

# Spotify Billboard Classifier

## Course Project - DATA1030

SAYAN SAMANTA

School of Engineering | Brown University

*Instructor: Dr. Andras Zsom | TA: Natalie Delworth*

*github*

December 3, 2019

## Contents

|          |   |           |
|----------|---|-----------|
| <b>1</b> | <b>Introduction</b>                               | <b>2</b>  |
| 1.1      | Dataset Description: . . . . .                    | 2         |
| 1.2      | Objectives . . . . .                              | 3         |
| 1.2.1    | Target Variable . . . . .                         | 3         |
| 1.2.2    | Motivation behind the problem . . . . .           | 3         |
| <b>2</b> | <b>Exploratory Data Analysis</b>                  | <b>3</b>  |
| 2.1      | Feature Correlation and Dataset Balance . . . . . | 3         |
| <b>3</b> | <b>Machine Learning Pipeline</b>                  | <b>5</b>  |
| 3.1      | Preprocessing . . . . .                           | 5         |
| 3.2      | Cross Validation Pipeline . . . . .               | 6         |
| 3.3      | Model Selection . . . . .                         | 6         |
| 3.3.1    | Random Forest Classification . . . . .            | 7         |
| 3.3.2    | XGBoost Classification . . . . .                  | 8         |
| 3.3.3    | AdaBoost Classification . . . . .                 | 9         |
| <b>4</b> | <b>Outlook</b>                                    | <b>10</b> |

### Abstract

In this report, an attempt has been made to predict whether a particular song will appear in the US top 200 billboard based on the songs given it's acoustic features. Song data that appeared on the top 200 billboard has been received from Kaggle while it was appended with songs that didn't appear on the billboard using the Spotify Python API. The data obtained was highly imbalanced, with the acoustic features highly overlapping. The baseline accuracy is at  $\approx 88\%$ . A no. of classification algorithms were performed with stratified k-fold cross validation. Among all the models analysed, only the ones with accuracy above baseline are discussed here. These include Random Forest, XGBoost and Adaboost. The entire machine learning pipeline including data description, exploratory data analysis, cross-validation, model inspection and feature importance have been performed and the results have been discussed with measures of improvement for future studies.

# 1 Introduction

The dataset under consideration is the **The Billboard 200 acoustic data**<sup>1</sup> which encompasses the entire chart from 1963-2019, along with the EchoNest acoustic features of as many songs as available. All the features of the songs were obtained using Spotify's python API<sup>2</sup>

The dataset mentioned above was curated for a piece<sup>3</sup> on the data science analytics website Components. However that piece dealt with the album length alone, there could be plethora of other questions and analysis that can be answered/performed with this data.

Since this dataset contain only songs that appeared on the billboard, further data (songs not on the billboard) was obtained using the Python Spotify API (aptly called Spotipy).

## 1.1 Dataset Description:

The dataset has the following features:

| Features         | Description   |
|------------------|---|
| duration_ms      | The duration of the track in milliseconds   |
| key              | The estimated overall key of the track. If no key was detected, the value is -1.  |
| mode             | Mode indicates the modality (major or minor) of a track.  |
| time_signature   | An estimated overall time signature (beats in each bar) of a track.   |
| acousticness     | A confidence measure from 0.0 to 1.0 of whether the track is acoustic.  |
| danceability     | describes how suitable a track is for dancing based on a combination of musical elements. A value of 0.0 is least danceable and 1.0 is most danceable.          |
| energy           | is a measure from 0.0 to 1.0 and represents a perceptual measure of intensity and activity.   |
| instrumentalness | Predicts whether a track contains no vocals. "Ooh" and "aah" sounds are treated as instrumental in this context. Rap or spoken word tracks are clearly "vocal". |
| liveness         | Detects the presence of an audience in the recording.   |
| loudness         | The overall loudness of a track in decibels (dB).   |
| speechiness      | Speechiness detects the presence of spoken words in a track.  |
| valence          | A measure from 0.0 to 1.0 describing the musical positiveness conveyed by a track.  |
| tempo            | The overall estimated tempo of a track in beats per minute (BPM).   |

Table 1: Acoustic Features Description

In addition, we have the name of the song, name of the artist, the album it appeared on, the date of release and the no. of tracks on the album.

<sup>1</sup><https://www.kaggle.com/snapcrack/the-billboard-200-acoustic-data>

<sup>2</sup><https://spotipy.readthedocs.io/en/latest/>

<sup>3</sup><https://components.one/posts/it-goes-on-album-length/>

## 1.2 Objectives

### 1.2.1 Target Variable

In this project, we would attempt to predict if a song would appear in the billboard or would it not. Hence it is a **classification** problem.

### 1.2.2 Motivation behind the problem

This classifier would help a no. of group both in the music industry or otherwise.

1. Would help artists gauge their music from an unbiased (?) and single metric.
2. Would help producers asses their clients music with a single metric.
3. Would help music listeners with their choices of what to listen to.

## 2 Exploratory Data Analysis

### 2.1 Feature Correlation and Dataset Balance

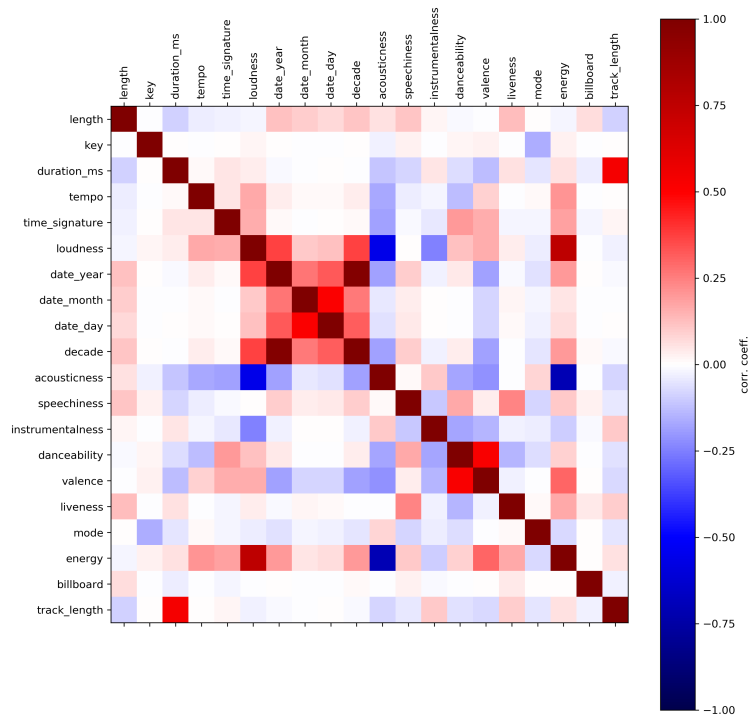


Figure 1: Feature correlation

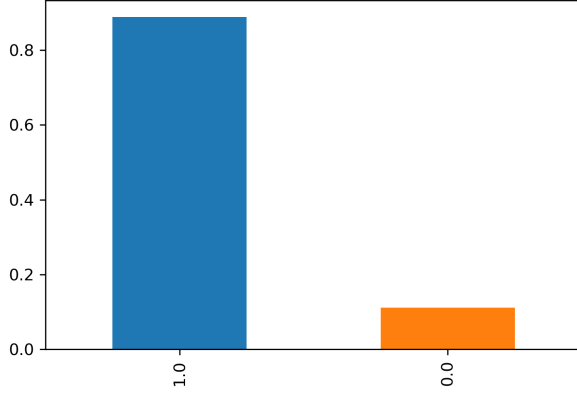


Figure 2: Dataset Class Balance

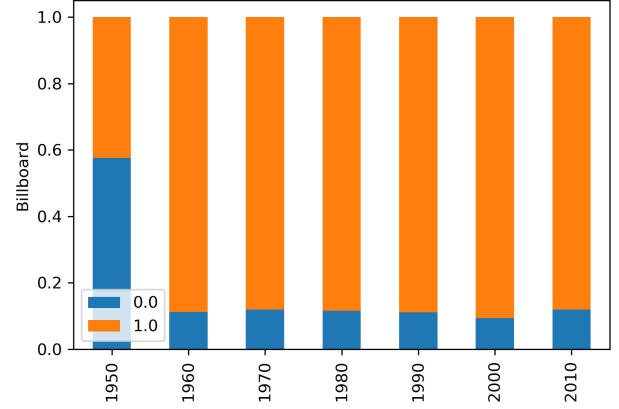


Figure 3: Track distribution over ages

The features are highly correlated from each other as observed in Fig. 1 and the dataset is imbalanced ( $\approx 88\%$ ) (refer to Fig. 2). However the data is distributed to match the overall dataset balance when observed on a yearly basis (Fig 3). To study the evolution and behaviour, we checked how the acoustic features evolve with time and with respect to songs on or not on the billboard.

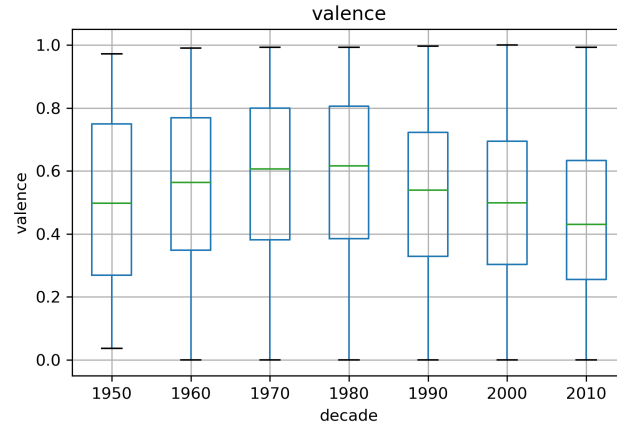


Figure 4: Valence evolution with time

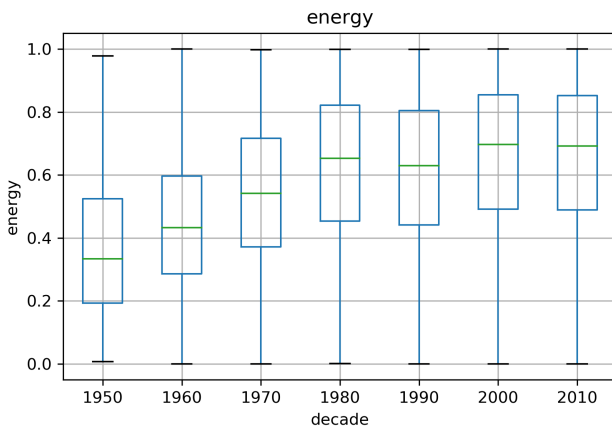


Figure 5: Energy evolution with time

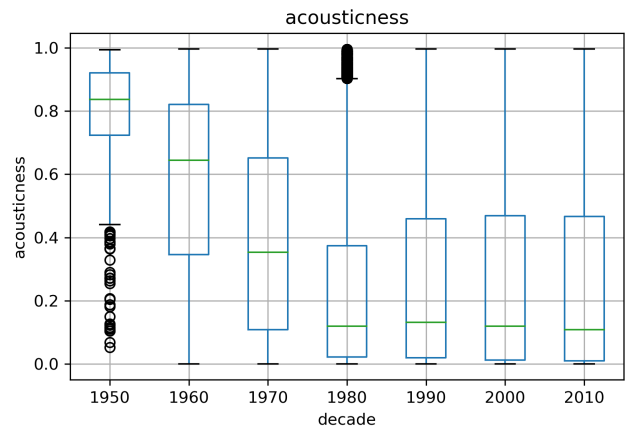


Figure 6: Acousticness evolution with time

Surprisingly, the range of values for most features do not vary with time, except acous-

ticness, valence and energy. With the advent of electronics, it makes sense that songs have departing from pure acoustic nature (Fig. 6) . Likewise, rock, metal, and more recently hip hop and electronic dance music have increased the energy levels of songs (Fig. 5) . Surprisingly, the valence (an indicator of the positiveness of a song) has slightly decreased over the years (Fig. 4).

Due the low variability of acoustic features and imbalance of the dataset, before we move to model evaluation, we expect the classifier to have a tough time (as can be guessed from the plot of the 1st two principal component (Fig. 7) where the points are highly overlapping)

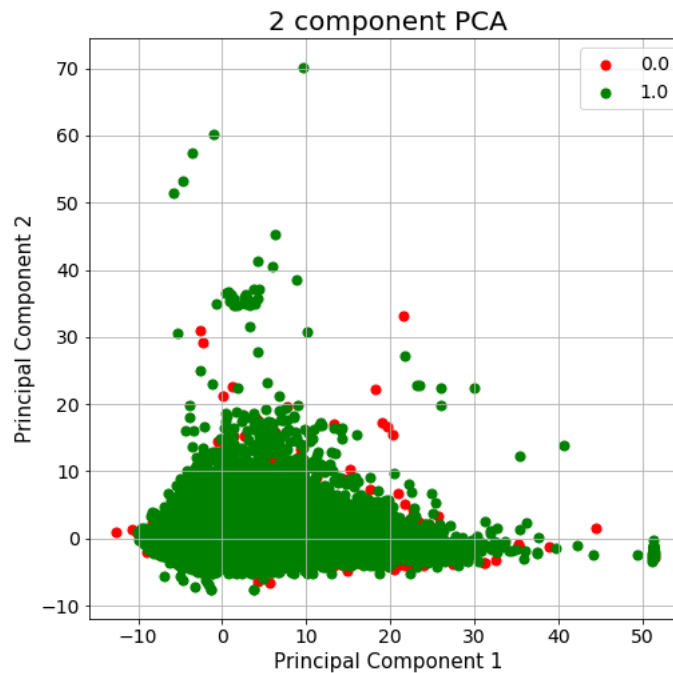


Figure 7: Plot of 1st two principal component

## 3 Machine Learning Pipeline

### 3.1 Preprocessing

The following are the broad categories of pre-processing applied to the datasets:

1. The date values (date released and dates on the chart) were each parsed to python datetime object and split to 3 columns - day, month, year. No scaling is done as yet on these features.
2. The song\_id and album\_id are kept in the spotify unique id format. To be scaled upon further discussion.
3. Acousticness, danceability, energy, instrumentalness, liveness, mode, speechiness, valence features are already scaled between  $[0, 1]$ . No further pre-processing was required.
4. Rank, track duration, key, tempo, time signature, loudness, album duration, album length, were all scaled using the StandardScaler.
5. The billboard label is added as 1 (if present) and 0 otherwise.

Note for the final calculation, I have ignored the month and day (since that is arbitrary and does not play any role when it comes to billboard consideration, and thus work only with the year. Since artist name, album name and song name was only used to curate songs, they are dropped from the final dataset before pre-processing.

### 3.2 Cross Validation Pipeline

In this study, we are trying to predict whether future songs would appear on the billboard, hence time of the song is not a factor in the requirements. We could have dropped the date feature but would anyway expect it have least importance.

However the data is highly imbalanced. Therefore the choice of cross-validation pipeline was chosen to be a stratified KFold Cross Validation, where the global distribution of classes is maintained in the K fold train, CV splits. The result of the splits are shown in Fig. 8.

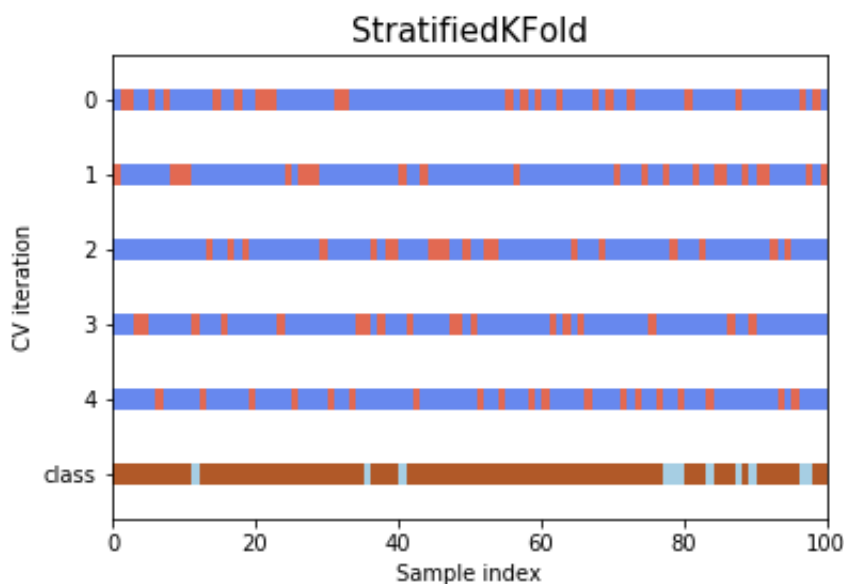


Figure 8: Cross Validation Folds. Training set is shown in blue, CV set in red.

### 3.3 Model Selection

Due to the large number of datapoints, certain popular techniques such as K-nearest Neighbours, and Support Vector Machine Classification failed to converge after 12 hours. Other methods such as Linear and Quadratic discriminant analysis failed to improve model performance over the baseline. For each of the model, we performed an exhaustive search for the parameter value in each estimator. Details of which is also described under each method.

**Evaluation Metric:** For all the models, we used the accuracy score since it is a binary classification problem.

In this report, we shall discuss 3 methods which showed improvement above the baseline. The 3 methods are:

1. Random Forest Classification
2. XGBoost Classification
3. AdaBoost Classification

### 3.3.1 Random Forest Classification

**Test Baseline:** 0.8807

**Test Score:** 0.8865

**Optimal Parameters:**

n\_estimators: 100, criterion: gini, max\_depth=10, min\_samples\_split: 3

The normalised confusion matrix is shown in Fig. 9

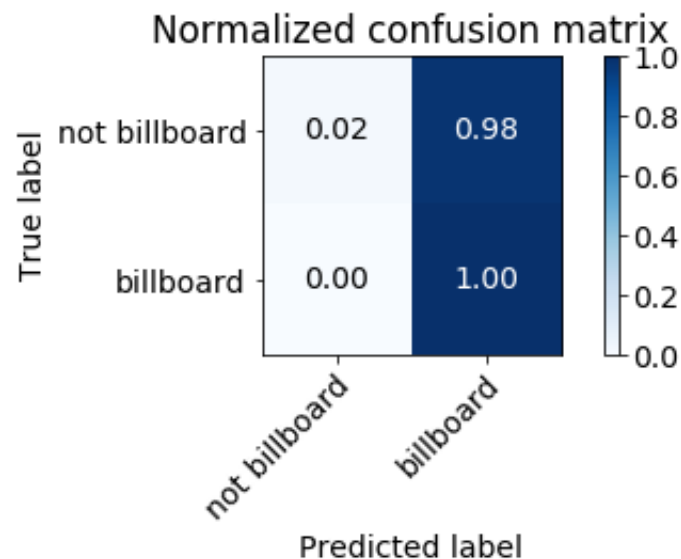


Figure 9: Normalised confusion matrix for Random Forest Classification

The feature importance under random forest is shown in Fig. 10

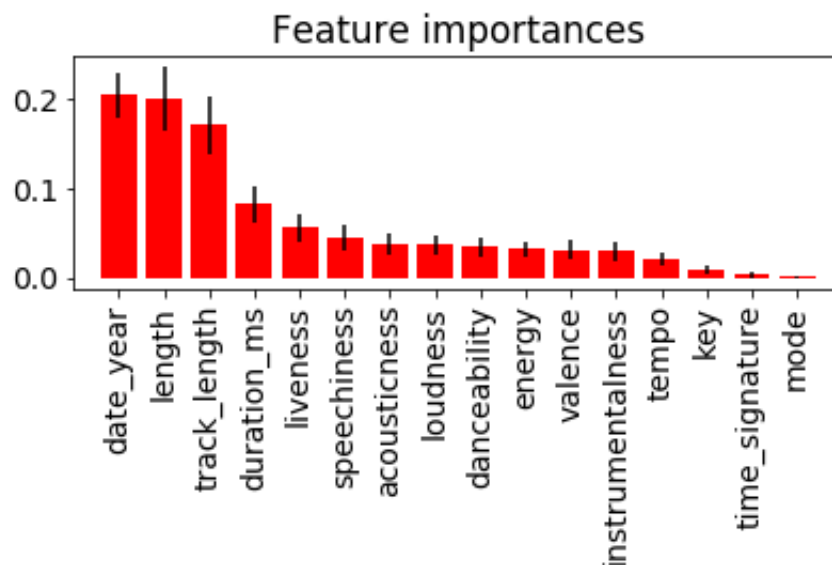


Figure 10: Feature importance for Random Forest Classification

### 3.3.2 XGBoost Classification

**Test Baseline:** 0.8807

**Test Score:** 0.8865

**Optimal Parameters:**

max\_depth=15, gamma: 0.4, min\_child\_weight: 1, learning\_rate: 0.3 ...

The normalised confusion matrix is shown in Fig. 11

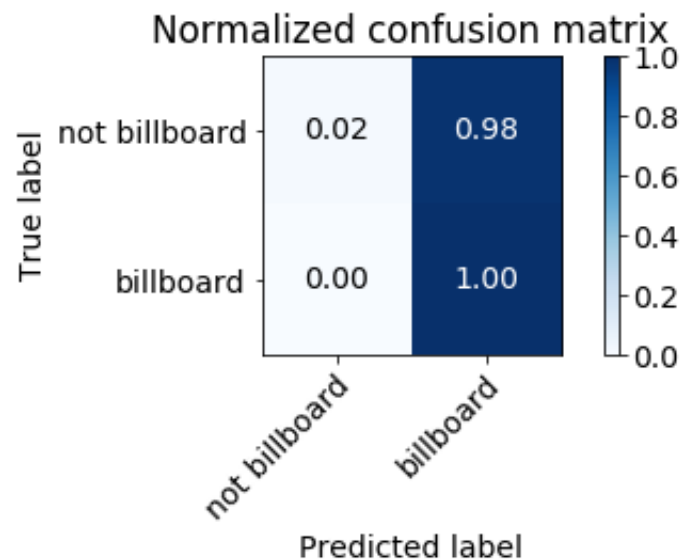


Figure 11: Normalised confusion matrix for XGBoost Classification

The feature importance by based on F score (as per SelectFromModel) under XGBoost is shown in Fig. 12

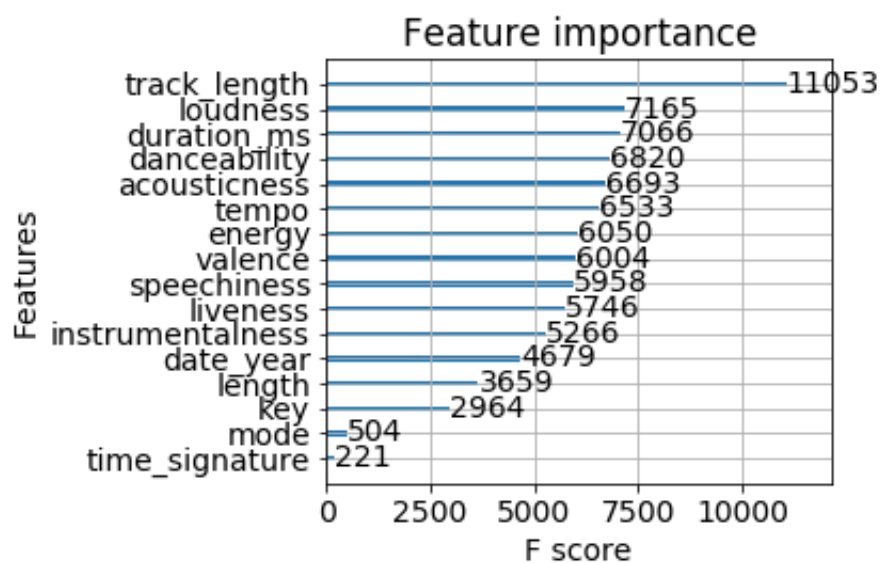


Figure 12: Feature importance by SelectFromModel for XGBoost Classification



### 3.3.3 AdaBoost Classification

**Test Baseline:** 0.8807

**Test Score:** 0.8865

**Optimal Parameters:**

algorithm: SAMME.R, gamma: 0.4, n\_estimator: 100, learning\_rate: 0.1

The normalised confusion matrix is shown in Fig. 13

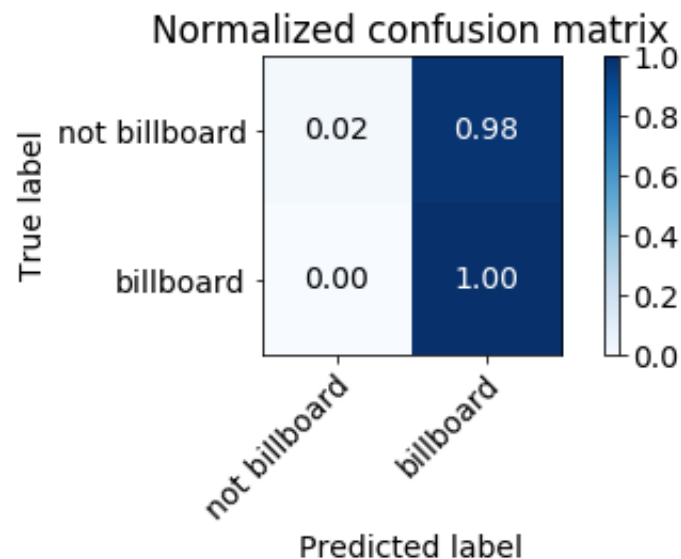


Figure 13: Normalised confusion matrix for AdaBoost Classification

The feature importance under AdaBoost is shown in Fig. 14

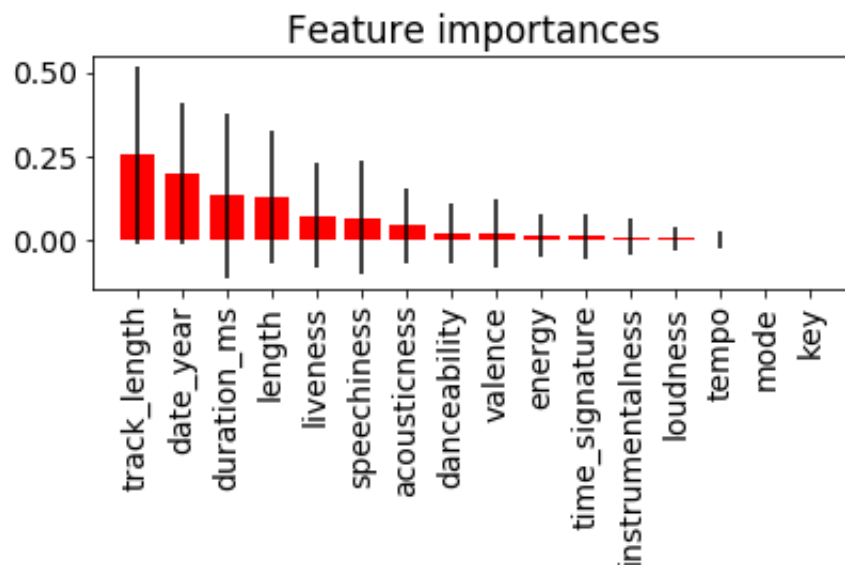


Figure 14: Feature importance by SelectFromModel for AdaBoost Classification

The tree estimator with least error of the adaboost classifier is show in Fig. 15

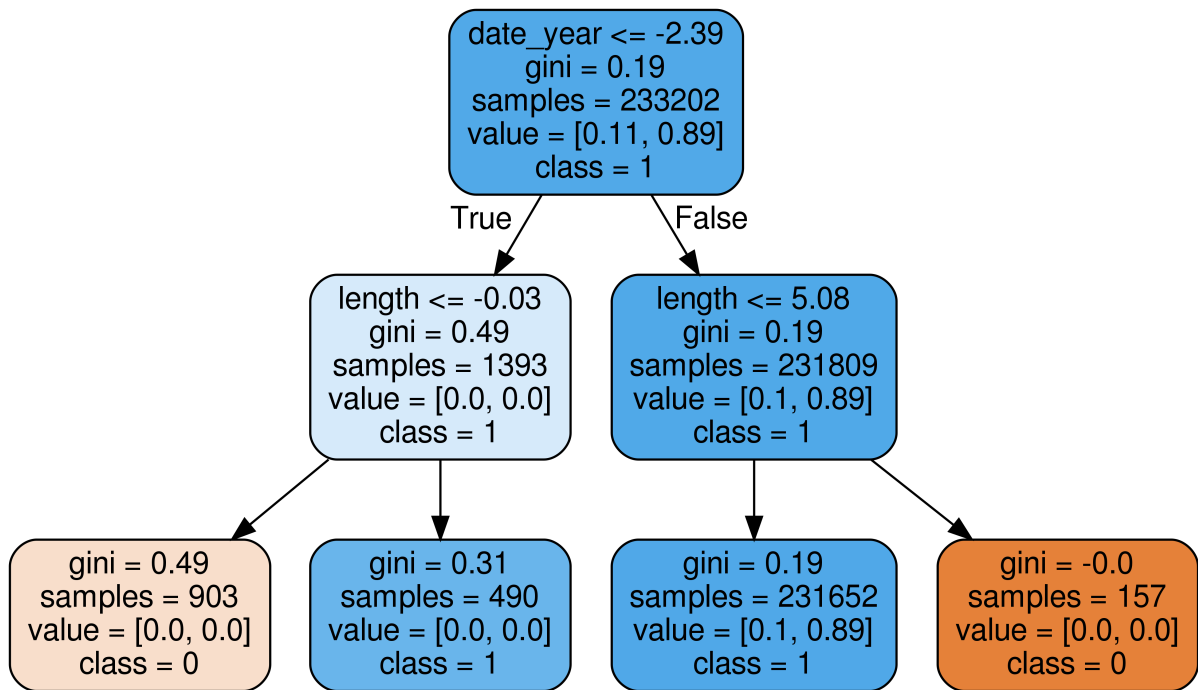


Figure 15: Estimator with least error

## 4 Outlook

The results of this project was not satisfactory. Except track length no other feature showed any importance that makes it stand out. The good news is that we can suggest solid steps to further improvement. The most important of them would be:

1. The dataset was highly imbalanced. Having a much more evenly divided set among the classes would improve the results.
2. The acoustic features of songs in both the classes overlap highly. To a large extent this was a fault of the method of data collection for songs that are not in the billboard. They were based on songs by artist that are similar to the artist in the billboard.
3. Advanced deep-learning methods could improve the result.
4. Due to lack of time, certain methods such as support vector machine classification or K-nearest neighbours couldn't be implemented. Among the models tested on, a wider grid search might have also led to some improvements.

However it was a great going through this project and the experience gained from this would surely translate to future endeavours.