# Spotify Billboard Classifier
## Course Project: DATA 1030 - Hands On Machine Learning

Sayan Samanta

Instructor: Dr. Andras Zsom
TA Advisor: Natalie Delworth
GitHub Repo: shorturl.at/BEN08

Figure: Sample data with Features

Figure: The 1st two principal component



Figure: Balance of the dataset

Figure: Correlation of features

Figure: Valence evolution with time



Figure: Energy evolution with time



Figure: Acousticness evolution with time

1 Since the dataset is highly imbalanced, we use stratified KFold split.



Figure: Different KFold Splits

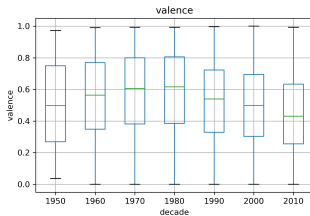## Random Forest

Parameters tuned:

| Parameter Name | Parameter Range | Optimal Parameter |
|---|---|---|
| max_depth | $1, 2 \ldots 10$ | 10 |
| min_sample_split | $1, 2 \ldots 10$ | 3 |

## XGBoost

Parameters tuned:

| Parameter Name | Parameter Range | Optimal Parameter |
|---|---|---|
| max_depth | 3, 4, 5, 6, 8, 10, 12, 15 | 15 |
| min_child_weight | 1, 3, 5, 7 | 3 |
| gamma | 0.0, 0.1, 0.2 , 0.3, 0.4 | 0.4 |

## AdaBoost

Parameters tuned:

| Parameter Name | Parameter Range | Optimal Parameter |
|---|---|---|
| learning_rate | $10^{-3}, \ldots, 10^{4}$ | 0.1 |
| algorithm | SAMME.R and SAMME | SAMME.R |

Figure: Confusion Matrix for Random Forest



Figure: Feature Importance for Random Forest

Figure: Confusion Matrix for XGBoost



Figure: Feature Importance for XGBoost

Figure: Confusion Matrix for AdaBoost



Figure: Feature Importance for AdaBoost

Figure: Estimator with least error

- The dataset was highly imbalanced. Having a much more evenly divided set among the classes would improve the results.
- The acoustic features of songs in both the classes overlap highly. Different data mining algorithm might do better
- Due to lack of time, certain methods such as support vector machine classification or K-nearest neighbours couldnt be implemented
- Advanced deep-learning methods could improve the result.

# Thank You. Question?