

Shrey Anand, Ansul Atriek, Andrew Emerson, Amanul Haque,  
Suvodeep Majumder, Vidhisha Jaswani

## 1. Introduction

- What makes a song a hit?
- Predict hit or not using the **acoustic properties** of the song.
- Data features:** acousticness, duration, danceability, energy, key.

## 2. Data Description

- Data Collection:** 3 Phases with 10,000 instances and **Spotify API** to collect features.

- Data Pre-processing:** Normalization, Imputation, removing duplicate data and feature selection.

## 3. Technical Section

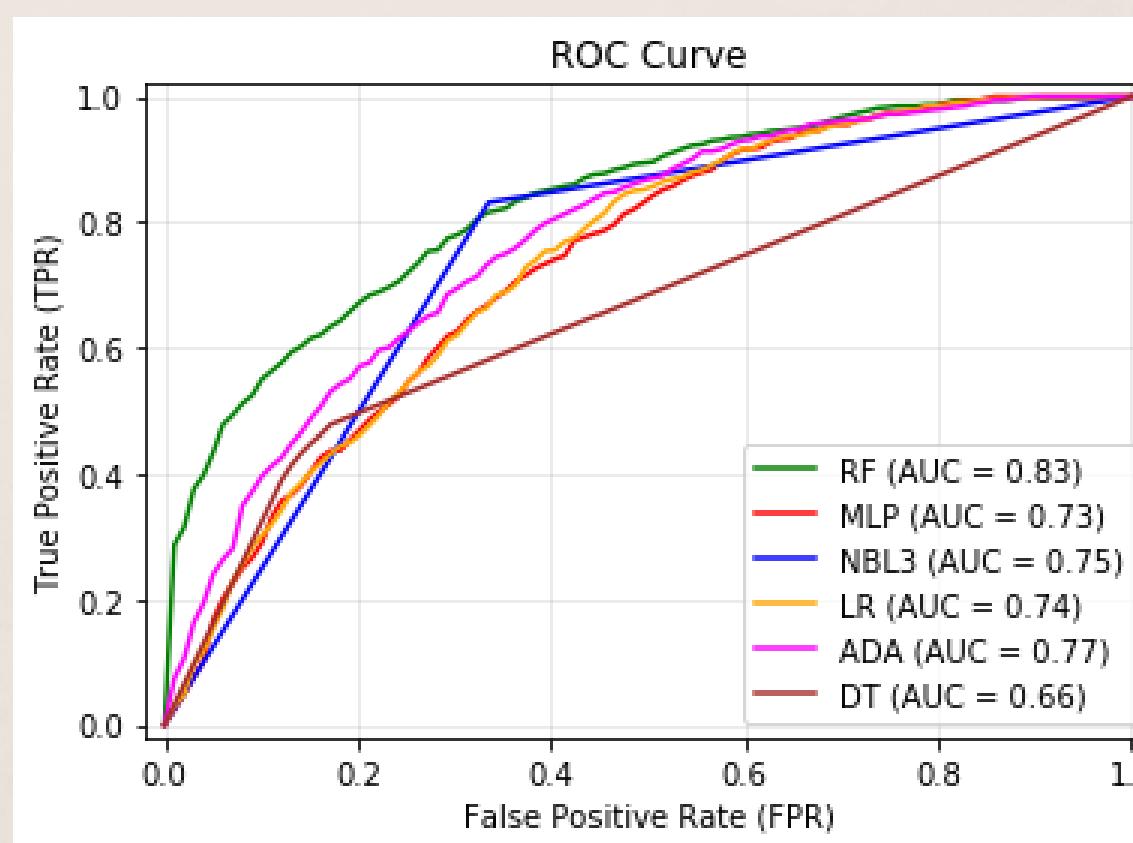
The algorithm from built-in libraries are:

- SVM
- Naive Bayes
- Logistic Regression
- Decision Trees
- Random Forests

Implemented the following algorithms from scratch:

- Neural Network
- Naive Bayes
- Multilayer tree based Naive Bayes using error correction.

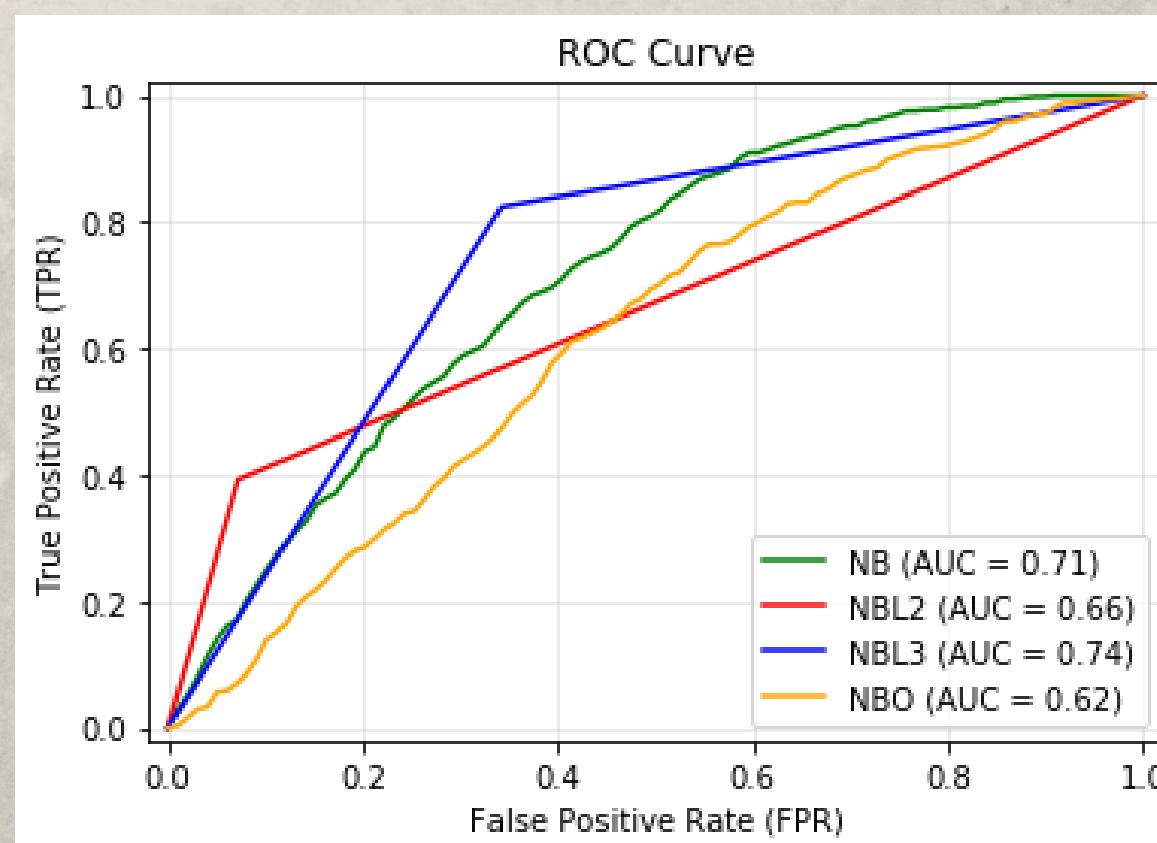
## 4. Results



Learner	Accuracy	F-measure	AUC
Naïve Bayes	0.697	0.689	0.691
Naïve Bayes <sup>**</sup>	0.575	0.621	0.617
2-Layer Naïve Bayes <sup>*</sup>	0.822	0.800	0.622
3-Layer Naïve Bayes <sup>*</sup>	0.693	0.727	0.744
Neural Network <sup>+</sup>	0.814	0.735	0.734
Neural Network <sup>**+</sup>	0.802	0.770	0.606
Decision Tree	0.767	0.770	0.628
SVM	0.816	0.733	0.552
Logistic Regression	0.816	0.733	0.707
AdaBoost	0.814	0.752	0.714
Random Forest <sup>+</sup>	<b>0.860</b>	<b>0.831</b>	<b>0.835</b>

<sup>\*</sup>Implementation from scratch

<sup>\*\*</sup>Implemented without Feature Selection.



## 5. Parameter Choices

### Naive Bayes

- priors: prior probabilities of the classes.

### Neural Network

- hidden\_layer\_sizes: ith element represents the number of neurons in ith hidden layer.
- activation: activation function for the hidden layer.
- learning\_rate: learning rate schedule for the weight updates.
- max\_iter: maximum number of iterations.

## 6. Conclusions

- The **feature selection** only marginally affected our results.
- Our implementations of Naïve Bayes and Neural Network **performed fairly** against those implemented in sklearn package.
- The **3-layer Naïve Bayes** method performs better than single-layer Naïve Bayes.
- SMOTE was also used to solve the class imbalance problem, but was found to not affect the overall results.
- Overall, **Random Forest** performs the best.