



Drifting Upstream

Predicting song popularity from
music features

The Problem

- Regression Records music label
- Extra promotional budget for Q3 and Q4 of 2024
- Hundreds of songs in catalog, must choose 3
- **Predict song popularity from music features**

Data Wrangling & EDA Pt. 1

945 Songs

Target Variable
Total # Streams

Median ~290 million

Min ~1.4 million

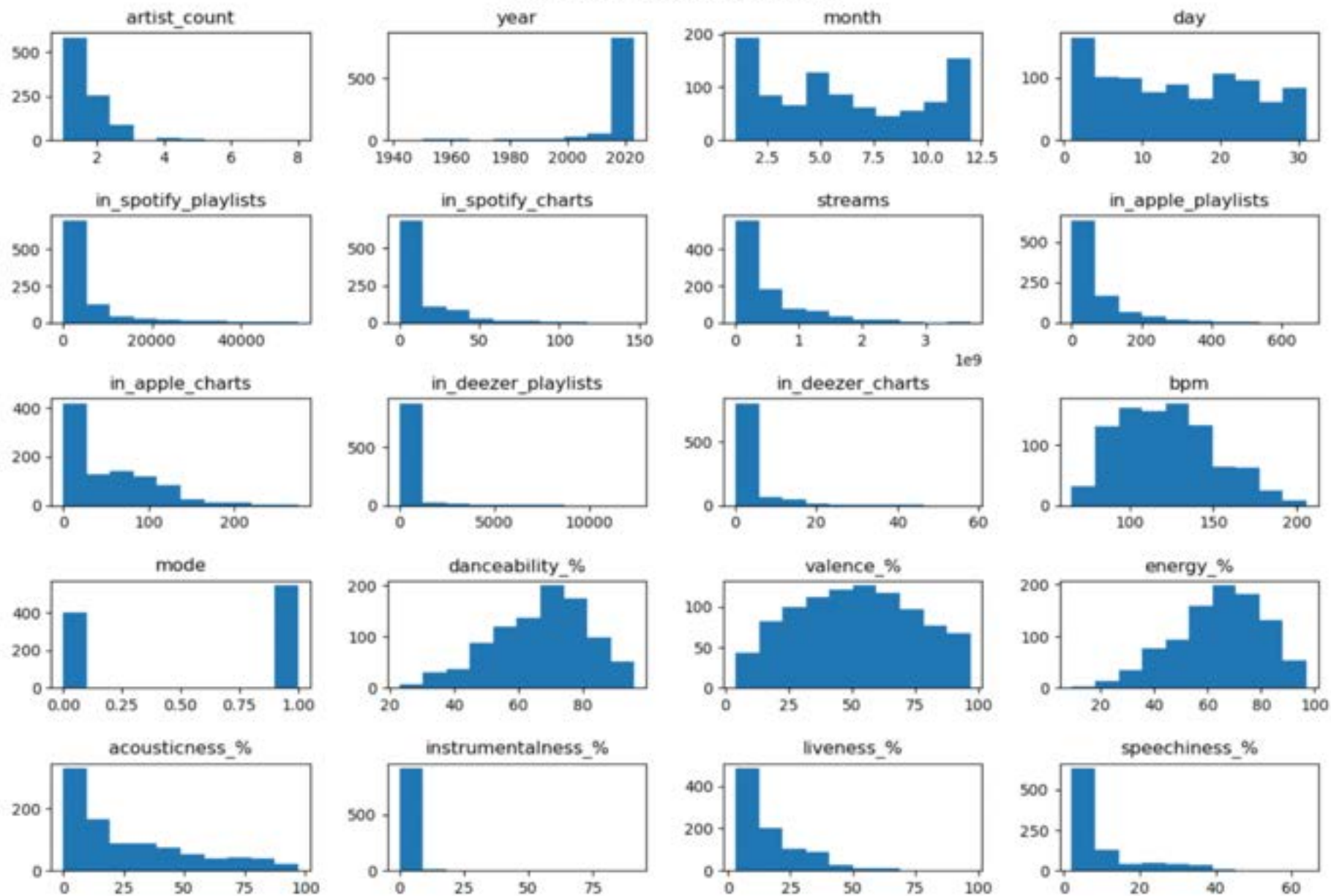
Max ~3.7 billion

Data Wrangling & EDA Pt. 1

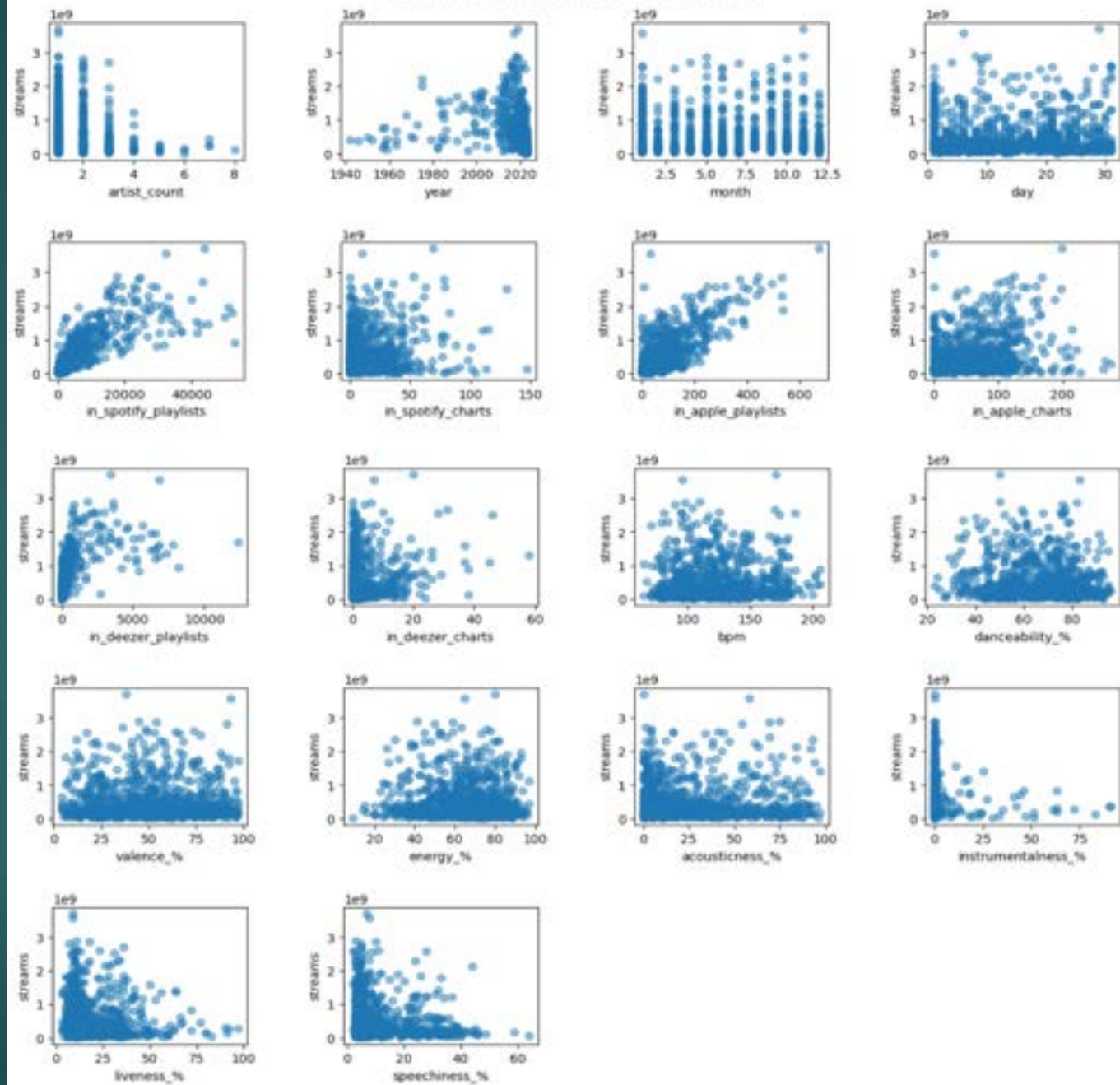
Features

- Tempo
- Danceability
- Valence
- Energy
- Acousticness
- Instrumentalness
- Liveness
- Speechiness
- Mode

Feature Distributions (Dataset 1)



Streams / Feature Correlations (Dataset 1)



Modeling Pt. 1

Random Forest Regressor

- Train / Test MAE 335.5M / 332.4M

Extra Trees Regressor

- Train / Test MAE 0 / 340M

Gradient Boosting Regressor

- Train / Test MAE 310M / 348M

(70 / 30 train test split)

Modeling Pt. 1

Random Forest Regressor

PowerTransformer()

- Train / Test MAE 335.5M / 332.2M

90 / 10 Split

- Train / Test MAE 332M / 268M

Remove top 5% outliers (80 / 20)

- Train / Test MAE 247M / 270M

Modeling Pt. 1

Random Forest Regressor

Hyperparameters

- `n_estimators = 200`
- `max_features = 6`
- `max_depth = 30`
- `min_samples_split = 7`
- `min_samples_leaf = 11`
- `random_state = 42`

Feature Importance

1. **Tempo (0.166281)**
2. **Danceability (0.161662)**
3. **Acousticness (0.155589)**
4. **Valence (0.146104)**
5. Speechiness (0.119351)
6. Liveness (0.115062)
7. Energy (0.109871)
8. Mode (0.023337)
9. Instrumentalness (0.002743)

Data Wrangling & EDA Pt. 2



32,833 Songs

Target Variable
Popularity Score (1-100)

Mean 42.48

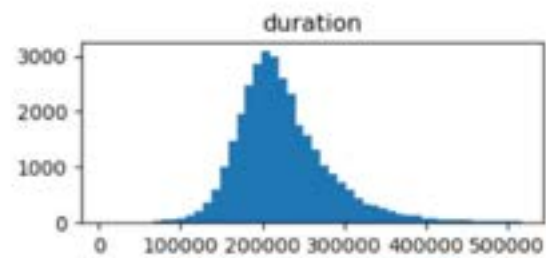
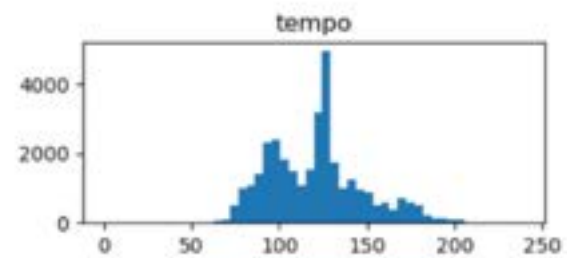
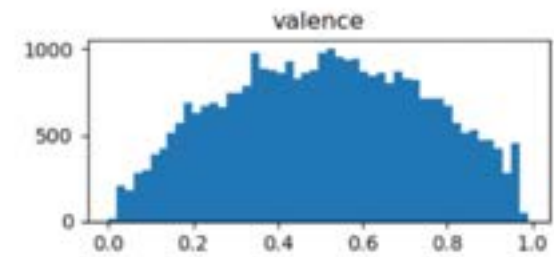
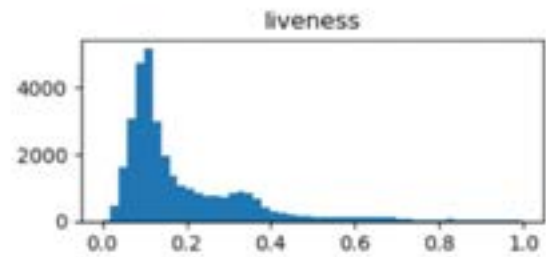
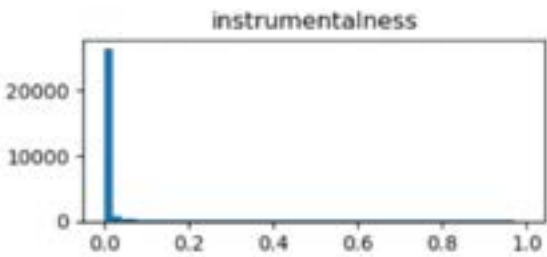
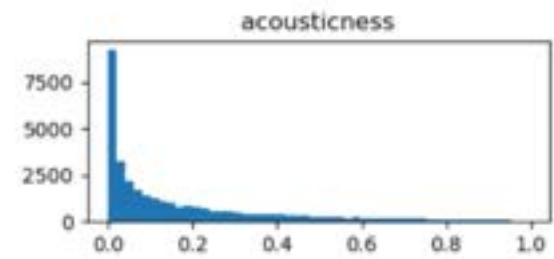
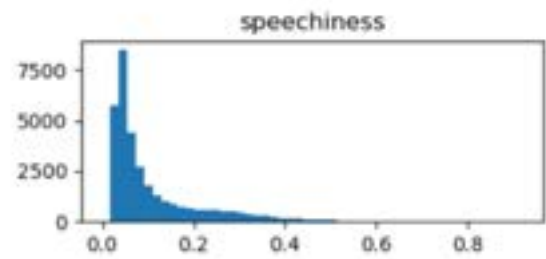
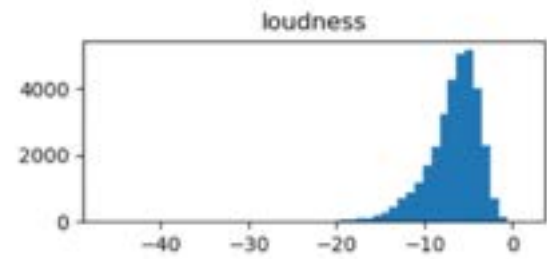
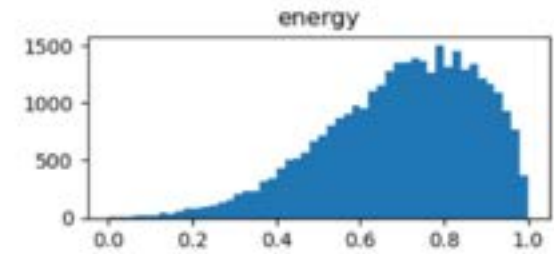
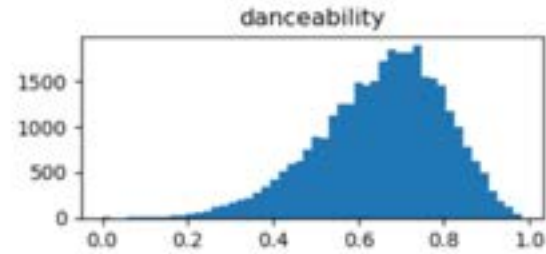
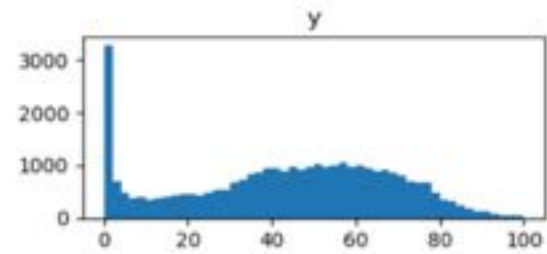
Median 45

Data Wrangling & EDA Pt. 2

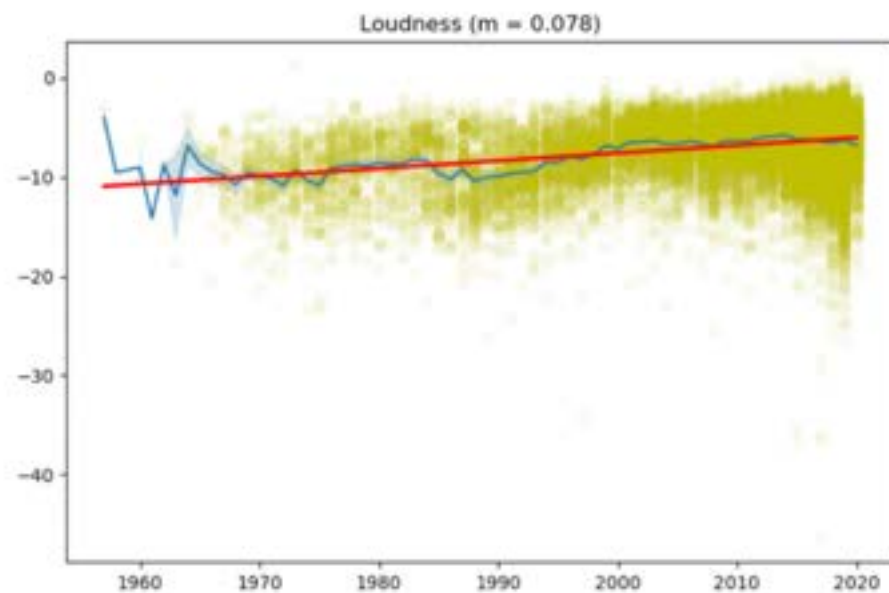
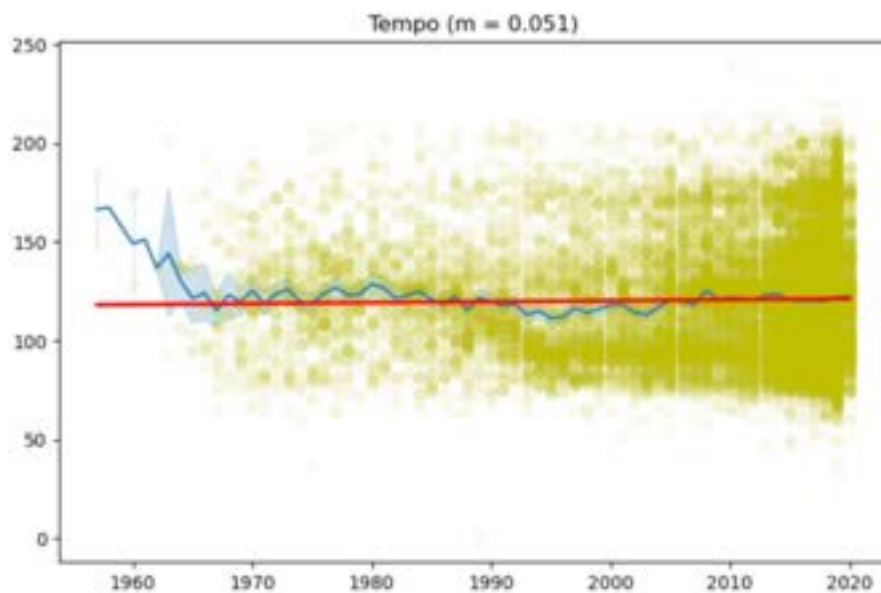
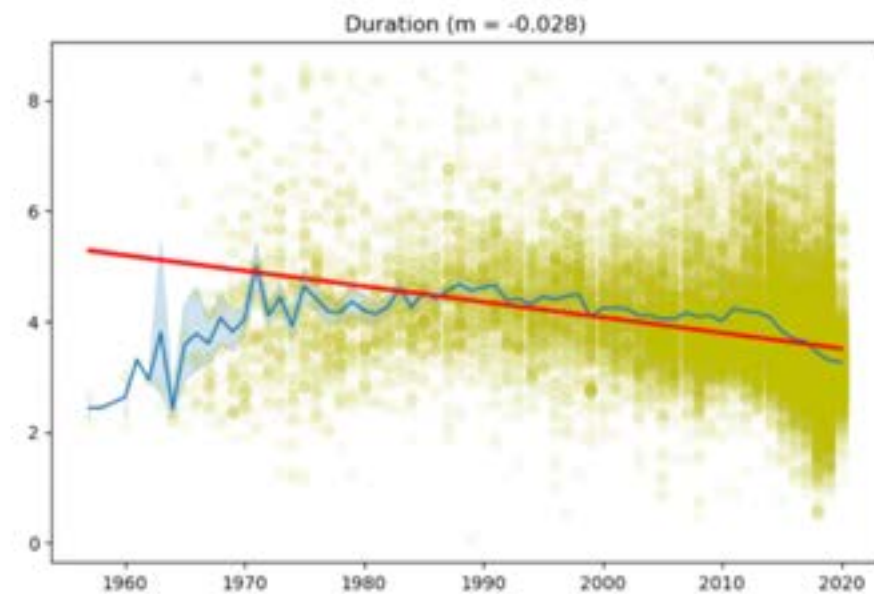
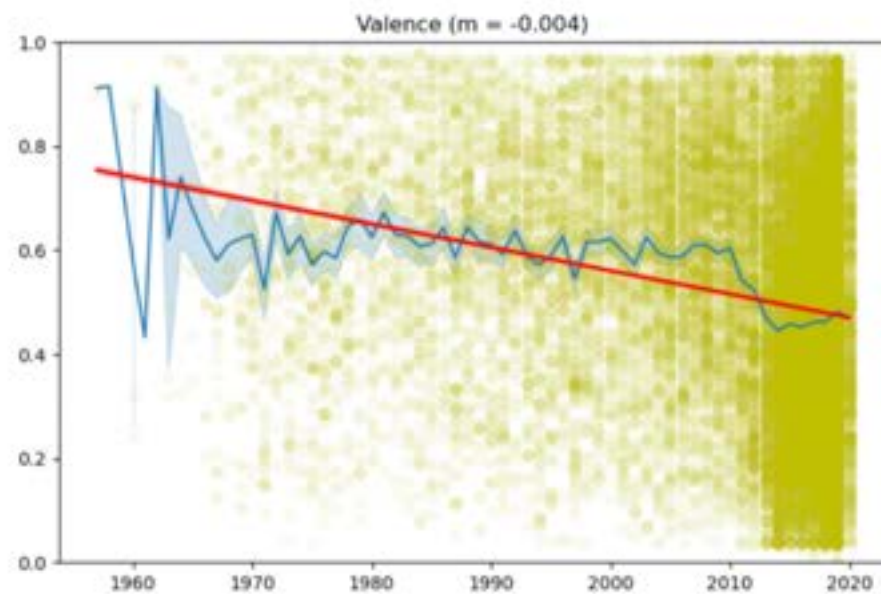
Features

- Tempo
- Danceability
- Valence
- Energy
- Acousticness
- Instrumentalness
- Liveness
- Speechiness
- **Duration**
- **Loudness**

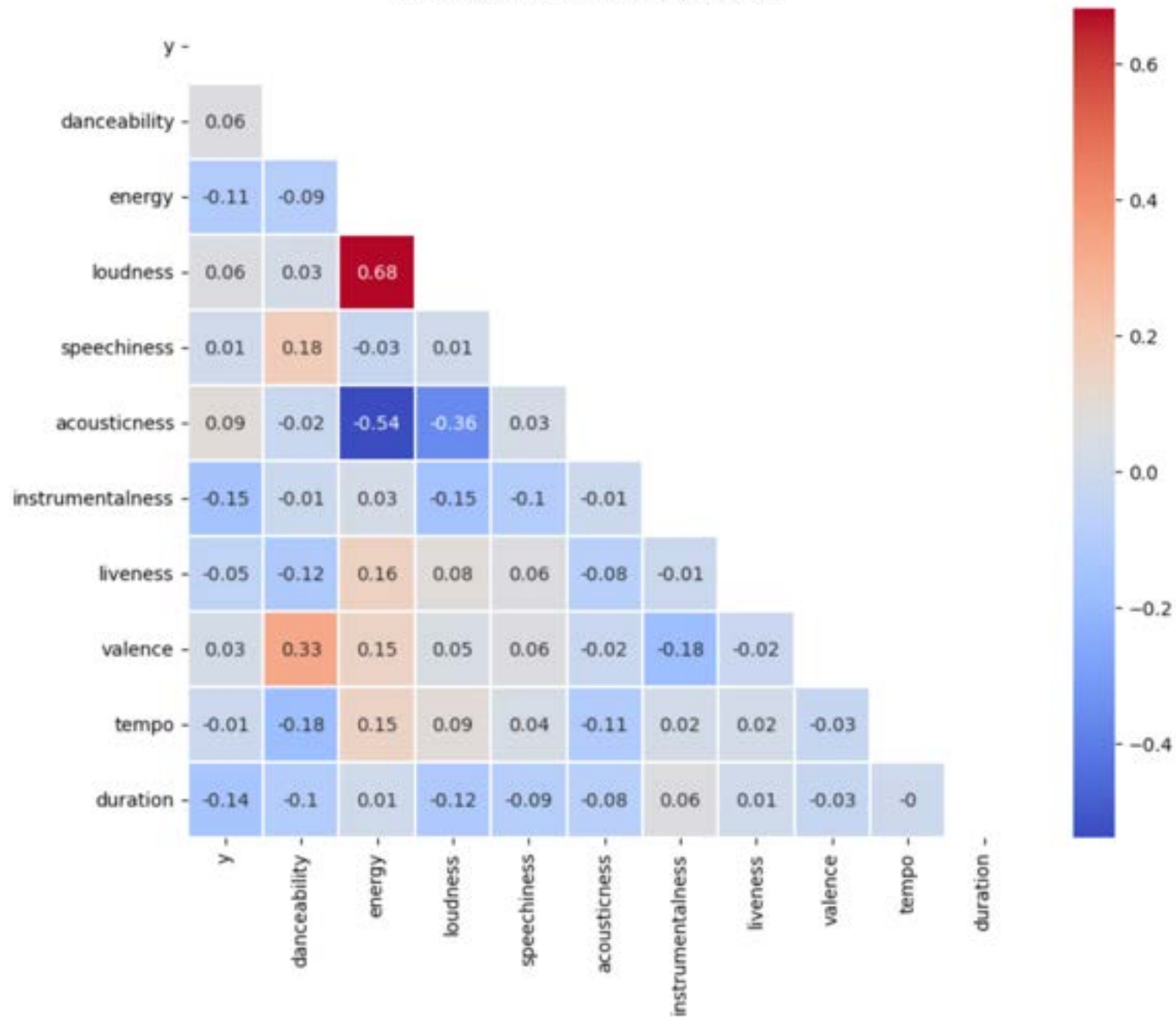
Feature Distributions (Dataset 2)



Music Feature Changes Over Time (1957 – 2020)



Feature Correlations Heatmap (Dataset 2)



Modeling Pt. 2

Random Forest Regressor

- Standard Scaler
- CV / Test MAE 16.86 / 16.31

Extra Trees Regressor

- Power Transformer
- CV / Test MAE 15.94 / 15.30

Gradient Boosting Regressor

- Min Max Scaler
- CV / Test MAE 16.83 / 16.00

Hist Gradient Boosting Regressor

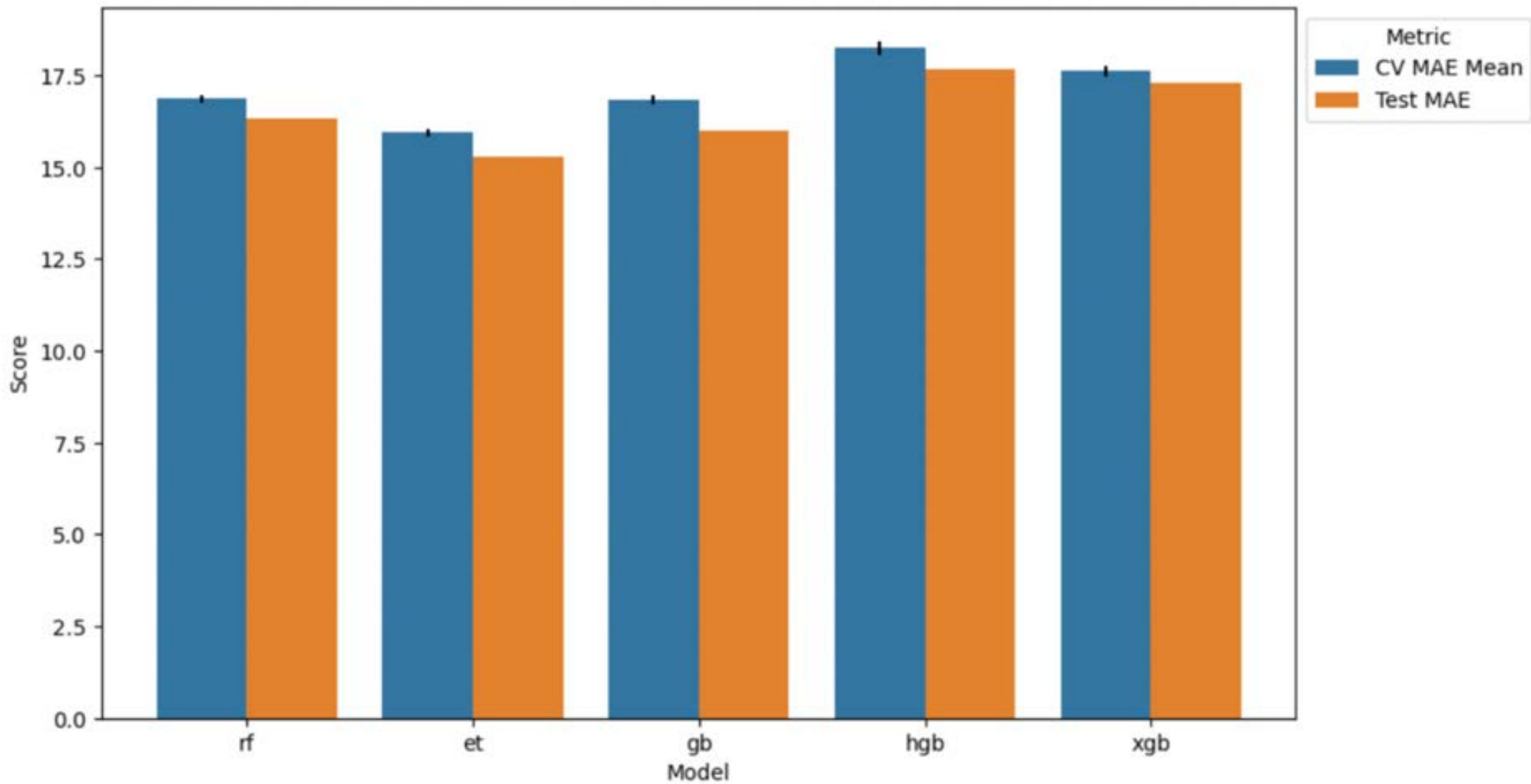
- Robust Scaler
- CV / Test MAE 18.24 / 17.68

XGBoost Regressor

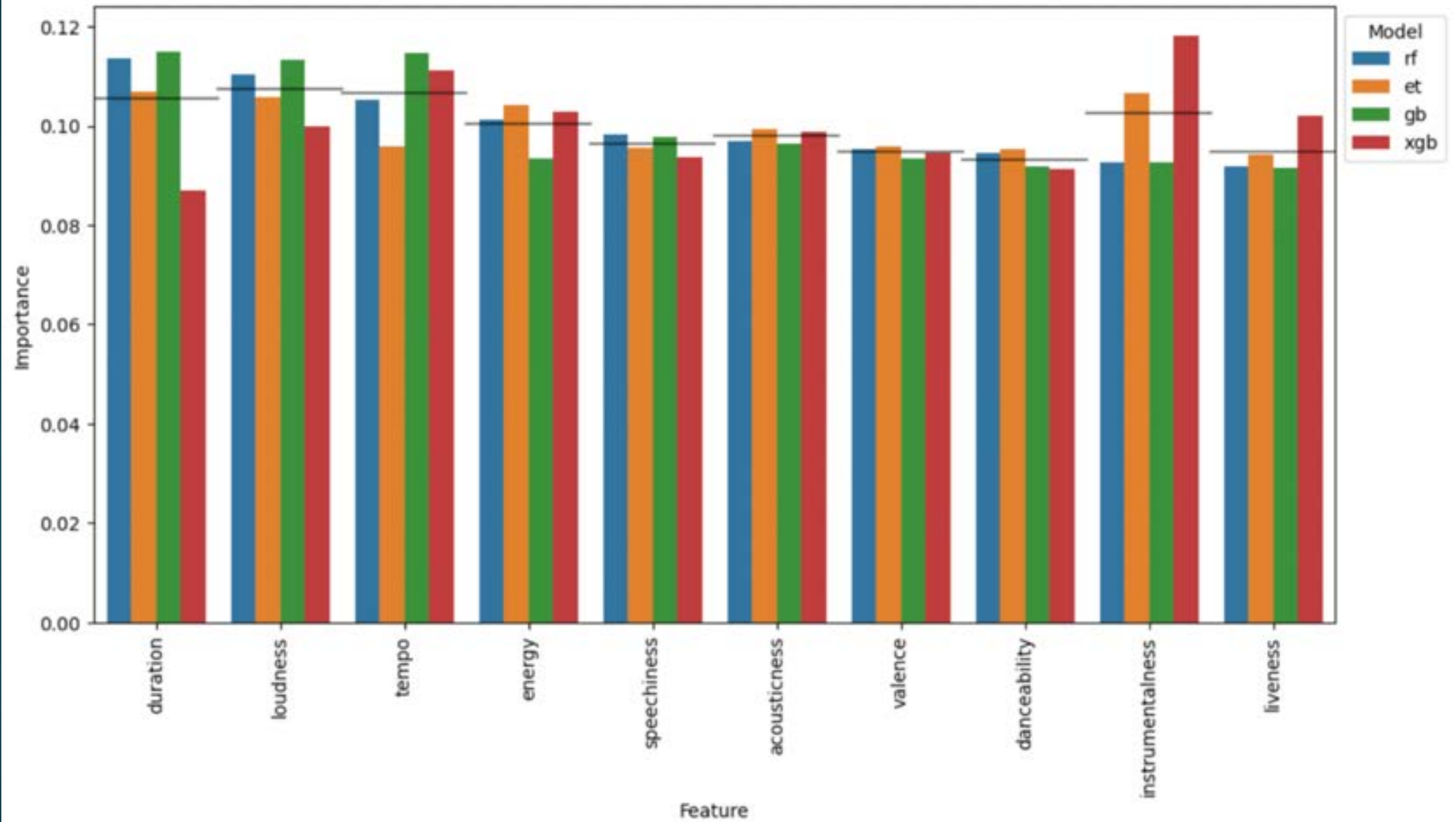
- Power Transformer
- CV / Test MAE 17.62 / 17.29

(80 / 20 train test split)

Model Performance Comparison



Feature Importance Comparison by Model



Modeling Pt. 2

Extra Trees Regressor

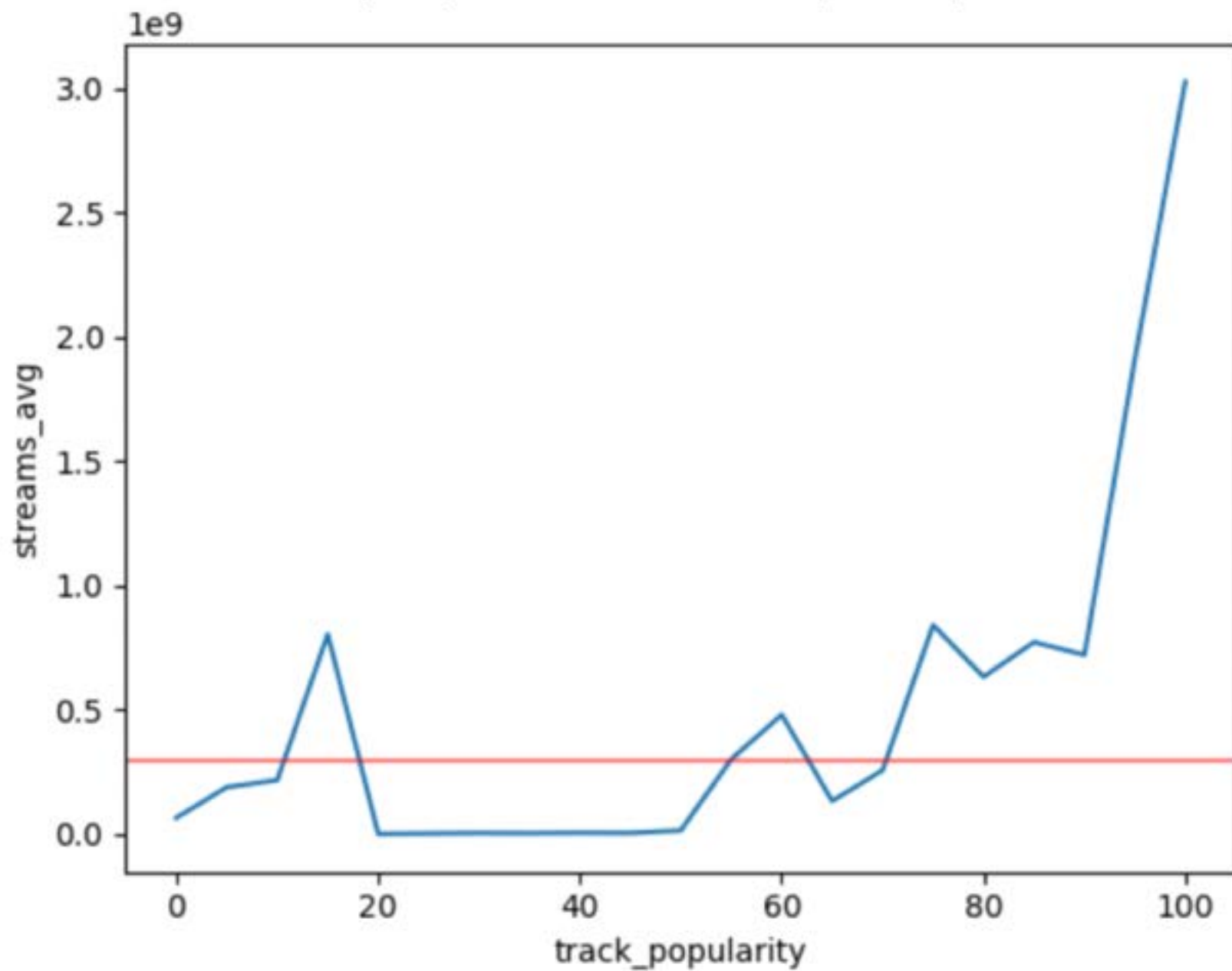
Hyperparameters

- `n_estimators` = 1458
- `max_depth` = 46
- `random_state` = 42

Feature Importance

1. **Duration (0.106840)**
2. **Instrumentalness (0.106628)**
3. **Loudness (0.105952)**
4. **Energy (0.104153)**
5. Acousticness (0.099422)
6. Tempo (0.095905)
7. Valence (0.095804)
8. Speechiness (0.095739)
9. Danceability (0.095383)
10. Liveness (0.094174)

Popularity Score and Stream Count (Dataset 2)



Conclusion

Streams / Popularity Score

Random Forest Regressor (Set 1)

- Test MAE **268M**
- Tempo, Danceability, Acousticness, Valence

Extra Trees Regressor (Set 2)

- Test MAE **15.30**
- Duration, Instrumentalness, Loudness, Energy

Future Work

- Use current models for baseline predictions
- Seek dataset with thorough observations and features
- Quantify success metric