

---

# Machine Learning Algorithms for Music Genre Classification

---

**Konstantinos Kalais**  
University of Thessaly, Volos  
**kkalais@inf.uth.gr**

[https://github.com/Kkalais/Music\\_Genre\\_Classification](https://github.com/Kkalais/Music_Genre_Classification)

## Abstract

Music genre classification is a broadly considered subject as it is an effective technique to structure and arrange the huge quantities of music records accessible on the Internet. Machine Learning methods have demonstrated to be able to distinguish such enormous pools of information, and eventually classify them.

A musical genre is described by the usual characteristics shared by its users. These attributes are connected most of the times with the rhythmic structure of the related instruments, as well as the harmonic content of the music that stands out from other genres' harmonies. Using an automatic musical genre classification mechanism can help the user into creating music playlist easily and would be a useful tool to music information retrieval frameworks.

In this paper, we describe the process of feature selection among a large set of features of a music archived dataset. Then, we explain how each supervised (the algorithms learn on a labeled

dataset) machine learning method - we choose- classifies our samples, after analysing tracks metadata, and finally summarize a performance comparison from where we should get the most suitable technique for our problem. More specifically, classical machine learning algorithms such as Decision Trees, Support Vector Machines, K Nearest-Neighbors and MultiLayer Perceptron Neural Nets are employed. Using the proposed algorithms, classification accuracy of 74,21% for sixteen musical genres is achieved.

**Keywords:** machine learning, music genre classification, feature selection

## 1. Introduction

### 1.1. The problem

More than ever, the web has turned into a spot for sharing inventive work - for example music - among a worldwide network of music artists. While music and music accumulations originate before the web, the web empowered a lot bigger scale accumulations.

While individuals used to possess a bunch of vinyls or CDs, now they have moment access to the entire of distributed musical library by online sites. Such sensational increment in the size of music accumulations made the need to naturally

sort out a gathering of music items, as publishers and users can't manage them physically any longer.

## 1.2. The motivation

Music genre classification (5) is a significant activity with numerous real time applications. It is recommended that ten thousands tunes were discharged each month on web sites, for example, Soundcloud and Spotify. So as the amount of music being discharged everyday keeps increasing, the requirement for exact metadata required for music database organisation and management is expanded.

Having the option to quickly arrange tunes by genre in some random playlist or library is a significant task for any music streaming-purchasing administration, and the area for metadata analysis that complete and accurate music and audio classification gives is basically restricted.

Additionally, by classifying successfully music genres we can prompt to more customized music choices and cleared music creation services.

## 1.3. Related Works

Lots of researches had been done for music genre classification purposes, using Machine Learning techniques.

In 2018, Hareesh Bahuleyan (1) utilizes hand-extracted features of an Audio dataset's samples and trained four traditional machine learning classifiers (namely Logistic Regression, Random Forests, Gradient Boosting and Support Vector Machines) with 65% accuracy.

Danny Diekroeger (2) tried to classify songs based on the lyrics using a Naive Bayes classifier, but concluded that analyzing solely lyrics is not enough.

Nowadays, companies use music classification tools in order to recommend songs to their cus-

tomers (eg Soundcloud, Spotify) or detect a song (eg Shazam). Spotify had invented a lot of research in a service called "Spotify's Discover Weekly", where Machine Learning is applied in the core of this mechanism. Nature Language Processing as well as Deep Learning are used in order to classify million of songs of Spotify's database into a specific genre.

In our case, we will see how we can improve the existing results from previous studies by using classical Machine Learning algorithms on a feature selected dataset.

## 2. Methodology

### 2.1. Tools used

The code was arranged around Spyder. Spyder consists a scientific code environment designed for data analysts, engineers and scientists, and written in Python. It has various features for editing, debugging, analysis and tools for interactive execution and visualization of the results, as well as data and variable exploration. Apart from this, Spyder offers further capabilities through scientific and mathematical libraries like Pandas, Matplotlib, SymPy and more. In our study, we mostly make use of Scikit-learn, a free Machine Learning library for Python. It features various algorithms Support Vector Machine, Random Forest, and K-Nearest Neighbors, and it also supports Python numerical and scientific libraries-like NumPy and SciPy. We can use separate working areas for the Variables values, the command line and code editor.

### 2.2. Dataset Construction and Preprocessing

Experiments were carried out on a dataset for music analysis called FMA (3) (Free Music Archive), an open and easily accessible dataset suitable for evaluating several tasks in Music Information Retrieval, a field concerned with browsing, searching, and organizing large music collections. FMA provides 106.574 music pieces - with pre-computed features and track

level metadata - categorized in 16 musical genres, which we split them into training and testing set with 30% test size.

Table 1. Number of songs per genre

GENRE	SONGS
HIP-HOP	891
POP	346
ROCK	3889
EXPERIMENTAL	17
FOLK	874
JAZZ	238
ELECTRONIC	2166
SPOKEN	0
INTERNATIONAL	129
SOUL-RNB	0
BLUES	66
COUNTRY	0
CLASSICAL	265
OLD-TIME / HISTORIC	357
INSTRUMENTAL	81
EASY LISTENING	0
<b>TOTAL</b>	<b>9319</b>

The dataset had some minor inconsistencies such as missing features that we removed during the dataset cleaning phase. It includes 252 features and after dropping the highly correlated ones - with at least 80% correlation, in order to make our algorithm more cost and time efficient - it consists of 24 features where the basic ones are explained below.

As for the **features dropping step**, when we examine the features of a dataset, some of them might not be useful to make the necessary prediction. Correlation (equation 1) refers to how close two variables (features for our problem) are to having a linear relationship with each other. Features with high correlation are more linearly dependent (e.g.  $y = 3 \cdot x$ ) and hence have almost the same effect on the dependent variable, in contrast with two features that are non-linearly dependent (e.g.  $y = x^2$ ). So, when two features have high correlation, we can drop one of the two features.

We use the LabelEncoder class provided by Scikit Learn library, in order to transform the nominal values(numeric codes used for labelling

or identification, such as strings, datetimes, etc) to numeric and enable the algorithm to perform arithmetic operations on these values. Next, we create the correlation matrix, which is computed as follows:

$$\text{Correlation Matrix} = [\text{Cor}(\vec{x}_j, \vec{x}_k)]_{jk}$$

where  $\vec{x}_j$  and  $\vec{x}_k$  are the vectors of two features' values and:

$$\text{Cor}(\vec{x}_j, \vec{x}_k) = \frac{\sum_{i=1}^n (x_{ij} - \text{mean}(\vec{x}_j)) \cdot (x_{ik} - \text{mean}(\vec{x}_k))}{(\sum_{i=1}^n (x_{ij} - \text{mean}(\vec{x}_j))^2)^{-1/2} \cdot (\sum_{i=1}^n (x_{ik} - \text{mean}(\vec{x}_k))^2)^{-1/2}} \quad (1)$$

where  $\text{mean}(\vec{x}_j)$  and  $\text{mean}(\vec{x}_k)$  stand for the average of the values of vectors  $\vec{x}_j$  and  $\vec{x}_k$ , and visualize it through the correlation heat-map, as we can see on figure 1.

- Instrumentalness : how many vocals are present in a song
- Acousticness : the probability that the song has acoustic instruments
- Liveness : the probability the song was played in a live performance
- Speechiness : the presence of spoken words in a track
- Energy : the amount of activity and liveliness in the song
- Danceability : the probability that the songs is made for dancing-entertaining purposes and not a calm version of music
- Valence : the musical positiveness conveyed by a track. Tracks with high value on the valence feature sound more positive and happy, while tracks with low value sound more depressed and negative.

- Tempo : the general evaluated rhythm of a track in beats per minute (BPM)

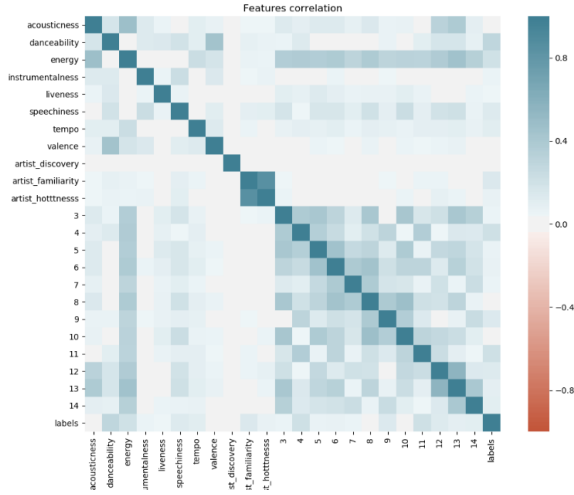


Figure 1. Plot of features correlation

Lastly, we scaled our data using the MinMaxScaler class of Scikit Learn in the range of (0,1) for easier convergence. MinMaxScaler is a good choice if we want our data to have a normal distribution or want outliers to have reduced influence. Specifically, if  $\max(X)$  and  $\min(X)$  are the maximum and minimum values that appear in all dataset for a given feature  $X$ , a value  $X$  is replaced by  $\text{new}X$  using the equation 2:

$$\text{new}X = \frac{X - \min(X)}{\max(X) - \min(X)} \quad (2)$$

This procedure was not implemented for the SVM methodology, as SVMs are scale invariant.

### 3. Classification Methods

#### 3.1. General Approach and Algorithms used

Based on the papers we mentioned, we decided on the use of eight algorithms to focus on the issue:

- Multilayer Perceptron

- Random Forest
- Support Vector Machine
- Gradient Boosting
- K-nearest neighbors
- Logistic Regression
- Linear Discriminant Analysis
- XGBoost

We will analyse below the work implemented in each one.

#### 3.2. Multilayer Perceptron

Neural networks (4) are built as the model of neurons present in the human brain. Neural networks (6) are typically stacked layers of interconnected neurons or nodes each having an activation function. The idea is to give our network several inputs (each input is a single feature). Then, we navigate through the network by applying weights to our inputs, summing them, and finally using an activation function to obtain the subsequent nodes. This process is repeated multiple times through the hidden layers, before applying a final activation function to obtain the output layer, our prediction. An example with  $n$  input nodes and 1 output node is shown on figure 2, with the step activation function given by equation 3.

$$f(x) = \begin{cases} 0 & x < 0 \\ 1 & x \geq 0 \end{cases} \quad (3)$$

The training takes place on the weights, which are initialized randomly. They are updated by trying to minimize training error using gradient descent (7), and are used after the training phase to make predictions.

**Gradient Descent.** Let's say we have an ex-

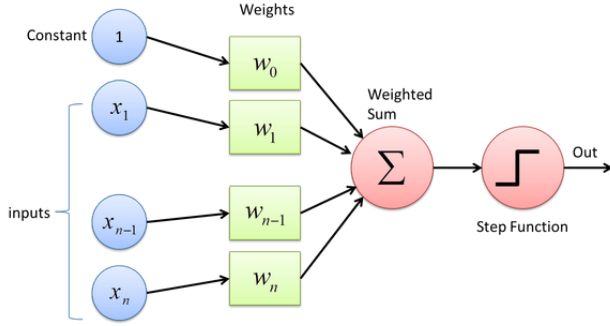


Figure 2. Multilayer Perceptron (4)

ample of a regression problem, so for given input  $x$  we want to predict the output. We present the hypothesis as (equation 4):

$$h_{\theta}(x) = \theta_0 + \theta_1 * x \quad (4)$$

where  $\theta_0$  and  $\theta_1$  are the parameters.

Then we need a function for minimization these parameters over our dataset. One common function that is often used is mean squared error, which measure the difference between the estimator  $Y_i$  (the dataset) and the estimated value  $\hat{Y}_i$  (the prediction). It looks like this (equation 5), supposing that  $n$  is the number of samples in our dataset:

$$MSE = \frac{1}{n} \cdot \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 \quad (5)$$

It turns out we can adjust the equation a little to make the calculation down the track a little more simple. We end up with (equation 6)

$$\frac{1}{2n} \cdot \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)})^2 \quad (6)$$

I will refer to the cost function as  $J(\theta)$ . In our problem we will need to minimize a cost function, called Mean Squared Error (MSE) (8). The use of Gradient Descent function is valuable for minimizing it, by changing the theta values, or parameters, step per step, until we hopefully achieved a minimum.

We start by initializing  $\theta_0$  and  $\theta_1$  to any two values, say 0 for both, and go from there. Formally, the algorithm is as follows:

$$\theta_j = \theta_j - \alpha \cdot \frac{\partial J(\theta_0, \theta_1)}{\partial \theta_j} \quad (for j = 0, 1) \quad (7)$$

and repeat until convergence:

$$\theta_0 = \theta_0 - \alpha \cdot \frac{1}{n} \cdot \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) \quad (8)$$

$$\theta_1 = \theta_1 - \alpha \cdot \frac{1}{n} \cdot \sum_{i=1}^n (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \quad (9)$$

where  $\alpha$  is the learning rate, or how quickly we want to move towards the minimum.

In our case, we construct a probability distribution of the 16 genres (figure 3) by running the outputs through a softmax activation (9) function.

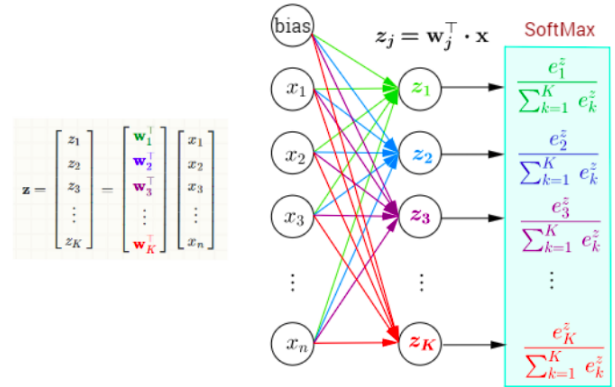


Figure 3. Softmax activation function (9)

Simply speaking, the sigmoid function only handle two classes, which is not what we expect. The softmax function squashes the outputs of each unit to be between 0 and 1, just like a sigmoid function (10). But it also divides each output such that the total sum of the outputs is equal to 1.

The output of the softmax function is equivalent to a categorical probability distribution (11),

it tells you the probability that any of the classes are true.

Mathematically the softmax function is shown below, where  $\vec{z}$  is a vector of the inputs to the output layer (if you have 10 output units, then there are 10 elements in  $\vec{z}$ ), and  $w_j$  are the weights. And again,  $j$  indexes the output units, so  $j = 1, 2, \dots, K$ .

Softmax activation function is defined as (equation 10) :

$$\sigma(j) = \frac{\exp(w_j^T \cdot x)}{\sum_{k=1}^K \exp(w_k^T \cdot x)} = \frac{\exp(z_j)}{\sum_{k=1}^K \exp(z_k)} \quad (10)$$

The model was compiled with Adam optimizer (12), since its is appropriate for problems with very noisy/or sparse gradients, and a categorical cross-entropy loss function (13), since we have multiple classes, and trained in batches of 150 instances for 50 iterations.

For the fully connected neural networks architecture we used:

- A Dense input layer of 150 neurons and rectifier linear unit activation function (equation 11)

$$ReLU(x) = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases} \quad (11)$$

- A Dense hidden layer of 100 neurons and rectifier linear unit activation function.
- A Dense output layer of 16 output classes and softmax activation function.

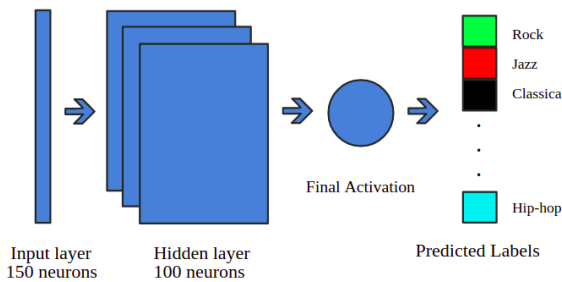


Figure 4. Neural Network (6)

To create our fully-connected neural network (figure 4), we begin with creating a Sequential model (14), a linear stack of 3 layers (input, hidden and output as mentioned above). Before training the model, we need to configure the learning process, which is done via the compile method. And then we train our model based on our training data, so as to learn our weights.

After trying various numbers of hidden units in the hidden layer, we conclude that the best accuracy is achieved with 140 neurons. However, we choose a hidden layer of 100 neurons since in that situation we have almost the same accuracy result and smaller execution time (figure 5).

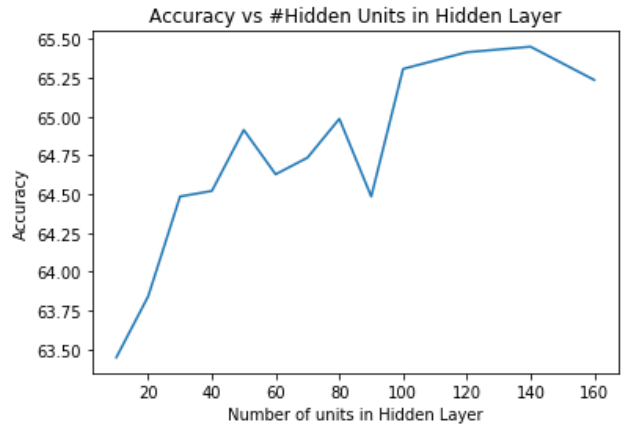


Figure 5. Accuracy vs Hidden Units in Hidden Layer

### 3.3. Support Vector Machine

In machine learning, Support Vector Machines (15) are supervised learning models used for data analysis and pattern recognition in classification and regression problems. SVMs (16) transforms the input feature space into higher-dimensional feature space using the kernel trick dot product, where  $(x, y)$  is replaced with  $k(x, y) = \langle f(x), f(y) \rangle$  ( $f$  is called kernel function). We can find each dataset's sample distance to a given dividing hyperplane (17). We call margin (18) the minimum distance from the samples to the hyperplane. The transformed data can be separated using a hyperplane, the dividing curve between distinct classes. The optimal hyperplane maximizes



the margin (figure 6).

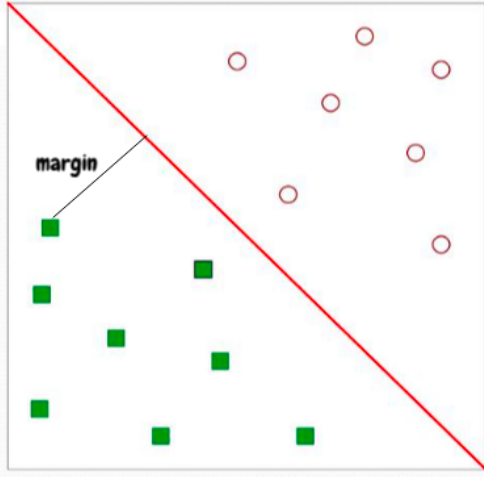


Figure 6. Hyperplane in 2d feature space (17)

In other words, an SVM is a linear separator that focuses on creating a hyperplane with the largest possible margin. Its goal is to classify a new sample by simply computing the distance from the hyperplane. On a two dimensional feature space, the hyperplane is a single line dividing the two classes (as we can see on figure 6). On a multi-dimensional feature space, where the data are non linearly separable an SVM cannot linearly classify the data. In this case it uses the kernel trick. The main concept has to do with the fact that the new multidimensional feature space could have a linear decision boundary which might not be linear in the original feature space. This can be done using various kernel functions such as :

- Polynomial kernel (19), popular in image processing. Equation is :

$$k(x, y) = (x \cdot y + 1)^d \quad (12)$$

where d is the degree of the polynomial

- Gaussian kernel (20), used when there is no prior knowledge about the data. Equation is :

$$k(x, y) = \exp\left(-\frac{\|x - y\|^2}{2\sigma^2}\right) \quad (13)$$

, where  $\sigma$  is standard deviation.

- Gaussian radial basis function (RBF) (21), also used when there is no prior knowledge about the data. Equation is :

$$k(x, y) = \exp(-\gamma\|x - y\|^2) \quad (14)$$

where  $\gamma = 1/2\sigma^2$

- Hyperbolic tangent kernel (22), which can be used in neural networks. Equation is :

$$k(x, y) = \tanh(kx \cdot y + c) \quad (15)$$

In our problem, a radial basis function (RBF) kernel is used to train the SVM because such a kernel would be required to address this non-linear problem.

Large value of parameter C should cause a small margin, and the opposite. Small value of parameter C should cause a large margin. There is no rule to choose a C value, it totally depends on our testing data. The only way is to try various different values of the parameters and go with the one that gives the highest classification accuracy on the testing phase. As for the gamma parameter, if its value is low then even the far away points can be taken into consideration when drawing the decision boundary and the opposite. With a high gamma parameter, the hyperplane is dependent on the very close points and ignores the ones that are far away from it. So, the best parameters we used in the Support Vector Machine technique, after testing several variants with trial-and-error approach, were 10 for the C value, and 1 for the gamma value.

### 3.4. Random Forest

A decision tree (24) is a graphical representation of data which applies a branching method to describe with illustrations any possible result of a decision. Each tree branch stands for a possible option that is available for taking a decision. The depth of the tree is extended as new decisions need to be made. So the decision tree (25)

as a tree structure, has few components; internal nodes for a test on an attribute, the branches that represents a test result and leaf nodes (terminal nodes) that contains information for the class labels. The outcome of a test -based on the node's attributes- creates a new branch, and as we go from root to the next node and so on we reach a leaf node, that belongs to a unique label (class).

As we can see in figure 7, we have to decide for example if a person is fit. Following a series of boolean tests we will end up on a "fit" or "Unfit!" leaf node, concluded for the fitness's level of a person.

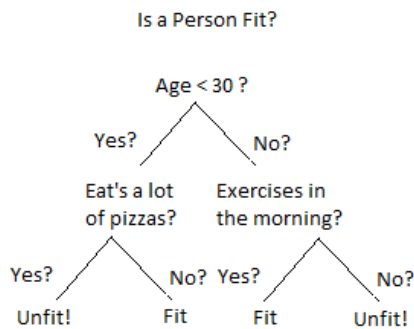


Figure 7. Decision tree toy example (24)

Random Forest (23) is a summation of Decision Trees. The general idea of this technique is that a mixture of learning models raises the general result. Random forest builds multiple decision trees and merges them together to achieve the preciseness and stability of the prediction. In that way, it prevents overfitting by creating random subsets of the features, building smaller trees using these subsets and combines them so as to increase the overall performance. With overfitting we mean the result of an overly complex model with too many parameters. If a model is overfitted, we can result that it is inaccurate because its behaviour does not represents real data information.

As we can see in the figure 8, a simplified Random Forest could categorize a sample to the class with the maximum "votes" among each subtree.

Random Forest makes the model more random, while developing the trees. Rather than looking for the most significant feature while splitting a node, it scans for the best element among a random subset of features. This outcomes in a wide variety that by and large results in a greater model.

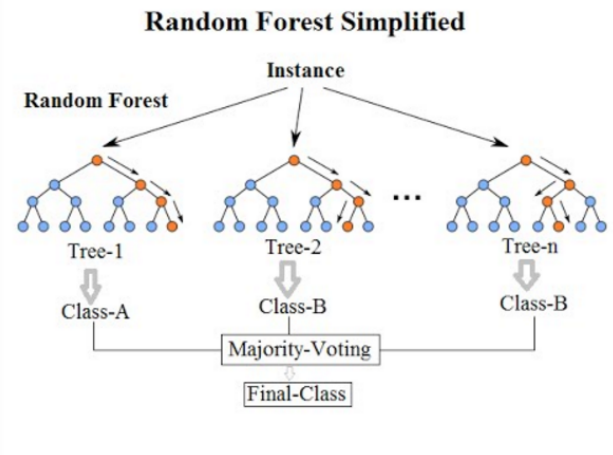


Figure 8. Random Forest Simplified (23)

The parameters we used after trying several variants, were 29 estimators - the number of trees in the forest - and 29 as the depth of the tree. The algorithm provided us with a relatively high accuracy result in a short time, specifically 68.74% in 1.29 seconds (figure 9). Generally our opinion was that the algorithm provided us with the best results in regard to how simple it was to develop and fast to execute.

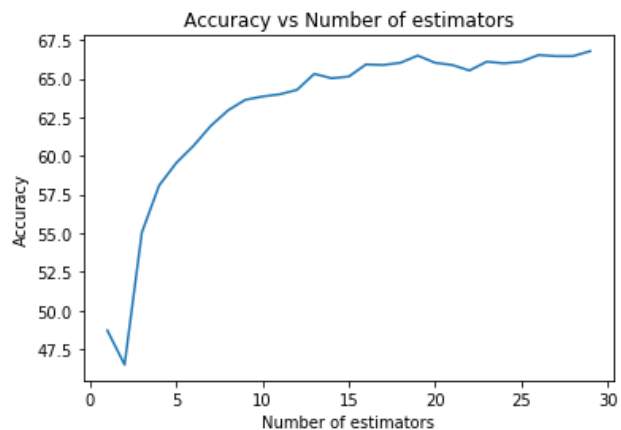


Figure 9. Accuracy vs Number of estimators



### 3.5. Gradient Boosting

Gradient boosting (26) as a machine learning method for classification and regression problems is trying to create a model for predictions in the form of a mixture of inefficient prediction models, that are called decision trees. In the boosting phase, every new tree is a fit on an altered version of the original data set. Firstly, the Gradient boosting trains a decision tree and assigns each observation an equal weight. After the first tree assessment, we lower the weights for the observations that are easy to classify and increase the weights of those that are hard to classify. Then, we grow the next tree on this weighted data where we try to improve the predictions of the first tree. Our new model is tree 1 together with tree 2. Then we compute the classification error from this combination of two-tree model and create a third one to predict the revised residuals. We repeat this process for a specified number of iterations. Subsequent trees help us to classify the non-well categorised observations by the former trees. The predictions of the final ensemble model is the weighted sum of the predictions made by the previous tree models.

After trying several variants, we set the values 0.1 and 150(with the smaller execution time) to the learning rate and the estimators - the number of boosting stages to perform - respectively and train our model with 74.21% accuracy (figure 10), the highest between the eight algorithms. Though, Gradient Boosting takes 39.77 seconds to execute, much slower than the above-mentioned techniques.

### 3.6. K-Nearest Neighbors

The k-nearest neighbors (KNN) (27) is a supervised machine learning algorithm useful for solving regression and classification problems in a simple and easy way. The KNN algorithm is based on the assumption that same things exist in a close area. In other words, similar things are close to one another.

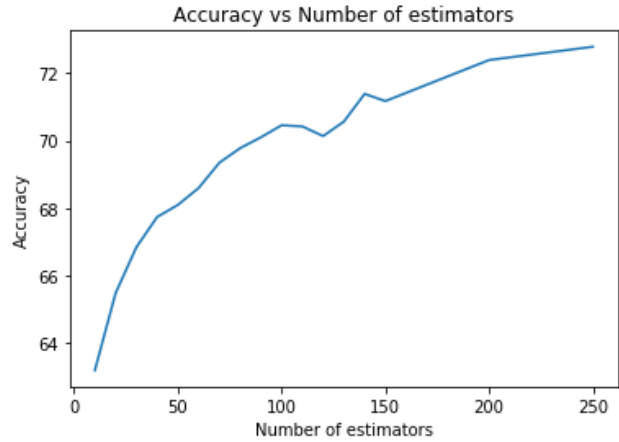


Figure 10. Accuracy vs Number of estimators

KNN (28) is based on the idea of similarity (also known as distance, proximity, or closeness) figuring the space between points on a graph. There are various methods of calculating distance, and one way might be preferable depending on the problem we are solving. For our problem, we use the most common distance metric, which is the Euclidean distance between two points  $x_1$  and  $x_2$ , and is given by the formula:

$$d(x_1, x_2) = ||x_1 - x_2||_2 \quad (16)$$

We will see now who does the KNN algorithm (29) work. First of all, we have to load the data - usually called X- and their target values - usually called y- we want to classify. Then we initialize K to a preferable number of neighbors and for each data sample we compute the distance between the sample whose target value we want to classify and the current sample from the data. We add both the index and the distance of the query example to an ordered list of indices and distances and sort this list in an ascending order(from smaller to bigger), with the distance as order criteria. Finally, we pick the first K entries from the sorted list and get the labels of the selected K entries, so we can return the form of the K labels.

For example, as we can see in figure 11 for K=3 the tested sample would have been classified

to class B, and for  $K=7$  would have been classified to class A.

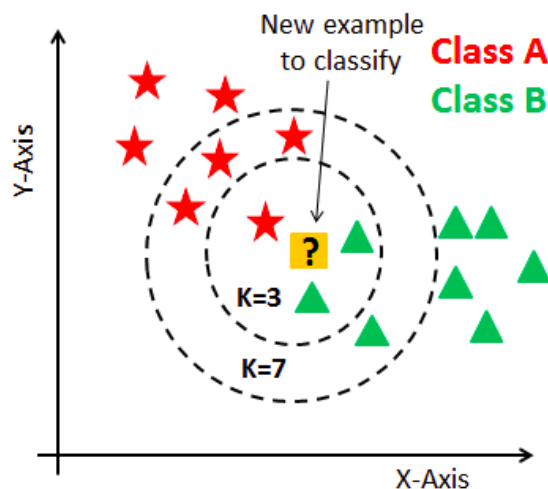


Figure 11. KNN example (28)

In order to choose the best number for the variable  $K$  to suit with our problem, we have to execute the  $K$ -Nearest Neighbors algorithm some times with different values of  $K$  and choose the  $K$  that decreases the number of errors we meet while keeping up the ability of the KNN algorithm to make accurate predictions when it is provided with data that has never seen before. As the value of  $K$  is reduced to one, our predictions' stability is also decreased. On the other, as we raise  $K$ 's value, our predictions' stability is increased to because of the majority voting in the compared samples. As a result, it is more possible to increase - until a specific level- the percentage of accuracy for our predictions.

In our problem, we choose 5 as the number of neighbors, despite the fact that we achieve the highest accuracy with 17 (figure 12), since we have better execution time.

### 3.7. Logistic Regression

Logistic regression (30) is a supervised machine learning algorithm used for classification problems, and specifically for categorizing observations into a group of discrete classes. Although

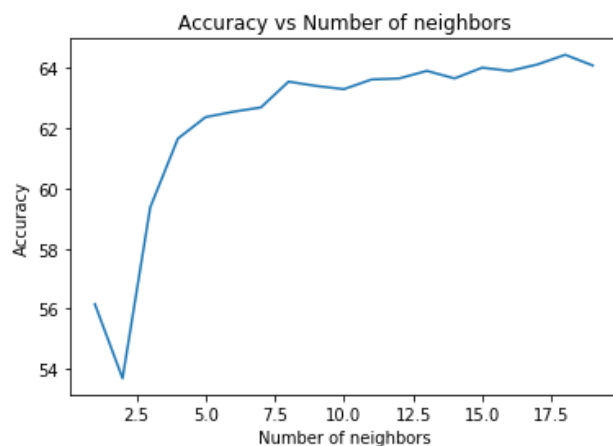


Figure 12. Accuracy vs Number of neighbors

linear regression assign observations to a continuous number values, logistic regression applies on its output a transformation - activation - function, called the logistic sigmoid function. In that way it returns a probability value which can then be matched with two or more classes. Logistic Regression (31) is used when the target - dependent - variable is categorical. For example, to predict whether an email is spam (1) or not (0) (binary logistic regression) or to predict whether a car with specific characteristics belongs to a model like BMW, Mercedes, Ford, etc (multi-class logistic regression).

**Binary logistic regression.** Suppose we are given a dataset (figure 13) that includes student exam results and we want to predict if a student will pass or fail the exam according to the combination of number of hours he has slept and hours spent for studying. In our data, we have to deal with the features of slept hours and studied hours and only two classes: passed (1) and failed (0), as a result of the pass or fail of a student's result.

Graphically we could represent these data with a scatter plot (figure 14).

**Sigmoid activation.** In machine learning, we make use of the sigmoid function so as to match predicted values with the probabilities. The function convert any real value into a probability, which is in fact another value between 0 and 1

Studied	Slept	Passed
4.85	9.63	1
8.62	3.23	0
5.43	8.23	1
9.21	6.34	0

Figure 13. The first four samples of these data

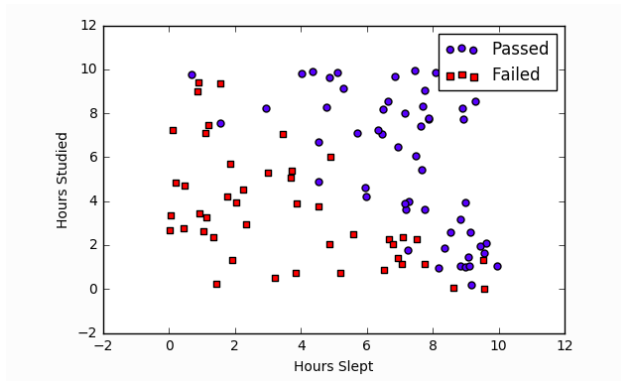


Figure 14. Scatter plot of our data

(figure 15).

- $s(z)$  = output between 0 and 1 (probability estimate)
- $z$  = input to the function (your algorithms prediction e.g.  $ax + b$ )
- $e$  = base of natural log

So, sigmoid function is given by the formula :

$$S(z) = \frac{1}{1 + e^{-z}} \quad (17)$$

**Decision Boundary.** The sigmoid function returns the predicted value of the probability in the range between 0 and 1. In order to map this to a value that belongs to a discrete class (true/false, BMW/Ford/etc), we select a threshold value below which we will classify values into the second

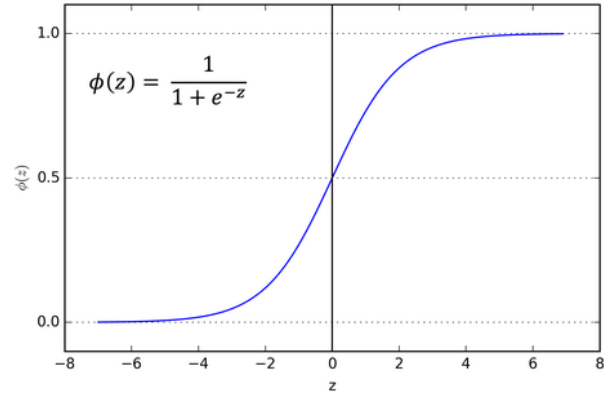


Figure 15. Sigmoid function (10)

class (in the case of binary logistic regression) and above which we will classify values into the first class. The same concept is followed when we have more than two discrete classes (multiclass logistic regression).

For example, if our tipping point (threshold value) was 0.5 and our prediction function returned 0.8, we would classify this observation to the class 1 (positive value). If our prediction was 0.3 we would classify the observation to the class 2 (negative value). In the case of multiclass logistic regression where we have to decide between multiple classes to map the observation, we could map the observation to the class with the highest predicted probability (figure 16).

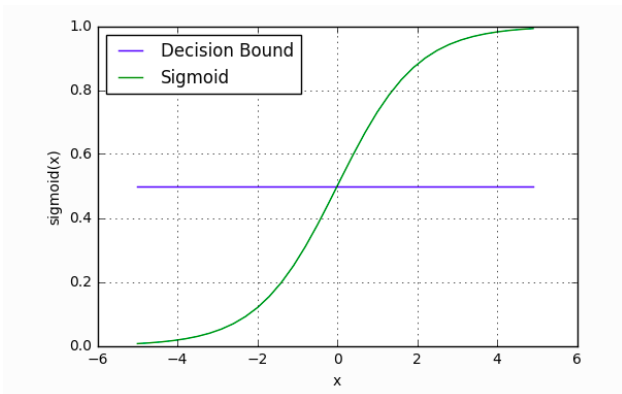


Figure 16. Decision boundary (37)

We can create now a prediction function by taking into consideration our knowledge of de-

cision boundaries and sigmoid functions. Such a function could compute the probability of our observation for being positive. In other words, the probability of our observation belongs to the class 1 and its notation is  $P(\text{class}=1)$ . As the probability increases and reaches 1, our model becomes more stable and we could be more confident that the observation is in class 1. On the other side, as the probability decreases and gets far away from 1, our model gives a higher probability that the observation is in class 2 (equation 18).

$$P(\text{class} = 1) = \frac{1}{1 + e^{-z}} \quad (18)$$

So, now using the sigmoid function we will apply the transformation to our output and it will return a probability value between 0 and 1.

$$z = w_0 + w_1 \text{ Studied} + w_2 \text{ Slept}$$

The final step is to map the predicted probabilities to the discrete classes with labels 0 or 1. So, we classify our sample at class with label = 1 if the probability  $\geq 0.5$ , else at class with label = 0.

In our FMA dataset, we use multiclass logistic regression since we have 16 discrete classes for our predicted labels. So we replace the definition of  $y = 0, 1$  and use a new definition of multiple labels,  $y = 0, 1, \dots, 15$ . In fact we run binary classification multiple times, once for each class. First, we divide the problem into  $n+1$  ( $n=16$ ) binary classification problems and for each class we predict the probability the observations are in that single class. The predicted class comes from the maximum probability between the classes.

In our problem, we use Sklearn's LogisticRegression classifier class, with lbfgs solver since we have a multiclass problem and can also handle multinomial loss. For the same reason multiclass parameter is equal to multinomial. In order for the method to converge we choose 300 as max iterations number.

### 3.8. Linear Discriminant Analysis

For binary classification problems, we are familiar with using Logistic Regression algorithm. However, in a multiclass classification problem with more than two discrete classes, Logistic Regression is not the best algorithm due to its complexity. Instead, we prefer to use another dimensionality reduction technique: Linear Discriminant Analysis (LDA) (32), useful for supervised machine learning problems.

Before we dig into the characteristics of Linear Discriminant Analysis, we will take a look of the Naive Bayes (33) classification algorithm that creates the basic details for LDA.

**Naive Bayes.** Given a feature vector  $X = (x_1, x_2, \dots, x_k)$  and a class variable  $C_n$ , Naive Bayes Theorem mentions that:

$$P(C_n|X) = \frac{P(X|C_n)P(C_n)}{P(X)}, \text{ for } n = 1, 2, \dots, k \quad (19)$$

$P(C_n|X)$  stands for the posterior probability,  $P(X|C_n)$  for the likelihood,  $P(C_n)$  for the prior probability of class, and  $P(X)$  for the prior probability of the predictor. We want to compute the posterior probability from the likelihood and prior probabilities.

Assuming the Naive independence, which claims that  $P(X|C_n) = P(x_1, \dots, x_k|C_n) = \prod_{i=1}^k P(x_i|C_n)$ , the posterior probability can then be thought as:

$$P(C_n|X) = \frac{P(C_n) \prod_{i=1}^k P(x_i|C_n)}{P(X)} \quad (20)$$

Then, the Naive Bayes classification problem turns into the statement: for different class values of  $C_n$ , find the maximum value of  $P(C_n) \prod_{i=1}^k P(x_i|C_n)$ . This can be written as:

$$C = \arg \max_{C_n} P(C_n) \prod_{i=1}^k P(x_i|C_n) \quad (21)$$

### Discriminant analysis.

Now we can examine a wider concept, Discriminant analysis. Discriminant analysis can be analyzed by 5 steps:

**1:** Compute prior probabilities. The prior probability of class  $P(C_n)$  is equals to the relative frequency of class  $C_n$  in the training set of our data.

**2:** Check the homogeneity of the samples' variances, actually check if  $K$  samples are from populations with equal covariance matrices. After this test we could decide which algorithm to use, Linear or Quadratic Discriminant Analysis.

- If there is variances' homogeneity, we should use Linear discriminant analysis:  $\sum_1 = \sum_2 = \dots = \sum_K = \sum$ .
- If there is variances' heterogeneous, we should use Quadratic discriminant analysis:  $\sum_i \neq \sum_j$  for some  $i \neq j$ .

**3:** Likelihoods' parameters estimation. We estimate the parameters (e.g.  $\mu_i$  and  $\sum$ ) of the conditional probability density functions  $P(X|C_n)$  from the training set of our data. In fact, we have to take the assumption that the data are multivariate normally distributed.

**4:** Calculate discriminant functions. This computation will determine into which of the known populations to classify the new object.

**5:** Compute the performance of our classification: Make use of the cross-validation method in order to calculate misclassification probabilities. And we can now predict the new observations.

**Linear Discriminant Analysis.** We take the standard assumption that in population  $i$  the probability density function of  $x$  is multivariate normal with mean vector  $\mu_i$  and variance-covariance matrix  $\sum_i$  (same for all populations). The equation to compute this normal probability density func-

tion is:

$$P(X|p_i) = \frac{1}{(2\pi)^{p/2} |\sum|^{1/2}} \cdot \exp[-0.5(X - \mu_i)' \sum^{-1} (X - \mu_i)] \quad (22)$$

Now, with our knowledge about Naive Bayes classification algorithm states we are concerned about the fact that we make the classification for the example to the population for which  $P(\pi_i)P(X|\pi_i)$  is maximized. To simplify our calculations, we apply to the previous equation a logarithmic transformation.

In Linear Discriminant Analysis, we take the decisions based on the **Linear Score Function**, which is defined as:

$$\begin{aligned} s_i^L(X) &= -0.5\mu_i' \sum^{-1} \mu_i + \mu_i' \sum^{-1} X \\ + \log P(\pi_i) &= d_{i0} + \sum_{j=1}^p d_{ij} x_j + \log P(\pi_i) = \\ &= d_i^L(X) + \log P(\pi_i) \end{aligned} \quad (23)$$

where  $d_{i0} = -0.5\mu_i' \sum^{-1} \mu_i$  and  $d_{ij}$  =  $j$ th element of  $\mu_i' \sum^{-1}$ . And we call  $d_i^L(X)$  the linear discriminant function.

From the above equation it is obvious that the right-hand expression seems like a linear regression with coefficients  $d_{ij}$  and intercept  $d_{i0}$ .

Given an example with a feature vector  $X = (x_1, x_2, \dots, x_p)$ , we compute for each population the linear score function, and we map the given example into the population with the largest value of this score. In fact, we classify the example to the population that has the largest posterior probability of membership.

In our problem we use Sklearns LinearDiscriminantAnalysis classifier class to label our samples with a music genre.

### 3.9. XGBoost

XGBoost (34) belongs to the family of the ensemble learning methods. Some times, it could



be insufficient to depend on the results of only one machine learning method applied to our data. Ensemble learning techniques use a systematic method to combine the predictive power of various learning methods. The output of this combination is a model that provides the totaled result from smaller-weaker- models (figure 17). Most of the times, we use XGBoost algorithm with decision trees.

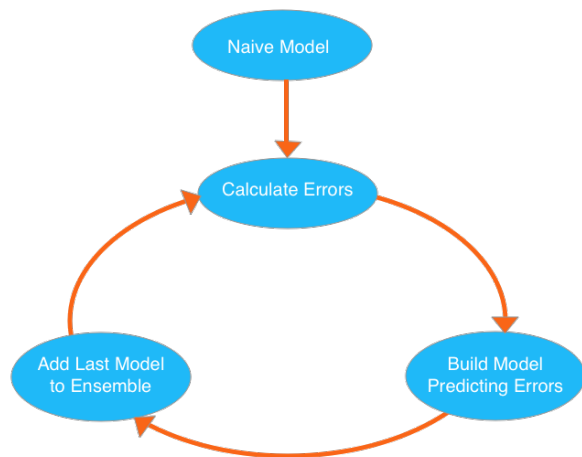


Figure 17. Gradient Boosted Trees diagram (35)

XGBoost (35) is the main model for dealing with the type of data you store in Pandas DataFrames, known as tabular data, and not other types of data like images and videos. XGBoost (36) belongs to the Gradient Boosted Decision Trees algorithms, which we will examine later on.

The basic concept follow the circular rule where the method repeatedly creates new models and blends them into an ensemble model. At the beginning of this procedure, we calculate the errors for each observation in the data. Then we create a new model to predict those errors. We add the error-predicting results of this model to the total of the model. Then we add the predictions from all previous models in order to make a new prediction. We can use the added predictions to calculate new errors, create the next model, and add it to the ensemble.

In the boosting phase, each tree is built in a subsequent way, trying to reduce the residual errors of the previous trees. The learning procedure follows the previous tree and updates the remaining errors. The next growing tree will learn from the last version of the residual errors.

The boosting method in its basis has learners with high bias and low - bit higher than random predictions - predictive strength. Each of these base learners produce crucial information that contributes to the ensemble model for the predictive power. As a result, by combining these base-weak learners the overall boosting method develops strong predictive power and becomes a strong learner (with lower bias and variance).

In the boosting technique, trees with fewer splits are usually used, and as a result such small trees can be easily managed and interpreted. On the other side, the bagging techniques such as Random Forest examine trees that are extended to the maximum of their depth. So it is preferable to go with a boosting instead of a bagging algorithm, since it is an updated version of the last one. Also, we have to select optimally the stopping criteria since a huge number of trees can cause overfitting. Applying validation techniques like k-fold cross validation we can carefully choose parameters such as the depth of the trees, the number of splits per level, etc.

The boosting phase consists of three steps:

- An initial model  $F_0$  is defined to predict the target variable  $y$ . This model will be associated with a residual  $(y - F_0)$
- A new model  $h_1$  is fit to the residuals from the previous step
- Now,  $F_0$  and  $h_1$  are combined to give  $F_1$ , the boosted version of  $F_0$ . The mean squared error from  $F_1$  will be lower than that from  $F_0$ :

$$F_1(x) < -F_0(x) + h_1(x) \quad (24)$$



To improve the performance of  $F_1$ , we could model after the residuals of  $F_1$  and create a new model  $F_2$ :

$$F_2(x) < -F_1(x) + h_2(x) \quad (25)$$

This can be done for  $m$  iterations, until residuals have been minimized as much as possible:

$$F_m(x) < -F_{m-1}(x) + h_m(x) \quad (26)$$

Here, the additive learners do not disturb the functions created in the previous steps. Instead, they impart information of their own to bring down the errors.

As the first step, the model should be initialized with a function  $F_0(x)$ .  $F_0(x)$  should be a function which minimizes the loss function or MSE (mean squared error), in this case:

$$\begin{aligned} F_0(x) &= \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma) \\ &= \arg \min_{\gamma} \sum_{i=1}^n (y_i - \gamma)^2 \end{aligned} \quad (27)$$

Taking the first differential of the above equation with respect to  $\gamma$ , it is seen that the function minimizes at the mean. So, the boosting model could be initiated with:

$$F_0(x) = \frac{\sum_{i=1}^n y_i}{n} \quad (28)$$

In our problem, we choose the default number of 100 estimators as it provides the highest accuracy in the shortest execution time (figure 18).

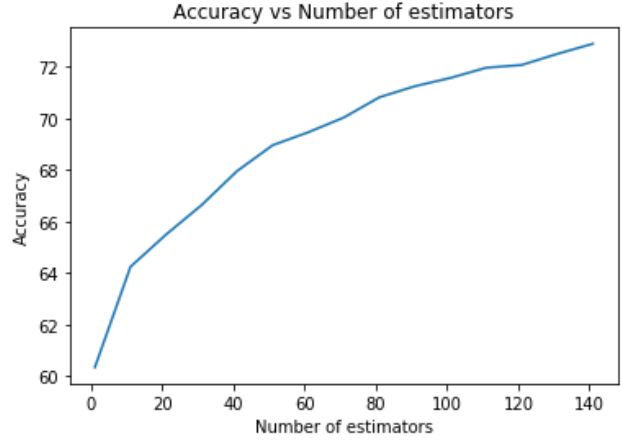


Figure 18. Accuracy vs Number of estimators

#### 4. Setup Installation

In order to test the algorithms over our dataset we need to install some tools, starting by updating the packages list and installing the prerequisites:

- `sudo apt update`
- `sudo apt install software-properties-common`
- `sudo apt-get -y install python3-pip`

Next, add the deadsnakes PPA to your sources list:

- `sudo add-apt-repository ppa:deadsnakes/ppa`

When prompted press Enter to continue.

Once the repository is enabled, install Python 3.7 with:

- `sudo apt install python3.7`

Then we will need to install the libraries below used by the project by entering in console the following command:

- `pip3 install pandas matplotlib keras scikit-learn numpy more-tertools seaborn xgboost`

Clone the repository locally by entering in con-

sole the following command:

- `git clone https://github.com/Kkalais/Music_Genre_Classification.git`

In order to run the code using the above-mentioned algorithms just enter in console the following commands :

- `python3 main.py name_of_the_algorithm`

There is also a mode that runs all four algorithms consecutively, and produces a bar plot to compare the algorithms' results. Please enter in console:

- `python3 main.py comparative`

## 5. Performance Comparison / Results

Our baseline was single-threaded CPU code. The libraries in python and sklearn are commonly used for machine learning purposes, so they are the de facto standard. The characteristics of our system are:

- Processor: Intel Core i7-4510U CPU @ 2.00GHz x4
- RAM: 8 Gb
- Cores: 4
- Sockets: 1
- OS type: Ubuntu 19.04

The quantitative metric which we used to judge our models is the combination of accuracy (that is, percentage of predicted labels which matched their true labels) and execution time. On the figures below (19, 20), we provide an (accuracy,time)-based comparison between the used algorithms in order to give a better look into which algorithm is considered to be more suitable.

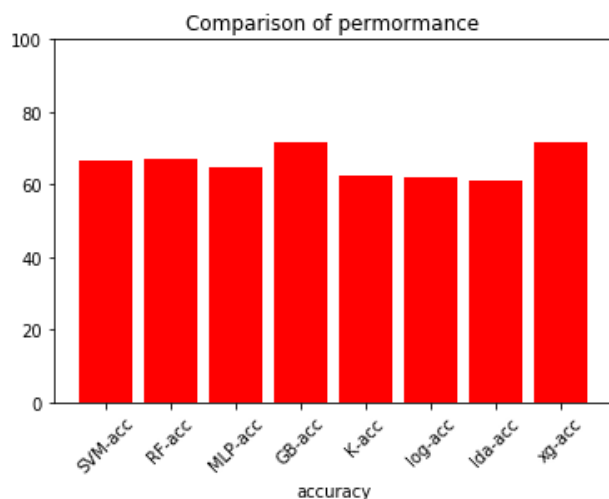


Figure 19. Accuracy comparison of the algorithms

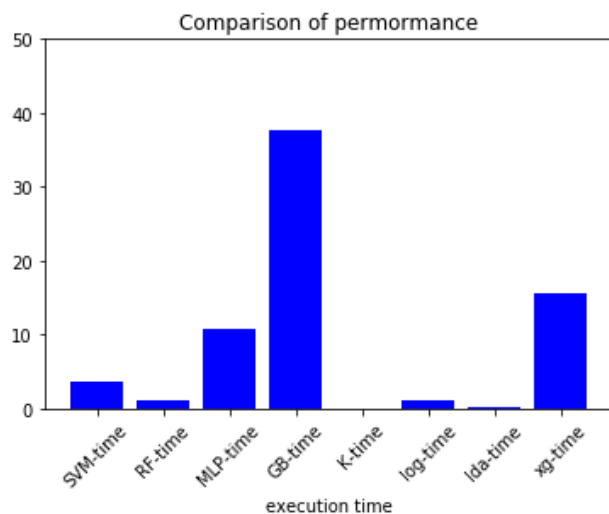


Figure 20. Execution time comparison of the algorithms

From the table below, it is obvious that Random Forest produces the better results combining the metrics of time and accuracy, and we also provide a bar plot to visually confirm that Random Forest performed better.

Another way of visualizing the performance of our best model is through a confusion matrix. A confusion matrix is a performance measurement for machine learning classification problem where output can be two or more classes. For a binary classification problem, it is a table with

Table 2. Classification accuracies and execution times

ALGORITHM	ACCURACY	EXECUTION TIME
GB	74.21%	39.77s
MLP	65.14%	11.21s
SVM	66.74%	3.97s
RF	68.74%	1.29s
KNN	62.37%	0.03s
LR	61.87%	1.09s
LDA	61.12%	0.29s
XGB	71.53%	15.57s

4 different combinations of predicted and actual values (figure 21).

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Figure 21. Confusion Matrix

- True positives (TP): In the example of predicting whether an email is spam (1) or (0), these are cases in which we predicted spam and it is also a spam email.
- True negatives (TN): We predicted non-spam, and it is not a spam email.
- False positives (FP): We predicted spam, but it is not spam.
- False negatives (FN): We predicted non-spam, but it is actually a spam one.

Here, as seen in the figure 22 Gradient Boosting algorithm did well with the rock, jazz and instrumental genre.

Additionally, it incorrectly classified some rock, as well as folk, as instrumental music. We should expect these misclassifications since rock

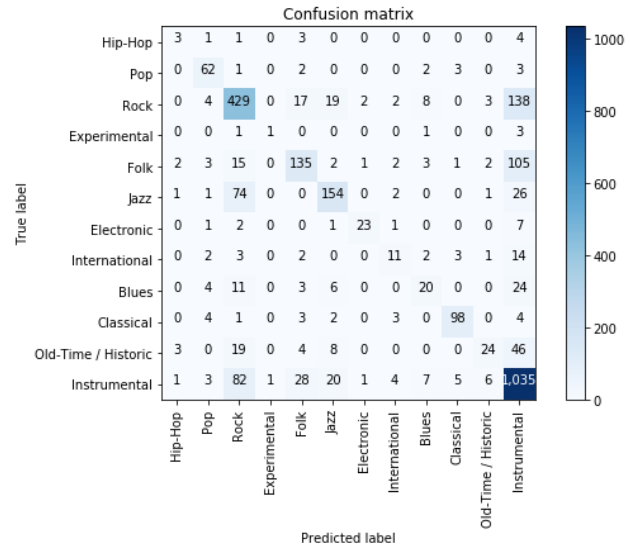


Figure 22. Confusion matrix for GB predictions

is a genre that both encapsulates many different styles (hard rock, progressive rock, indie rock, alternative rock, etc.) and heavily influences many other derivative genres. Finally, only few samples can be labeled as Hip-hop or Experimental music - classes with the least importance - so it would not make any difference if we had not included them in our dataset.

## 6. Future work

There are several directions that would be interesting to pursue in the future:

- Experiment with other types of deep learning methods - for example, Convolutional Neural Networks - that they will produce higher classification accuracies
- Try classifying music by another target-feature, e.g. by artist or by decade
- Including additional metadata text features such as album, song title, or lyrics could allow us to make an extension to music mood classification as well.
- If we are able to train a much deeper net-

work that attains high accuracy, it would be interesting to plot the accuracy of a classical algorithm using different activation layers.

- Examine the application of Long Short Term Memory (LSTM) model in music genre classification

## 7. Acknowledgments

The project was implemented in the spectrum of a special issue course of the University of Thessaly, Electronic and Computer Engineering Department. The professor and supervisor of the project, prof. Yota Tsompanopoulou guided us to the majority of the knowledge, tools and skills needed for this purpose.

## References

- [1] Hareesh Bahuleyan, *Music Genre Classification using Machine Learning Techniques*, University of Waterloo, ON, Canada, 2018
- [2] Danny Diekreoger, *Can Song Lyrics Predict Genre?*, Stanford University, 2012
- [3] Michal Defferrady, Kirell Benziy, Pierre Vanderghensty, Xavier Bresson, *Fma: a dataset for music analysis*, LTS2, EPFL, Switzerland / SCSE, NTU, Singapore, 2017
- [4] MARIUS-CONSTANTIN POPESCU, *Multilayer Perceptron and Neural Networks*, Faculty of Electromechanical and Environmental Engineering, University of Craiova
- [5] Cory McKay, Ichiro Fujinaga, *Musical genre classification: Is it worth pursuing and how can it be improved?*, Music Technology, Schulich School of Music, McGill University Montreal, Quebec, Canada
- [6] Jamal M. Nazzal, Ibrahim M. El-Emary, Salam A. Najim, *Multilayer Perceptron Neural Network (MLPs) For Analyzing the Properties of Jordan Oil Shale*, Al Ahliyya Amman University, P.O. Box 19328, Amman, Jordan, King Abdulaziz University, P.O. Box 18388, Jeddah, King Saudi Arabia
- [7] Marcin Andrychowicz, Misha Denil, Sergio Gomez Colmenarejo, Matthew W. Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, Nando de Freitas, *Learning to learn by gradient descent*, Google DeepMind, University of Oxford, Canadian Institute for Advanced Research
- [8] Weijie Wang, Yanmin Lu, *Analysis of the Mean Absolute Error (MAE) and the Root Mean Square Error (RMSE) in Assessing Rounding Model*, CFETS Information Technology (Shanghai) Co., Ltd. China Foreign Exchange Trade System, Shanghai, China
- [9] Bolin Gao, Laca Pavel, *On the Properties of the Softmax Function with Application in Game Theory and Reinforcement Learning*
- [10] Thamer M. Jamel, Ban M. Khammas, *IMPLEMENTATION OF A SIGMOID ACTIVATION FUNCTION FOR NEURAL NETWORK USING FPGA*, Electrical Engineering Department, University of Technology, Baghdad, Iraq
- [11] ANTONY UGONI BSc (Hons). BRUCE F. WALKER DC., *AN INTRODUCTION TO PROBABILITY DISTRIBUTIONS*
- [12] Diederik P. Kingma, *ADAM: A METHOD FOR STOCHASTIC OPTIMIZATION*, University of Amsterdam, OpenAI
- [13] Zhilu Zhang, Mert R. Sabuncu, *Generalized Cross Entropy Loss for Training Deep Neural Networks with Noisy Labels*, Electrical and Computer Engineering, Meinig School of Biomedical Engineering, Cornell University

- [14] Ilya Sutskever, Oriol Vinyals, Quoc V. Le, *Sequence to Sequence Learning with Neural Networks*, Google
- [15] DURGESH K. SRIVASTAVA, LEKHA BHAMBHU, *DATA CLASSIFICATION USING SUPPORT VECTOR MACHINE*, Department of CSE/IT, BRCM CET, Bahal, Bhiwani, Haryana, India
- [16] Theodoros Evgeniou, Massimiliano Pontil, *WORKSHOP ON SUPPORT VECTOR MACHINES: THEORY AND APPLICATIONS*, Center for Biological and Computational Learning, Artificial Intelligence Laboratory
- [17] L.Gonzalez-Abril, *A note on the bias in SVMs for multi-classification*, Grupo COSDE, Departamento de Economia Aplicada I, Universidad de Sevilla
- [18] Claudio Gentile, *A New Approximate Maximal Margin Classification Algorithm*, DSI, Università di Milano, Via Comelico 39, 20135 Milano, Italy
- [19] Rimah Amami, Dorra Ben Ayed, Nouredine Ellouze, *Practical Selection of SVM Supervised Parameters with Different Feature Representations for Vowel Recognition*, Department of Electrical Engineering, National School of Engineering of Tunis University of Tunis El Manar, Tunisia
- [20] Marc G. Genton, *Classes of Kernels for Machine Learning: A Statistics Perspective*, Department of Statistics, North Carolina State University, Raleigh, NC 27695-8203, USA
- [21] Vasileios Apostolidis-Afentoulis, Konstantina-Ina Lioufi, *SVM CLASSIFICATION WITH LINEAR AND RBF KERNELS*, Department of Information Technology, Alexander TEI of Thessaloniki
- [22] Mathini Sellathurai, *The separability theory of hyperbolic tangent kernels and support vector machines for pattern classification*, Acoustics, Speech, and Signal Processing, IEEE International Conference
- [23] Jehad Ali, Rehanullah Khan, Nasir Ahmad, Imran Maqsood, *Random Forests and Decision Trees*, Computer Systems Engineering, UET Peshawar, Pakistan, Sarhad University of Science and Information Technology, Peshawar, Pakistan
- [24] Lior Rokach, Oded Maimon, *DECISION TREES*, Department of Industrial Engineering, Tel-Aviv University
- [25] J.R. QUINLAN, *Induction of Decision Trees*, Centre for Advanced Computing Sciences, New South Wales Institute of Technology, Sydney 2007, Australia
- [26] Alexey Natekin, Alois Knoll, *Gradient boosting machines, a tutorial*, Department of Informatics, Technical University Munich, Garching, Munich, Germany
- [27] Gongde Guo, Hui Wang, David Bell, Yaxin Bi, Kieran Greer, *KNN Model-Based Approach in Classification*, School of Computing and Mathematics, University of Ulster
- [28] Sadegh Bafandeh Imandoust, Mohammad Bolandraftar, *Application of K-Nearest Neighbor (KNN) Approach for Predicting Economic Events: Theoretical Background*, Department of Economics, Payame Noor University, Tehran, Iran
- [29] P. Viswanath, T. Hitendra Sarma, *An Improvement to K-Nearest Neighbor Classifier*, Department of Computer Science and Engineering, Rajeev Gandhi Memorial College of Eng. and Technology
- [30] Park, Hyeoun-Ae, *An Introduction to Logistic Regression: From Basic Concepts to Interpretation with Particular Attention to Nursing Domain*, College of Nursing and

System Biomedical Informatics National  
Core Research Center, Seoul National Uni-  
versity, Seoul, Korea

- [31] CHAO-YING JOANNE PENG, KUK LIDA LEE, GARY M. INGERSOLL, *An Introduction to Logistic Regression: Analysis and Reporting*, Indiana University-Bloomington
- [32] Alaa Tharwata, Tarek Gaberc, Abdelhameed Ibrahim, Aboul Ella Hassanie, *Linear discriminant analysis: A detailed tutorial*, Department of Computer Science and Engineering, Frankfurt University of Applied Sciences, Frankfurt am Main, Germany
- [33] Pouria Kaviani, Mrs. Sunita Dhotr, *Short Survey on Naive Bayes Algorithm*, International Journal of Advance Engineering and Research Development, Volume 4, Issue 11, November-2017
- [34] Tianqi Chen and, Carlos Guestrin, *XGBoost: Reliable Large-scale Tree Boosting System*, University of Washington
- [35] Tianqi Chen, Carlos Guestrin, *XGBoost: A Scalable Tree Boosting System*, University of Washington
- [36] Jinghui Ma, *Application of the XGBoost Machine Learning Method in PM2.5 Prediction: A Case Study of Shanghai*, January 2019, Aerosol and Air Quality Research 20(1)
- [37] Zhiyong Yan, Congfu Xu, *Using Decision Boundary to Analyze Classifiers*, Intelligent System and Knowledge Engineering, 2008. ISKE 2008. 3rd International