

Spotify Popularity Prediction

Robert Bull, Uriel Garcia, Ahsin Saleem, Ross Vrbanac
DS w207- Tanya Roosta
Dec. 12th, 2024





Table of Contents

1. Data Source
2. Research Question
3. Feature Overview
4. Data Processing
5. Genre 1 - Death Metal
6. Genre 2 - Hip Hop
7. Genre 3 - Salsa
8. Genre 4 - Piano
9. Conclusion



Data Source:

- Switched our original dataset from crop yield to Spotify popularity prediction.
- **New data:**
 - Kaggle Spotify 1 million tracks dataset.
 - Extracted from the Spotify using their API
 - <https://www.kaggle.com/datasets/amitanshjoshi/spotify-1million-tracks/data>
- **Key Highlights:**
 - Size: ~1 Million tracks.
 - Features: 19 musical and metadata features.
 - Artists: 61,445 unique artists.
 - Genres: 82 distinct genres.



Research Question:

What is the best performing model for classifying Spotify song popularity per genre, irrespective of artist?

Sub-Question: What features define each genre?

Sub-Question: How well can the popularity per song of different genres be predicted?





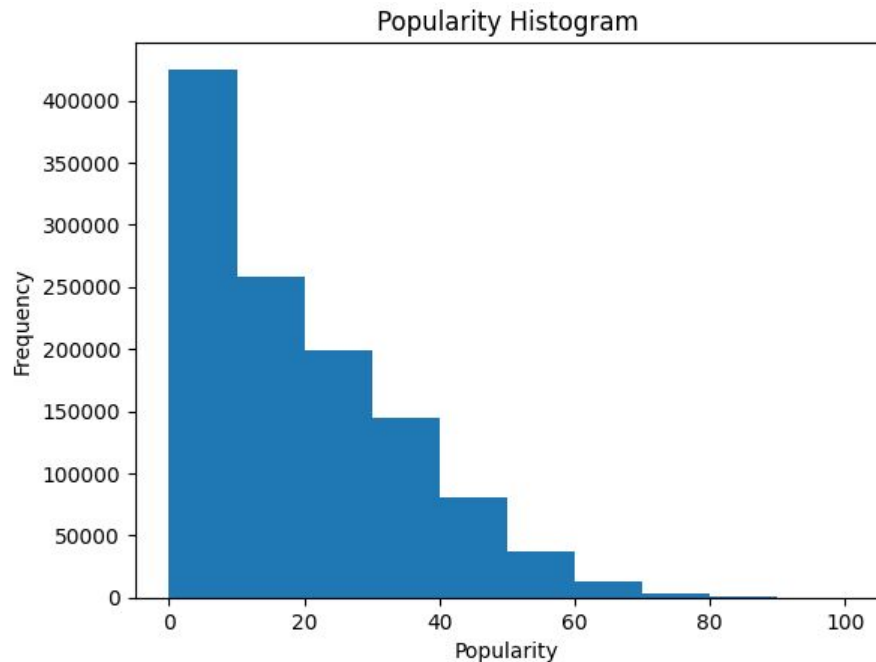
Features:

Audio Features	Description
Popularity	Track popularity (0 to 100)
Year	Year released (2000 to 2023)
Danceability	Track suitability for dancing (0.0 to 1.0)
Energy	The perceptual measure of intensity and activity (0.0 to 1.0)
Key	The key, the track is in (-1 to -11)
Loudness	Overall loudness of track in decibels (-60 to 0 dB)
Mode	Modality of the track (Major '1' / Minor '0')
Speechiness	Presence of spoken words in the track
Acousticness	Confidence measure from 0 to 1 of whether the track is acoustic
Instrumentalness	Whether tracks contain vocals. (0.0 to 1.0)
Liveness	Presence of audience in the recording (0.0 – 1.0)
Valence	Musical positiveness (0.0 to 1.0)
Tempo	Tempo of the track in beats per minute (BPM)
Time_signature	Estimated time signature (3 to 7)
Duration_ms	Duration of track in milliseconds



Data Processing

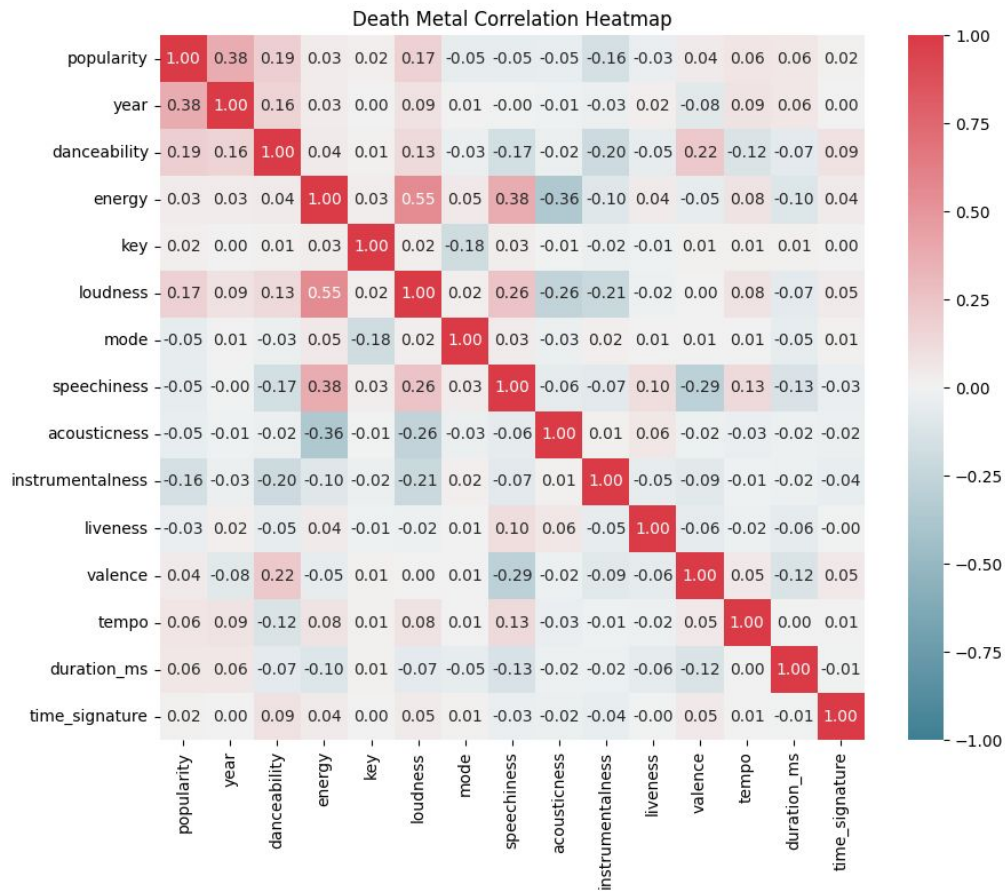
- Popularity separated into binary classes
 - 0 for popularity < median popularity per genre
 - 1 for popularity \geq median popularity per genre
- Min-max standardization of numeric features
- One-hot encoding of categorical features
- 60/20/20 train/validate/test splits with shuffled data



Death Metal



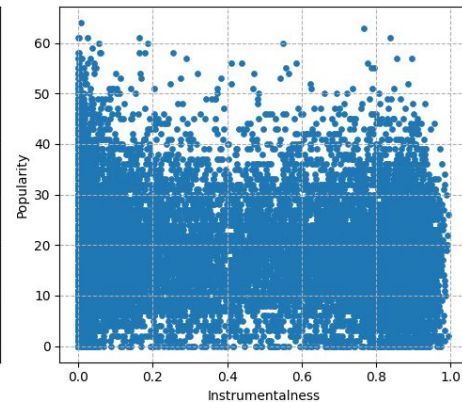
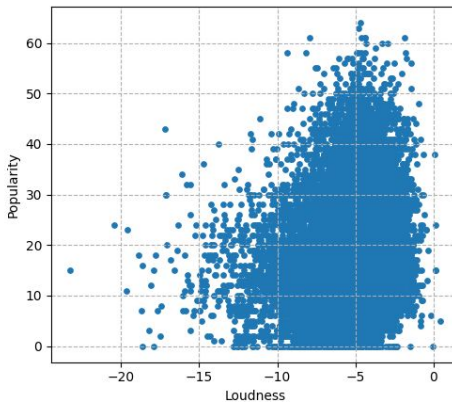
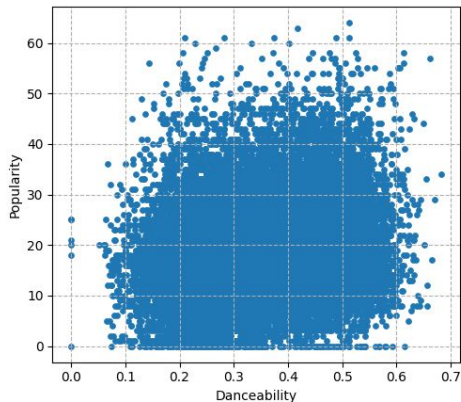
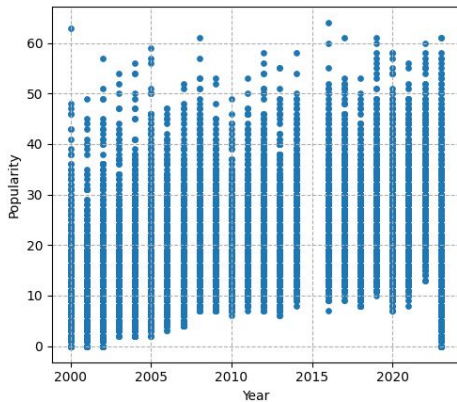
Death Metal: EDA





Death Metal: EDA

Popularity vs Key Features

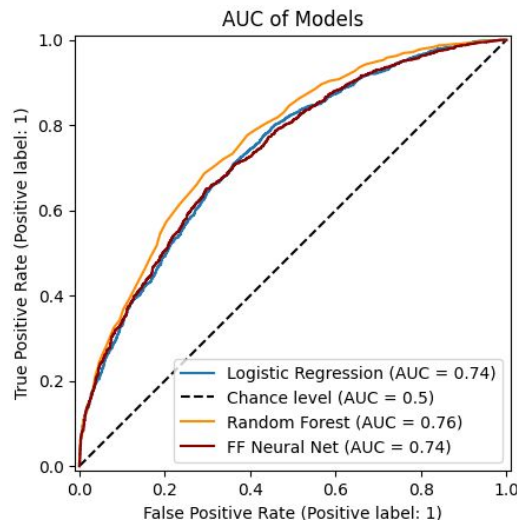




Death Metal: Modeling

Model	Accuracy	Precision	AUC Score
Baseline	0.52	0.52	0.5
Logistic Regression	0.67	0.67	0.74
Random Forest Classifier	0.69	0.69	0.76
FF Neural Net	0.67	0.68	0.74

- Comparisons on validation set
- Baseline - majority class classifier
- Logistic Regression - L2 penalization
- Random Forest Classifier - no maximum depth, splits based on Gini impurity
- FF Neural Net - 4 hidden layers of sizes [204, 204, 102, 102], 0.1 dropout layer for each, and tuned hyperparameters



Death Metal: Results (Random Forest)

	Feature	Feature Importance
0	danceability	0.083296
1	loudness	0.081751
2	instrumentalness	0.076436
3	duration_ms	0.072134
4	tempo	0.069519
5	speechiness	0.069384
6	acousticness	0.068290
7	valence	0.068242
8	liveness	0.065552
9	energy	0.063458

Death Metal RFC Test Accuracy:

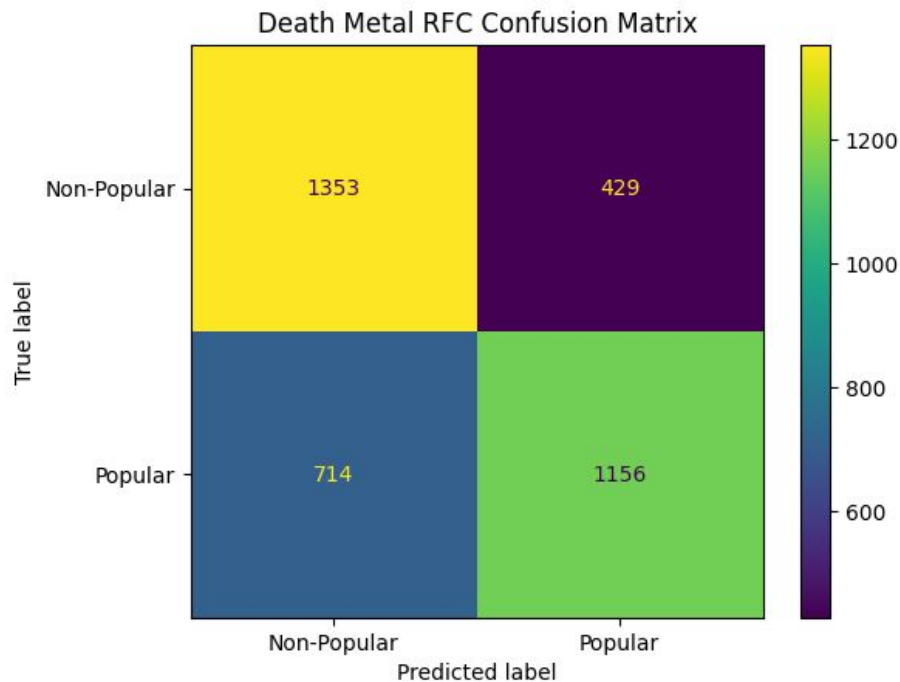
0.7009423503325942

Death Metal RFC Test Precision:

0.7036852589641435

Death Metal RFC Test AUC Score:

0.7622558858969235





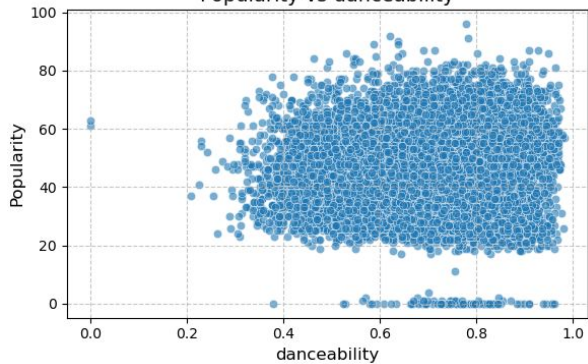
Hip-Hop



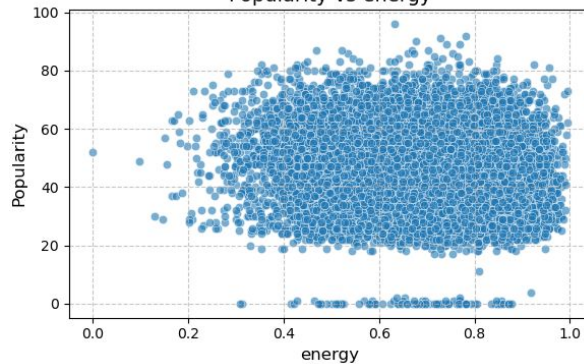
Hip-Hop: EDA



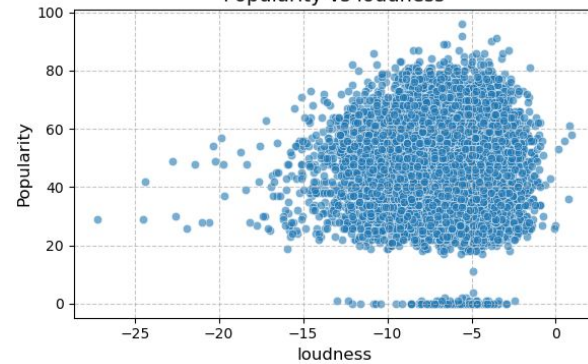
Popularity vs danceability



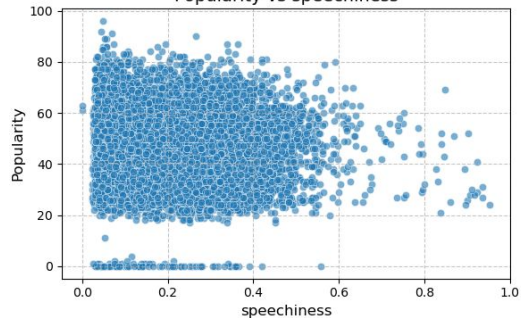
Popularity vs energy



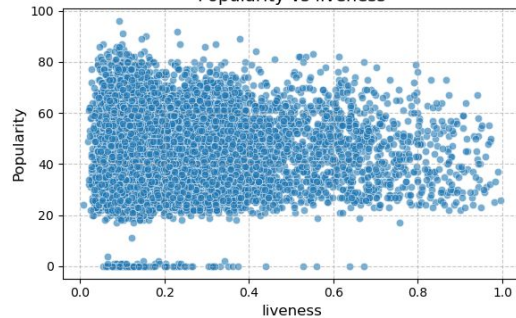
Popularity vs loudness



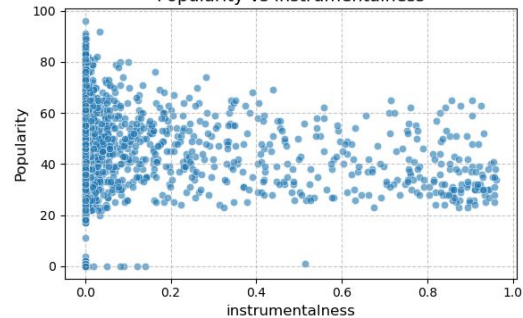
Popularity vs speechiness



Popularity vs liveness



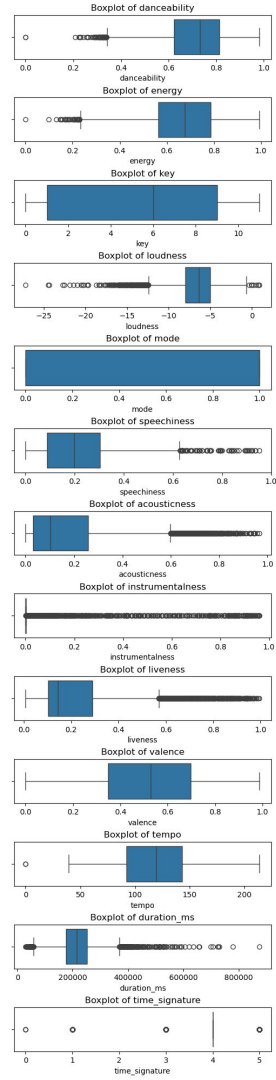
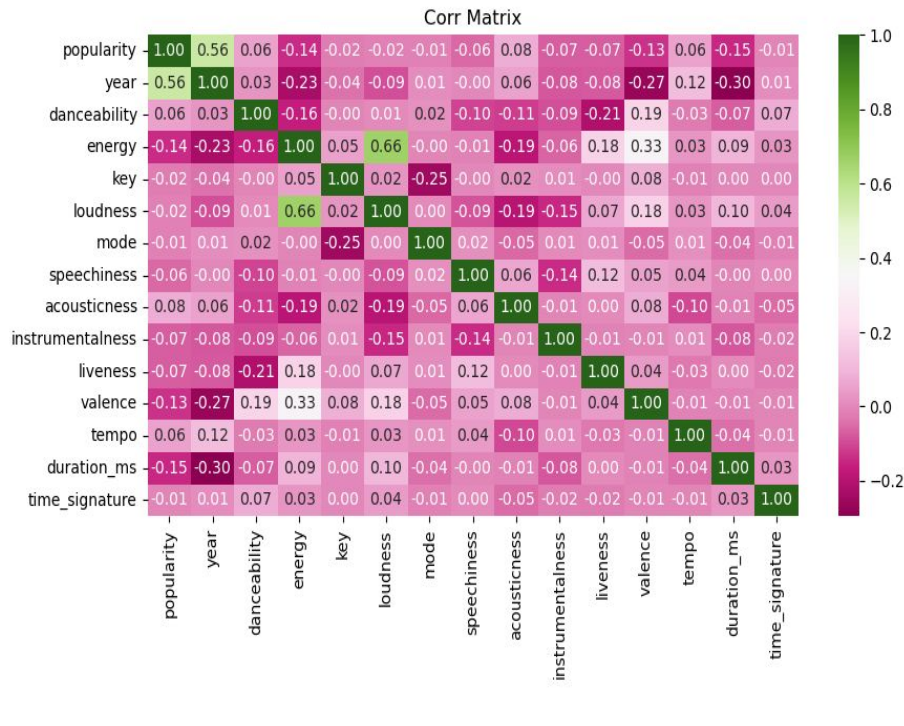
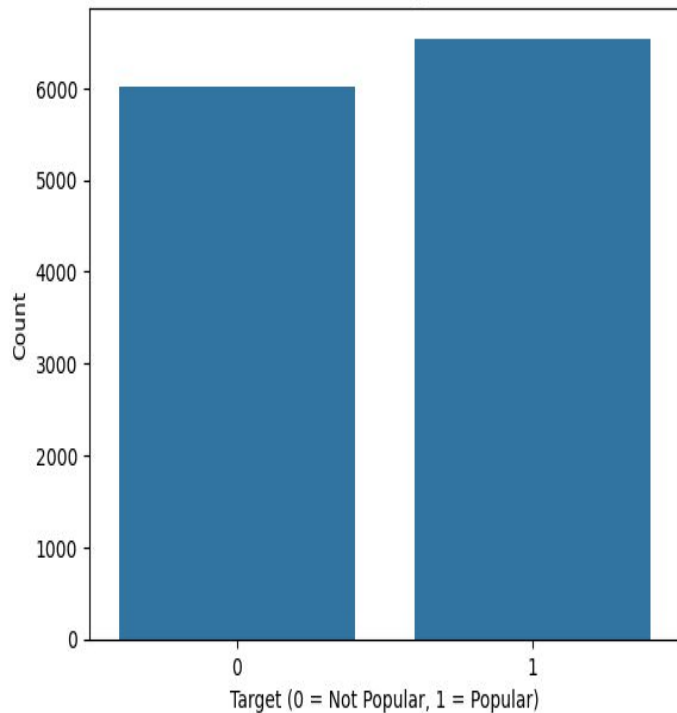
Popularity vs instrumentalness



Hip-Hop: EDA



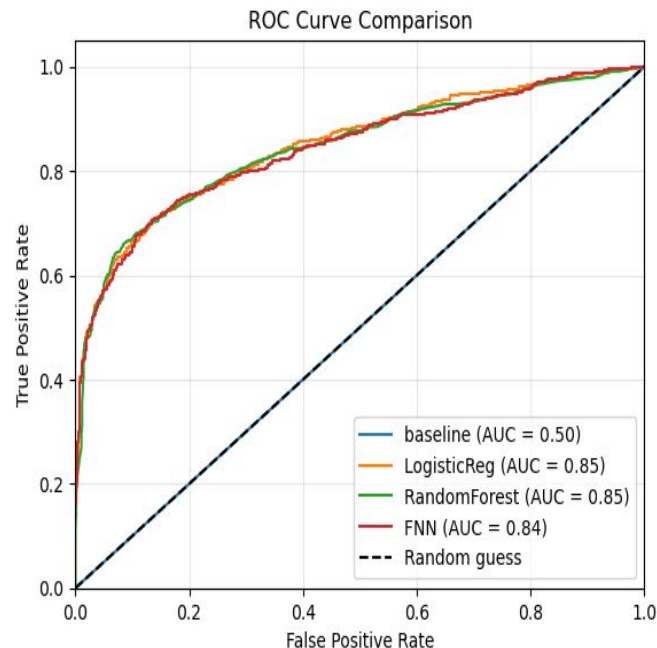
Distribution of Target Variable



Hip-Hop: Modeling

Model	Accuracy	Precision	AUC Score
Baseline	0.52	0.52	0.5
Logistic Regression	0.78	0.83	0.85
Random Forest Classifier	0.78	0.85	0.85
FF Neural Net	0.77	0.81	0.84

- Comparisons on validation set
- Baseline - majority class classifier
- Logistic Regression - L1 penalization
- Random Forest Classifier - no maximum depth, splits based on Gini impurity
- FF Neural Net - 4 hidden layers of sizes [256, 128, 64, 32], 0.4 dropout layer for each, and tuned hyperparameters

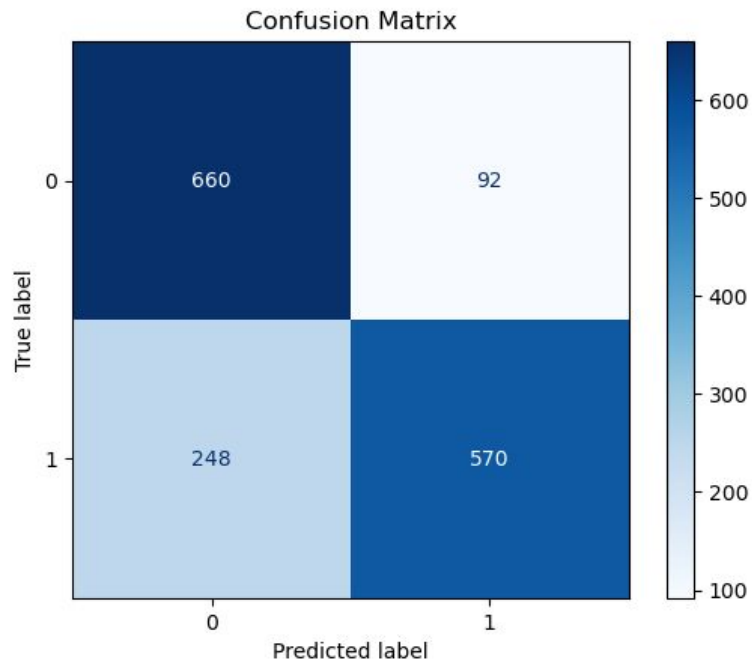


Hip-Hop: Results (Random Forest)

	Feature	Coefficient
8	duration_ms	0.075731
6	valence	0.061783
1	energy	0.061680
4	acousticness	0.060087
7	tempo	0.059545
3	speechiness	0.058403
0	danceability	0.058317
5	liveness	0.057761
2	loudness	0.057374
30	year_2022	0.050634

Final Results:

- Test Accuracy: 0.78
- Test Precision: 0.86
- Test_AUC: 0.78

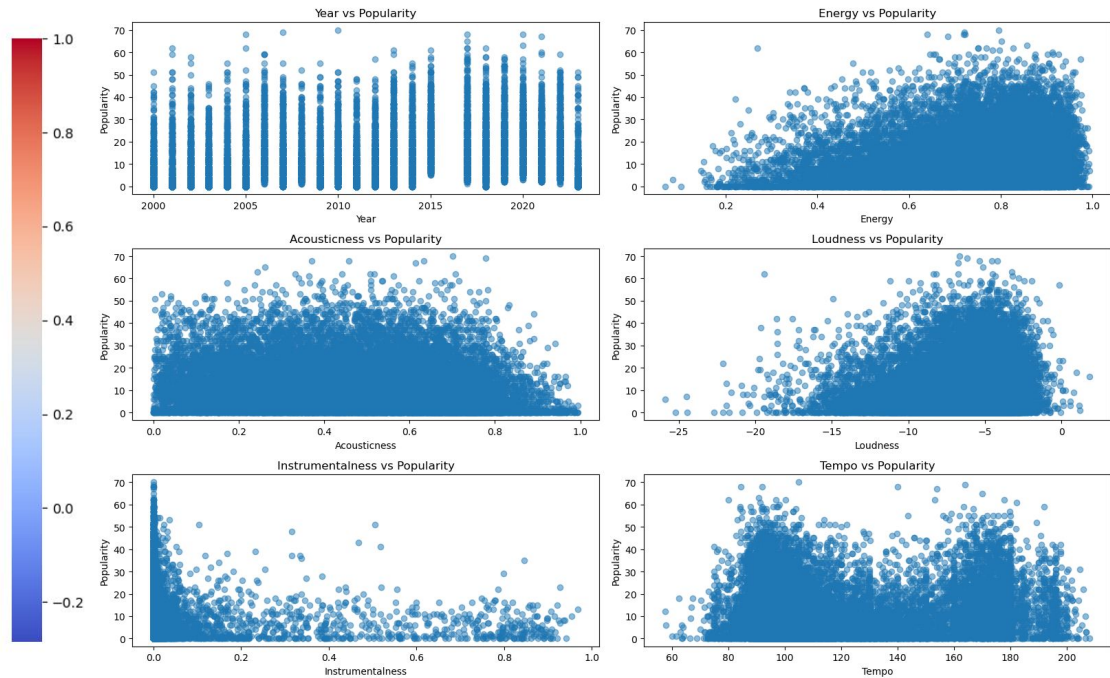
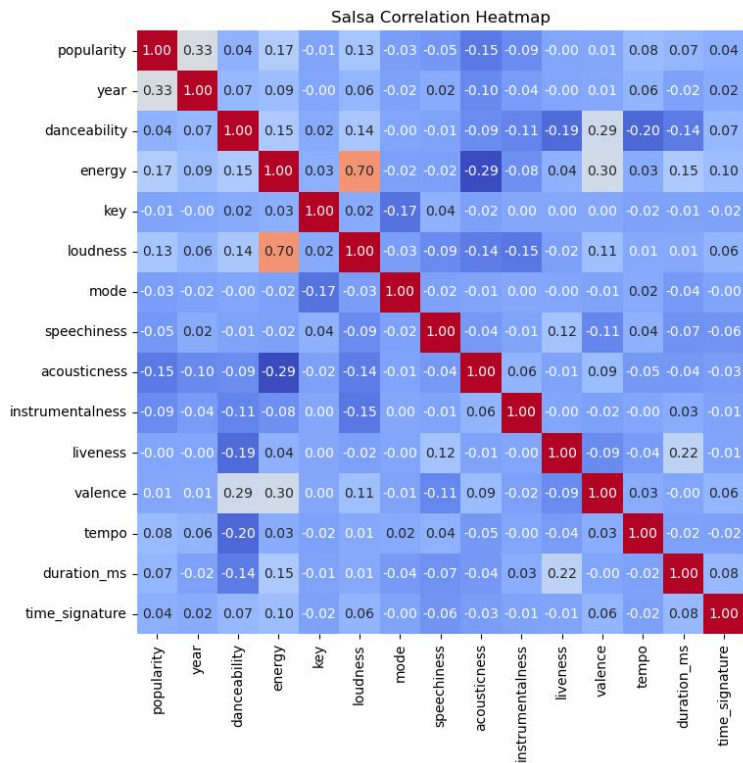




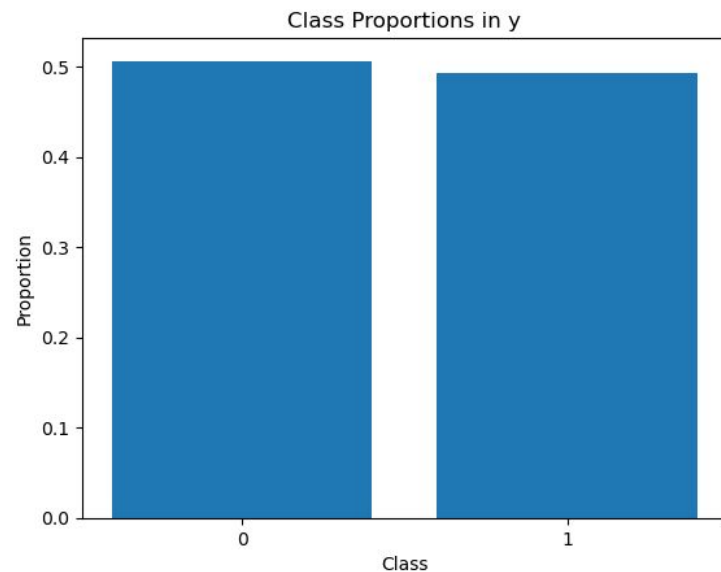
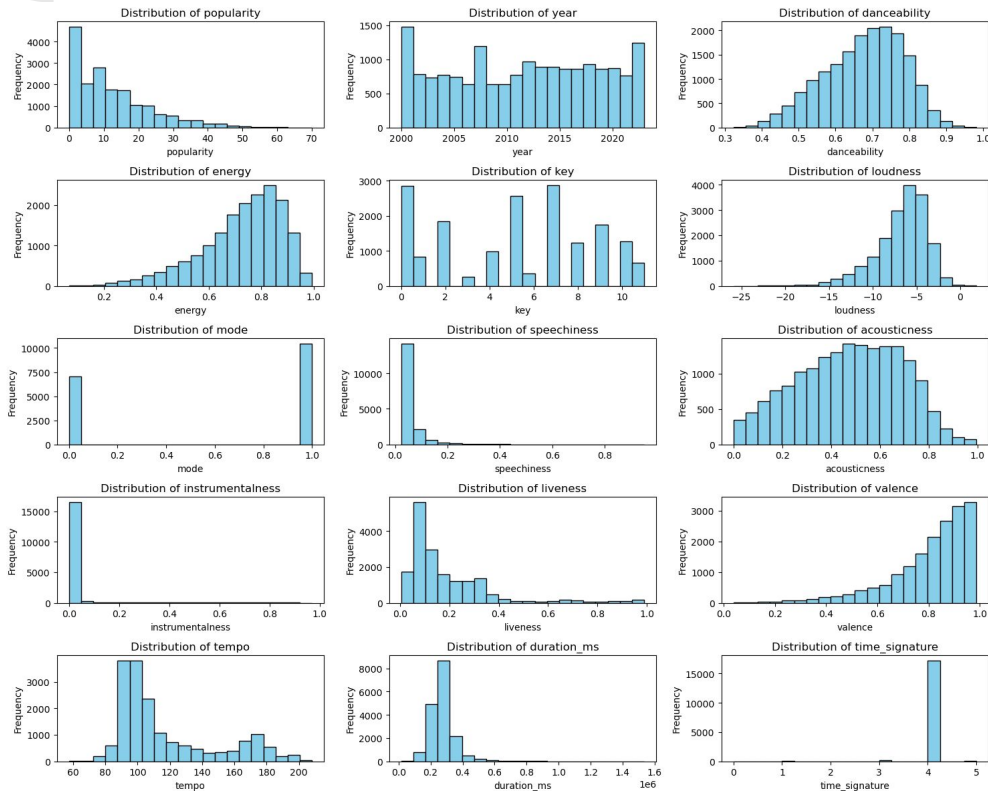
Salsa



Salsa: EDA



Salsa: EDA (part 2)

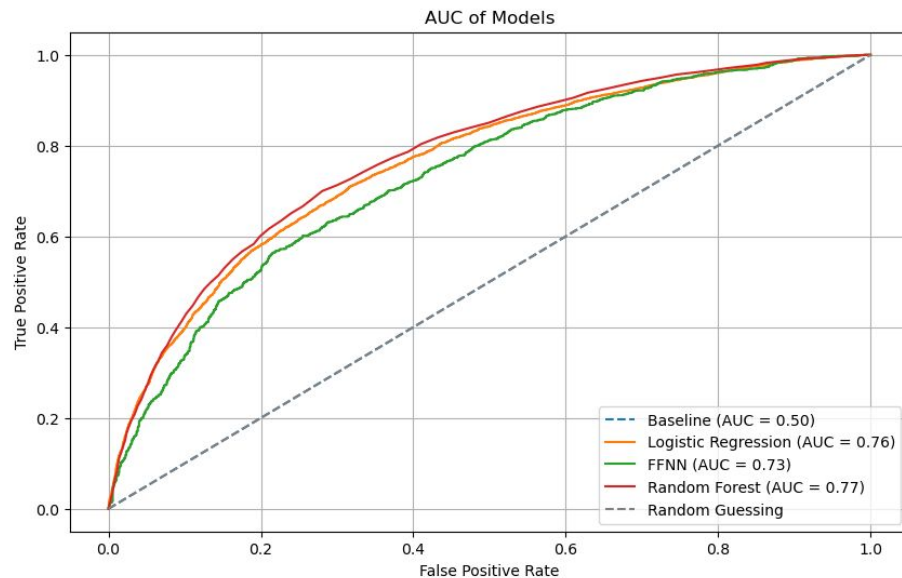




Salsa: Modeling

- Comparisons on validation set
- Baseline - majority class classifier
- Logistic Regression - L1 penalization
- Random Forest Classifier - no maximum depth, splits based on Gini impurity
- FF Neural Net - 2 hidden layers of sizes [128, 64], 0.3 dropout layer and tuned hyperparameters

Model	Accuracy	Precision	AUC
Baseline	0.50	0.50	0.50
Logistic Regression	0.70	0.70	0.76
FNN	0.67	0.68	0.73
Random Forest	0.70	0.71	0.77

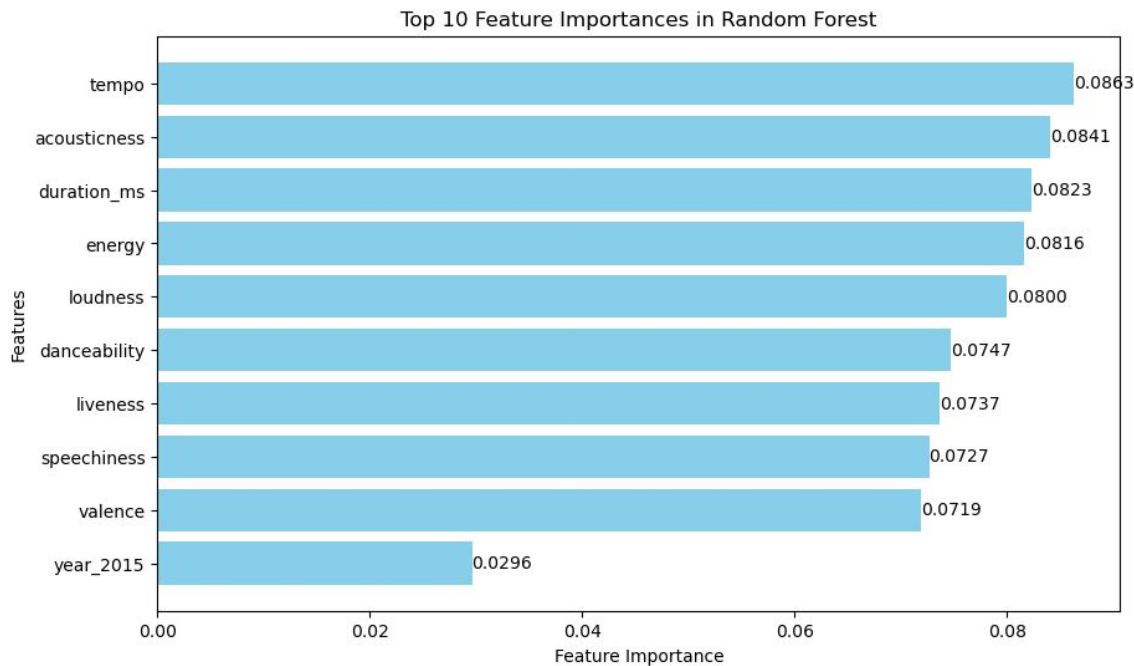
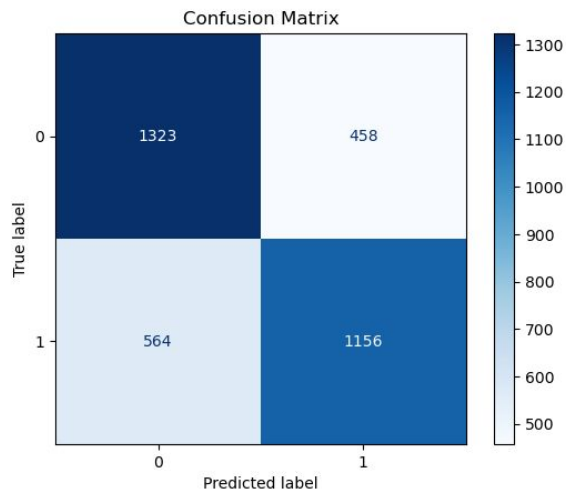




Salsa: Results(Random Forest)

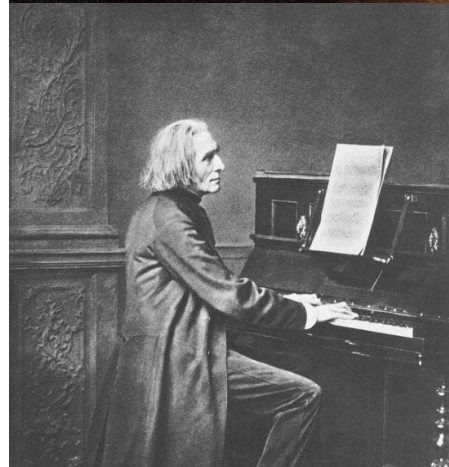
Final Results:

- Test Accuracy: 0.70
- Test Precision: 0.71
- Test AUC Score: 0.77

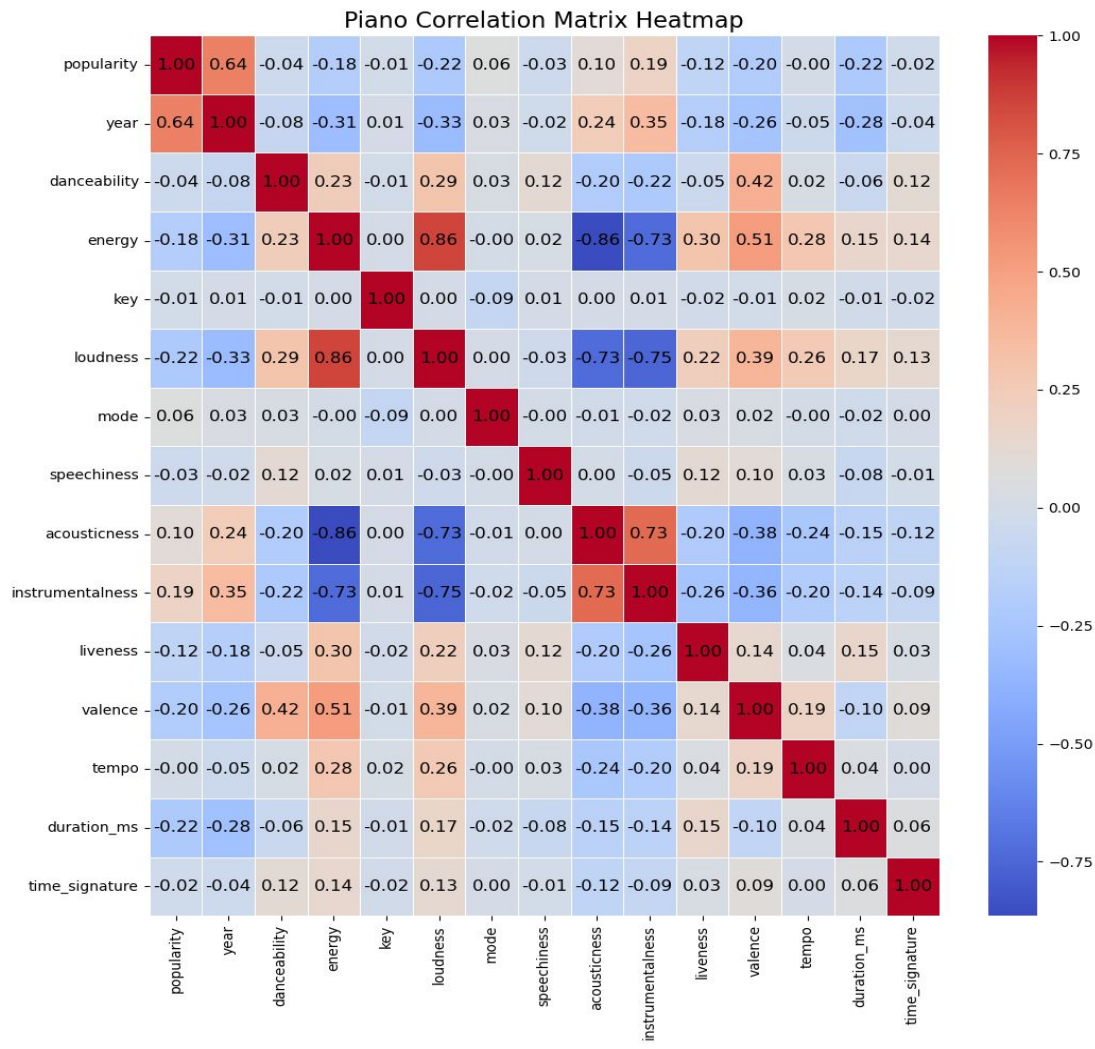




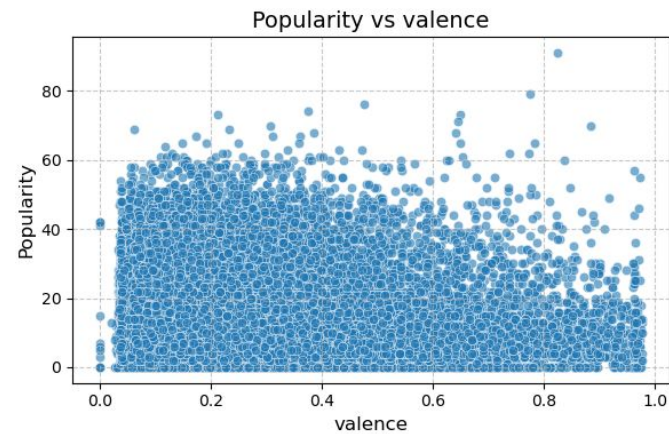
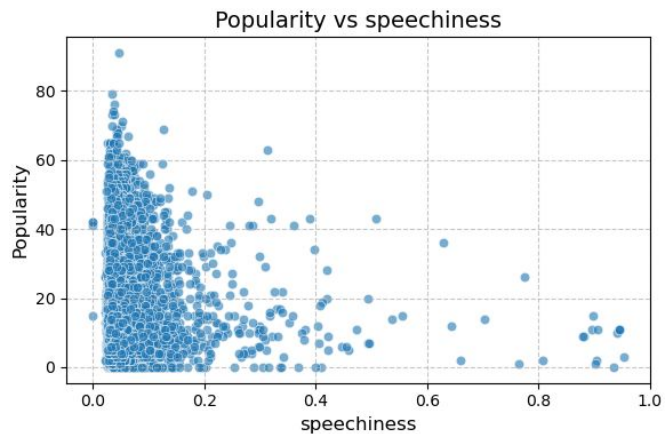
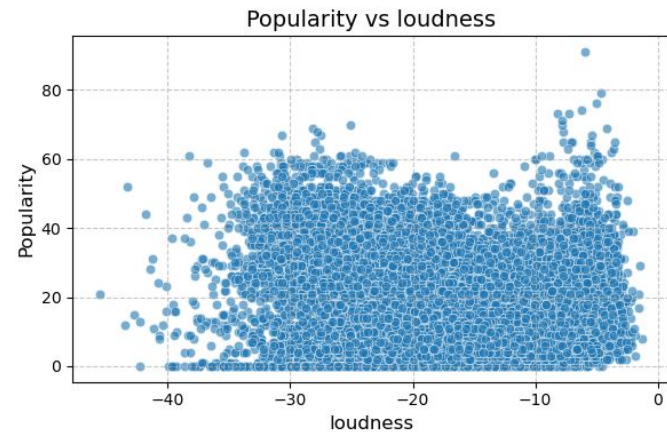
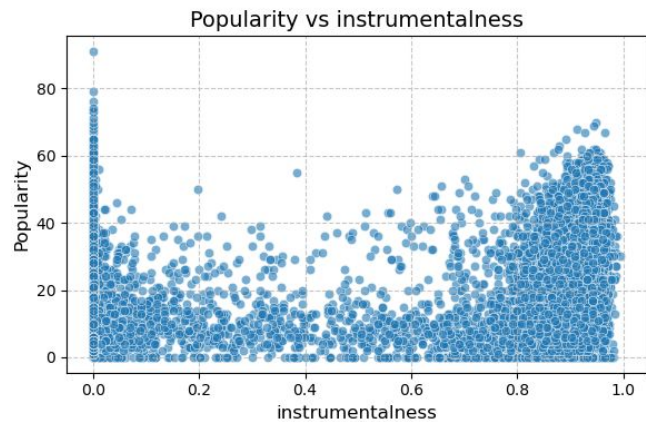
Piano



Piano: EDA



Piano: EDA

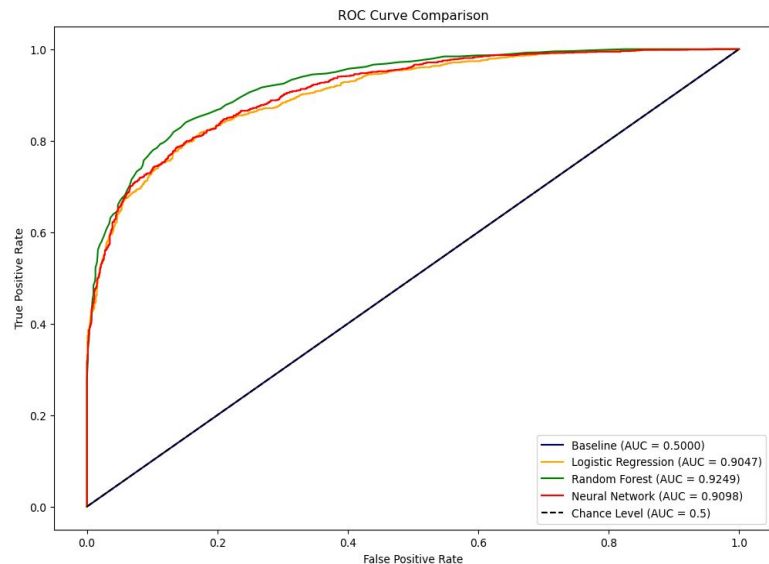




Piano: Modeling

Model	Accuracy	Precision	AUC
Baseline	0.52	0.52	0.50
Logistic Regression	0.82	0.87	0.90
Random Forest Classifier	0.84	0.88	0.92
FF Neural Net	0.82	0.84	0.91

- Comparisons on validation set
- Baseline - majority class classifier
- Logistic Regression - L1 penalization
- Random Forest Classifier - no maximum depth, splits based on Gini impurity
- FF Neural Net - 2 hidden layers of sizes [64, 64], 0.2 dropout layer for each, and tuned hyperparameters. Tuned.





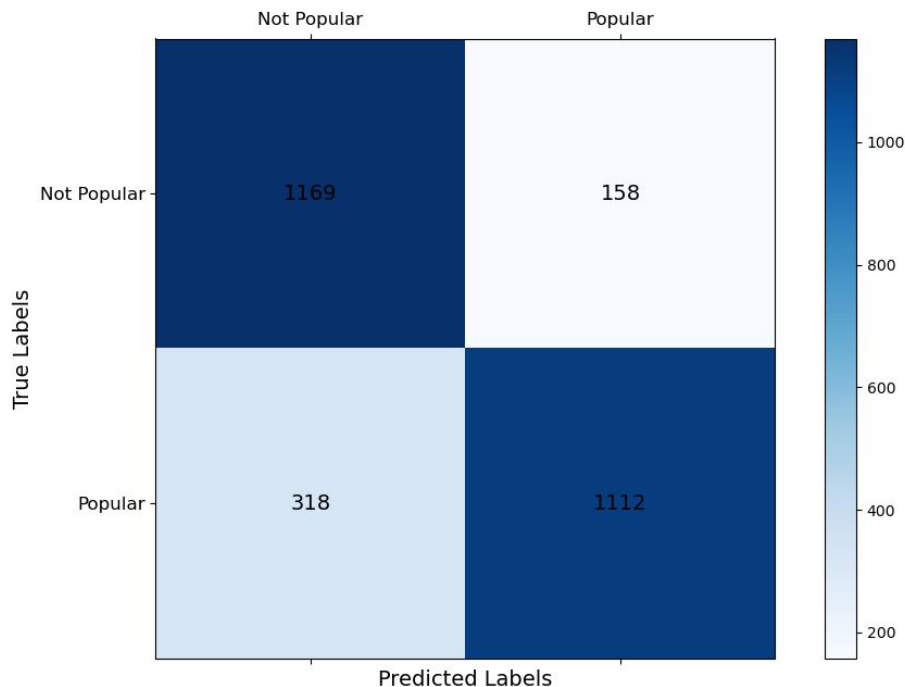
Piano: Results (Random Forest)

	Feature	Importance
0	instrumentalness	0.072622
1	loudness	0.063618
2	duration_ms	0.063003
3	valence	0.060654
4	energy	0.060022
5	acousticness	0.056234
6	liveness	0.047553
7	speechiness	0.045470
8	tempo	0.044045
9	danceability	0.043561
10	year_2020	0.043095

Final Results:

- Test Accuracy: 0.83
- Test Precision: 0.88
- Test AUC: 0.92

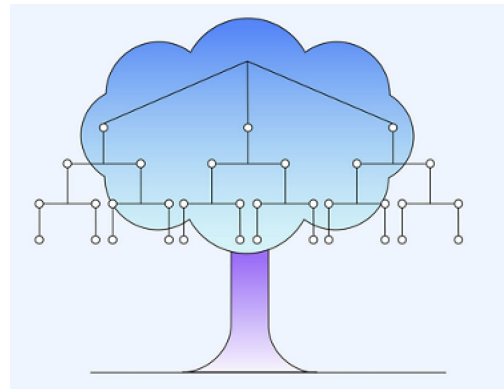
Piano Confusion Matrix (Random Forest, Test Data)





Conclusion

- Random Forest performed the best on all genres
 - Handles irrelevant features well
- Feature Importance varied across genres
 - Loudness varied significantly



Thank you for listening to this presentation today!



Contributions

Name	Presentation Contribution	Code Contribution
Robert Bull	Data Features, Data Processing, Death Metal	All modeling and EDA with respect to the Death Metal genre
Uriel Garcia	Motivation, Salsa	All modeling and EDA with respect to the Salsa genre
Ahsin Saleem	Introduction, Hip-Hop	All modeling and EDA with respect to the Hip-Hop genre
Ross Vrbanac	Conclusion, Piano	All modeling and EDA with respect to the Piano genre



Github

<https://github.com/rfbull/mids-w207-final-project>