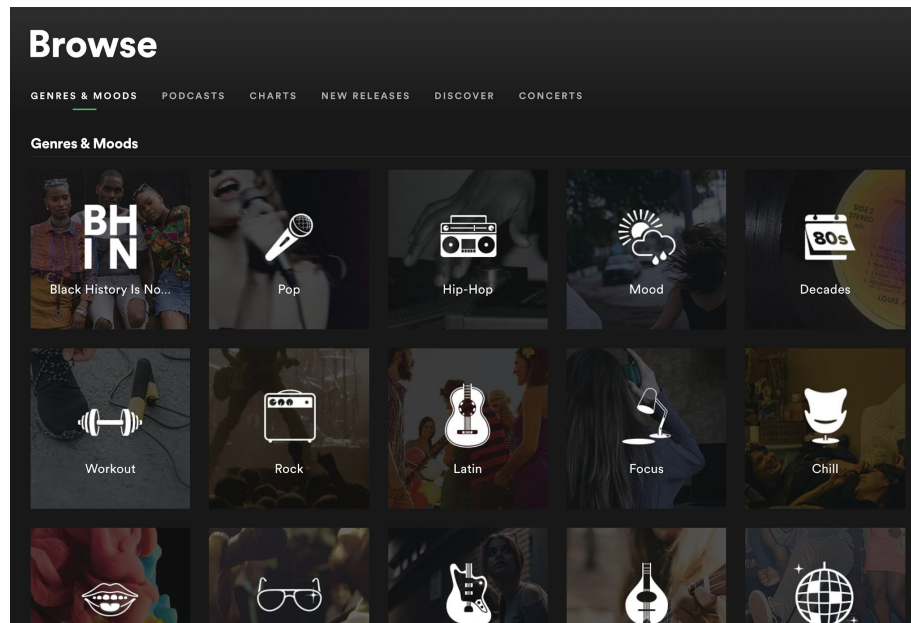
# Methodology

# The Data

- 10-20 song sessions
- Song and user data was anonymized
- Balanced classes (0.517 skips)
- 350GB of CSV data
- Used subset of the data
  - 100k Session Rows
  - 50k Track Rows
- Stored data in PostgreSQL DB

| training_set | |
|---|---|
| **session_id** | text |
| **session_position** | bigint |
| session_length | bigint |
| track_id_clean | text |
| skip_1 | boolean |
| skip_2 | boolean |
| skip_3 | boolean |
| not_skipped | boolean |
| context_switch | bigint |
| no_pause_before_play | bigint |
| short_pause_before_play | bigint |
| long_pause_before_play | bigint |
| hist_user_behavior_n_seekfwd | bigint |
| hist_user_behavior_n_seekback | bigint |
| hist_user_behavior_is_shuffle | boolean |
| hour_of_day | bigint |
| date | text |
| premium | boolean |
| context_type | text |
| hist_user_behavior_reason_start | text |
| hist_user_behavior_reason_end | text |

| track_features | |
|---|---|
| **track_id** | text |
| duration | double |
| release_year | bigint |
| us_popularity_estimate | double |
| acousticness | double |
| beat_strength | double |
| bounciness | double |
| danceability | double |
| dyn_range_mean | double |
| energy | double |
| flatness | double |
| instrumentalness | double |
| key | bigint |
| liveness | double |
| loudness | double |
| mechanism | double |
| mode | text |
| organism | double |
| speechiness | double |
| tempo | double |
| time_signature | bigint |
| valence | double |

| acoustic_vectors | |
|---|---|
| **track_id** | text |
| acoustic_vector_0 | double |
| acoustic_vector_1 | double |
| acoustic_vector_2 | double |
| acoustic_vector_3 | double |
| acoustic_vector_4 | double |
| acoustic_vector_5 | double |
| acoustic_vector_6 | double |
| acoustic_vector_7 | double |

# Feature Engineering

Generate features to account for a user's listening history

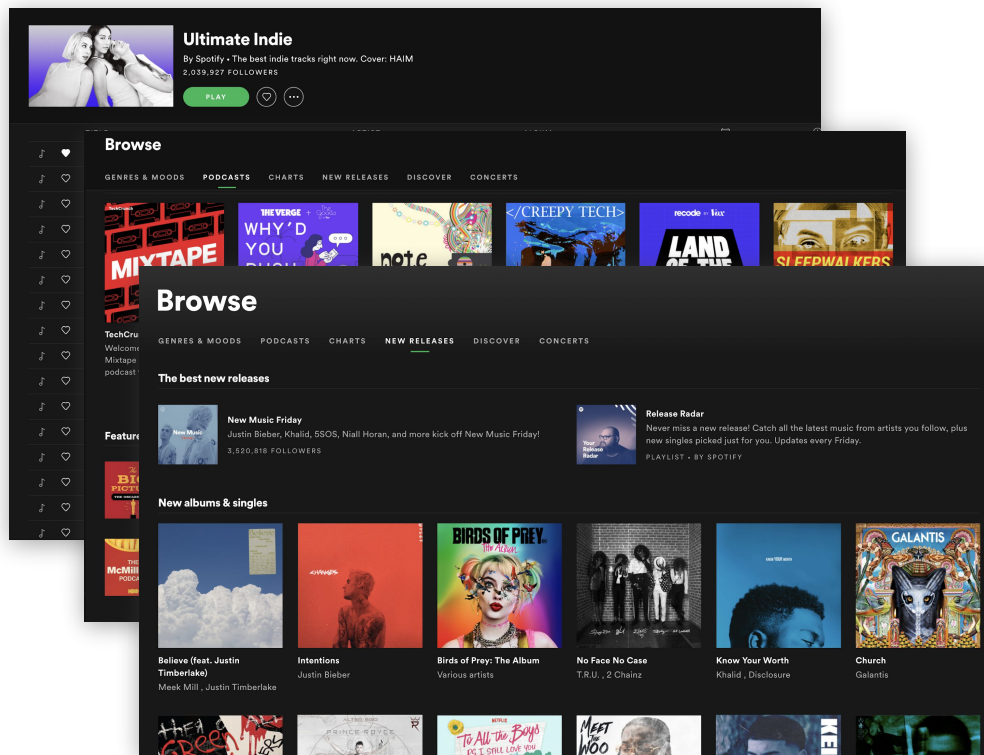Added previous track features, including if that track was skipped

# Model Selection

Target Metric **Accuracy** per competition guidelines

Baselined with Logistic Regression

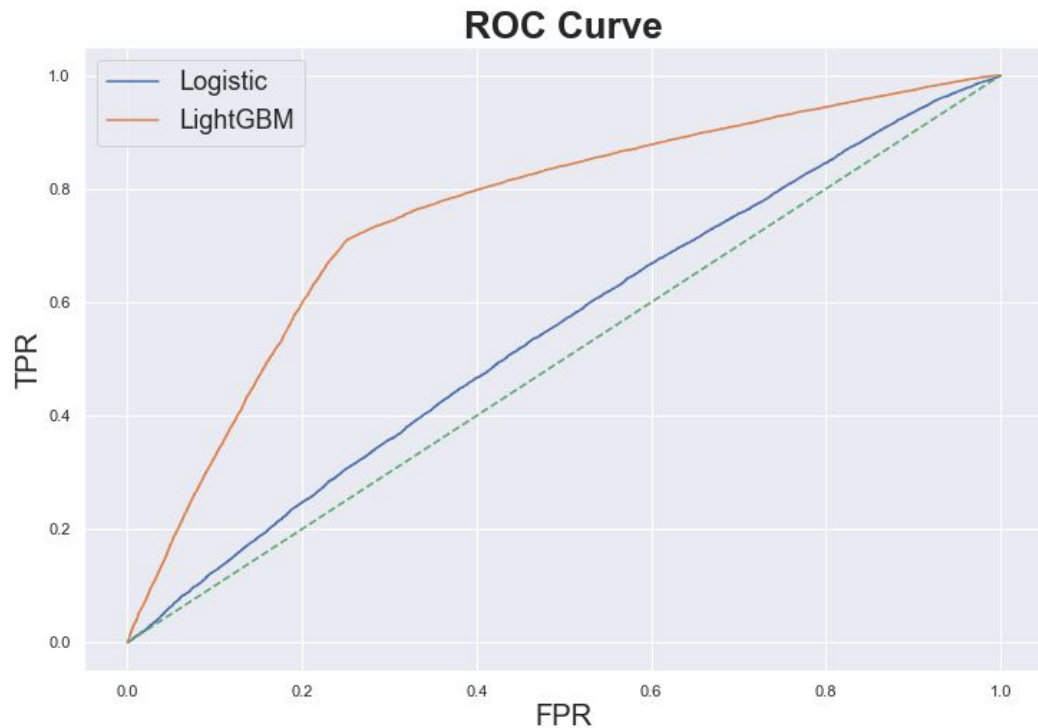Moved to tree-based models which would automatically handle feature interactions

# Results

# Model Results

Best Test Accuracy: **0.73** (with LightGBM)

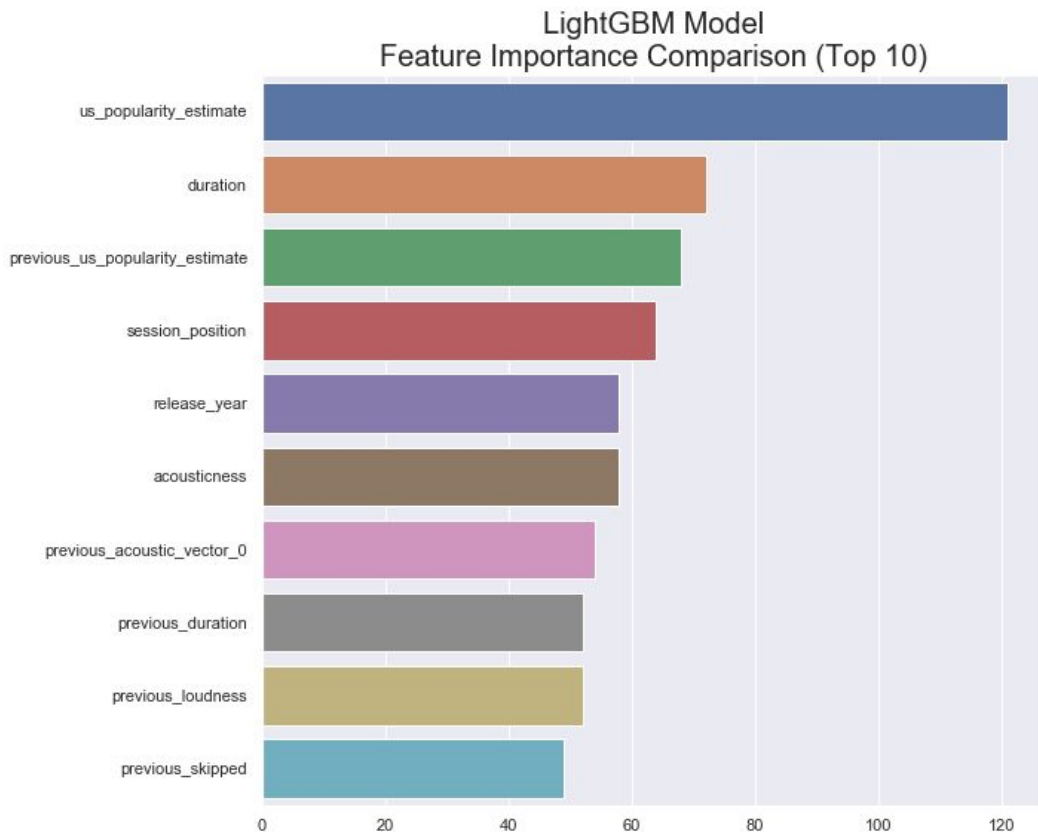Error Analysis: No clear trend in the model residuals

# Conclusions

# Conclusions

us_popularity_estimate had the highest
feature importance followed by duration
and then previous_us_popularity_estimate

Model results are pretty good but with
room for improvement



LightGBM Model
Feature Importance Comparison (Top 10)

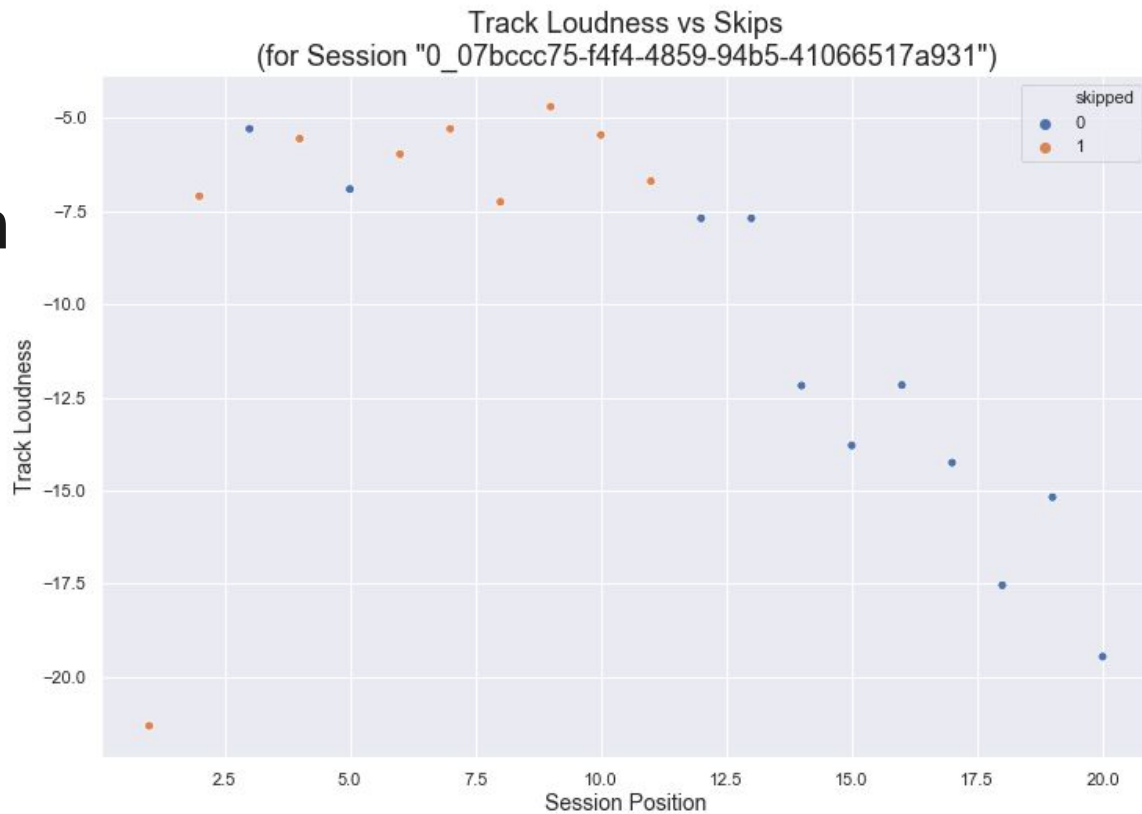# Future Work

# Future Work

- Test other types of algorithms:
    - Unsupervised Learning to cluster songs
    - RNN to predict based on the sequence of tracks
- Supplement the dataset with more data from the Spotify API
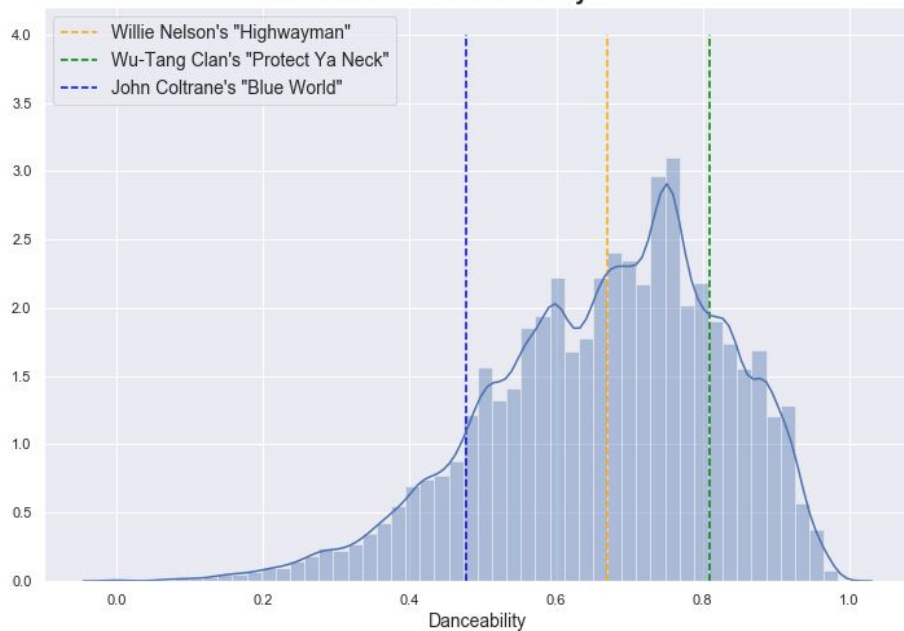- Create a Flask app to visualize predictions using D3

# Thank you

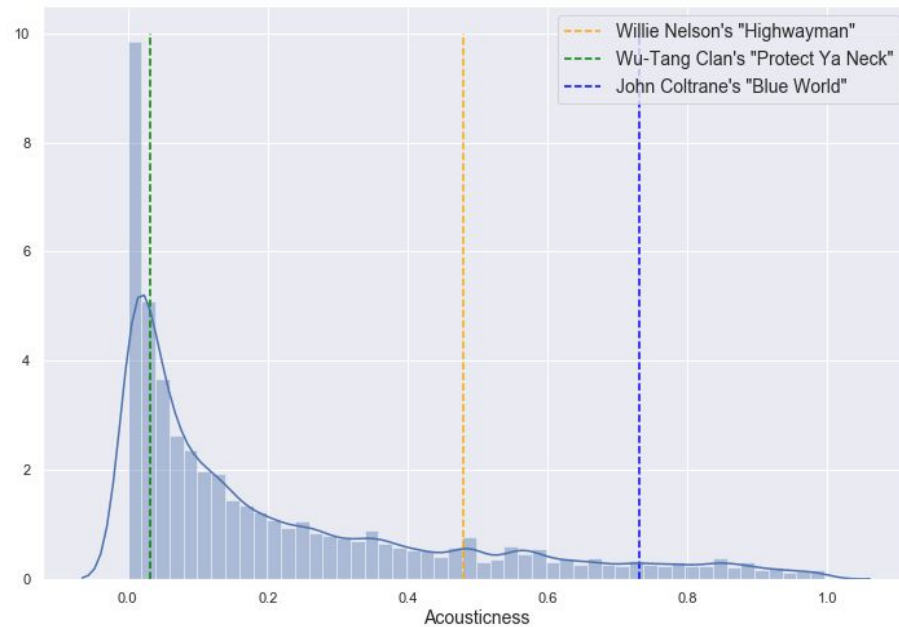# Appendix

# Visualizing a Listening Session



Track Loudness vs Skips
(for Session "0_07bccc75-f4f4-4859-94b5-41066517a931")

# Exploring Features



Distribution of "Danceability" Values

- Willie Nelson's "Highwayman"
- Wu-Tang Clan's "Protect Ya Neck"
- John Coltrane's "Blue World"

Danceability

Distribution of "Acousticness" Values

- Willie Nelson's "Highwayman"
- Wu-Tang Clan's "Protect Ya Neck"
- John Coltrane's "Blue World"

Acousticness

# Exploring Features



Distribution of "Energy" Values — Legend: Willie Nelson's "Highwayman", Wu-Tang Clan's "Protect Ya Neck", John Coltrane's "Blue World"

Distribution of "Loudness" Values — Legend: Willie Nelson's "Highwayman", Wu-Tang Clan's "Protect Ya Neck", John Coltrane's "Blue World"

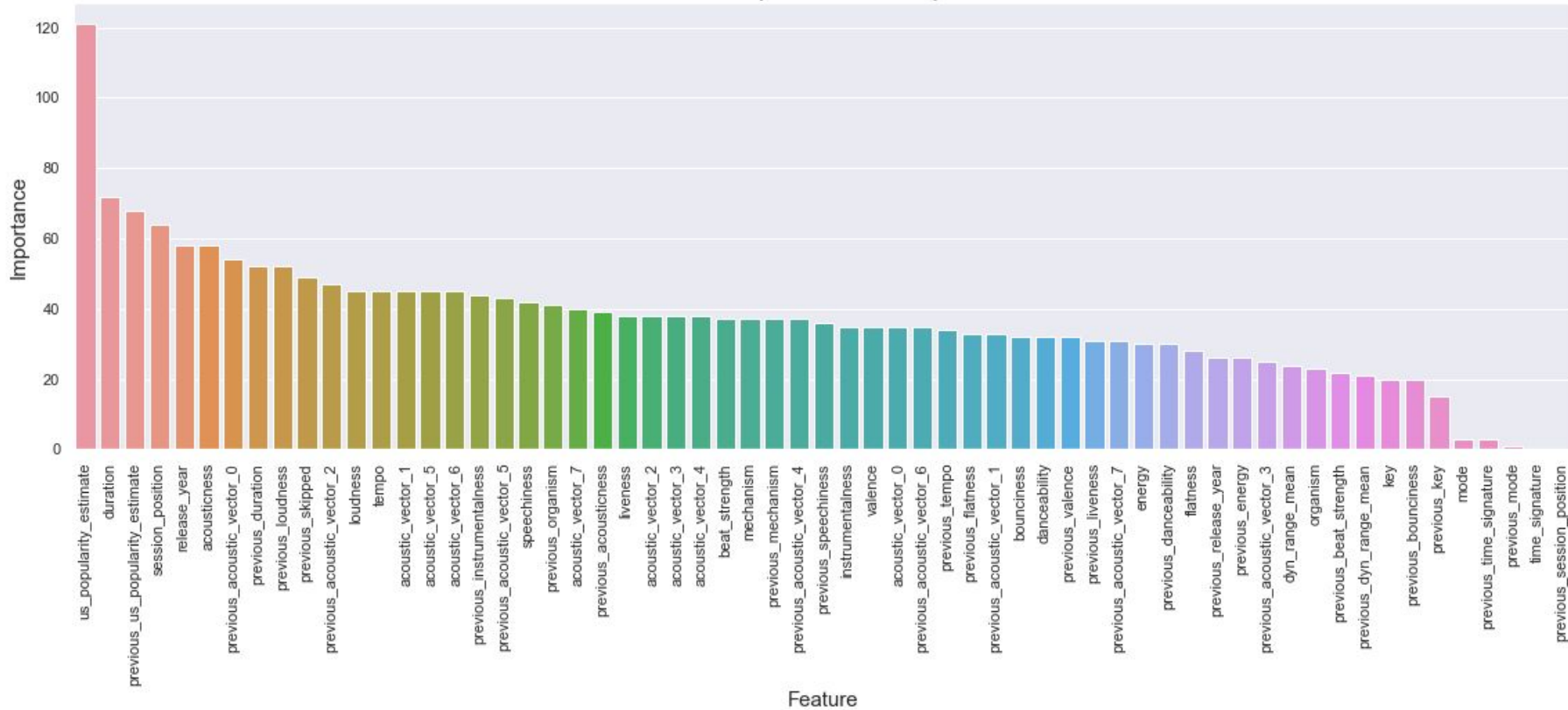# Exploring Features

| Feature | Description |
| --- | --- |
| *Acousticness* | Likelihood a track is acoustic |
| *Danceability* | Describes how suitable a track is for dancing |
| *Energy* | Measure of a track's intensity |
| *Valence* | Level of "positiveness" conveyed by a track |
| *Speachiness* | Detects the presence of spoken word |

# Acoustic Vector Correlations



Correlation Matrix for Acoustic Vectors

LightGBM Model
Feature Importance Comparison

# Flask App with D3 Visualizations