

# **Automatic Music Genre Classification using Ensemble Learning**

by

**Orgest Xhelili**

Bachelor Thesis in Computer Science

Prof. Dr. Herbert Jaeger  
Bachelor Thesis Supervisor

Date of Submission: May 7, 2018

With my signature, I certify that this thesis has been written by me using only the indicates resources and materials. Where I have presented data and results, the data and results are complete, genuine, and have been obtained by me unless otherwise acknowledged; where my results derive from computer programs, these computer programs have been written by me unless otherwise acknowledged. I further confirm that this thesis has not been submitted, either in part or as a whole, for any other academic degree at this or another institution.

Signature

Place, Date

## **Abstract**

Organizing music libraries and databases requires automatic classification of music. Music genres are essential descriptors created by humans to label pieces of music. This paper presents a novel approach of ensemble learning for solving the problem of genre classification. Ensemble learning is a machine learning paradigm where multiple models are trained to solve the same problem and the final result is a combination of the individual results by each model. Our system consists of three different classifiers: echo state networks, support vector machines and logistic regression.

Echo state networks (ESNs) provide a supervised learning principle for analyzing and training recurrent neural networks (RNNs). They present a more practical approach to the usage of RNNs because of their computational efficiency and simple implementation. Support vector machines (SVMs) are supervised learning models which can be used for regression analysis or classification problems.

Logistic regression is a regression model which performs linear classification by estimating probabilities of each of the outcomes.

Contents

1 Introduction 1

2 Statement and Motivation of Research 2

2.1 Music Genre Classification 2

2.2 Audio Features 2

2.3 Extracted features 3

2.4 Research objectives 4

3 Ensemble of Classifiers 5

3.1 Echo State Networks 5

3.2 Support Vector Machines 6

3.3 Logistic Regression 9

3.4 Ensemble Learning 9

4 Documentation of Methods 9

4.1 Problem and Dataset Description 9

4.2 State of the Art 10

4.3 Pre-processing 10

4.4 Cross Validation of Parameters 11

4.5 ESN Procedure 11

# 1 Introduction

Music genres are labels created by humans to classify pieces of music giving rise to organized music libraries. However, the distinction between genres remains not well defined, as does their definition. Because of the vast amount of musical pieces, automatic genre classification is crucial for organizing huge music databases allowing users to find the music they want to. Automatic genre classification is part of automatic music information retrieval where much research has been ongoing in the 21st century. A summary of state of the art automatic music genre classification may be found on this survey [1].

Echo state networks provide an architecture and supervised learning principle for recurrent neural networks. They present a practical approach to training RNNs because of their simple implementation and computational efficiency. RNNs represent a large and varied class of computational models that are designed by more or less detailed analogy with biological brain modules. In an RNN numerous abstract neurons (also called units or processing elements) are interconnected by likewise abstracted synaptic connections (or links), which enable activations to propagate through the network [2]. ESNs differ from classic RNNs because of the usage of a fixed random reservoir and the training of only output neurons weights [3]. Because there are no cyclic dependencies between the trained readout connections, training an ESN becomes a simple linear regression task [4].

Support vector machines are a set of supervised learning methods used for classification, regression and outliers detection. Given training data labeled in two classes, the SVM algorithm computes an optimal hyperplane to categorize new data. Thus a SVM behaves as a non-probabilistic binary linear classifier. However, they can be extended to support multi class classification of non-linearly separable data in a probabilistic setting. To support non-linear classification, SVMs use the so called kernel trick, implicitly mapping their inputs into high-dimensional feature spaces. To support multi class classification, we use the one-vs.-one strategy, reducing a  $K$  class classification problem to  $K(K - 1)/2$  binary classifiers. Probabilities can be estimated for binary classification using Platt Scaling [5]. For multi class classification the probability scores can be estimated by using a pairwise coupling strategy [6]. It should be noted that the estimation of probabilities is rather an expensive procedure. SVMs are widely used nowadays for different machine learning problems because of their effectiveness, versatility and efficiency.

Logistic regression sometimes called the logistic model or logit model, analyzes the relationship between multiple independent variables and a categorical dependent variable by estimating probabilities of occurrence of an event [7]. In its classical form, logistic regression is used as a binary classifier. In this case the dependent variable is dichotomous. When the dependent variable is not dichotomous, a multinomial logistic regression can be applied.

The report is organized as follows. The second section states the motivation of this research and describes in more detail the problem of the automatic music genre classification problem. The third section gives a more in-depth analysis of ESNs, SVMs and logistic regression. The fourth section documents the ensemble learning strategy used to solve the problem of automatic music genre classification. This section also gives a summary of the dataset used and the evaluation criteria. The fifth section documents the results of our experiment. The last section discusses the results and potential focus areas for future research in the topic.

## 2 Statement and Motivation of Research

### 2.1 Music Genre Classification

The basis of automatic music genre classification is the representation of musical pieces and the extraction of feature vectors from the agreed representation. A symbolic representation of musical pieces is rarely available so one has to deal with audio samples. This representation presents a good strategy as it places music genre classification problem in the sphere of speech recognition where more research has been done. However, using the exact waveform of a musical piece for automatic classification is not feasible, because of the low level information contained in the audio samples. Therefore the first step in automatic genre classification is the extraction of useful features from the audio representation of musical pieces. Much work on extraction of features from music has been devoted to timbral texture features, rhythmic content features and pitch content features [8].

### 2.2 Audio Features

- Timbral texture features are used to differentiate sounds with the same pitch and loudness from each other. The use of timbral texture features originates from speech recognition [9]. A detailed list of features used to characterize timbre may be found in [10]. These features are usually referred to as being low-level as they represent sound on samples of milliseconds. We summarize here the main low-level features used in genre classification problems as given in [1]:
  - temporal features: features computed from the audio signal frame (zero-crossing rate and linear prediction coefficients)
  - energy features: features that describe the energy content of the signal
  - spectral shape features: features that describe the shape of the power spectrum of a signal frame: centroid, spread, skewness, kurtosis, slope, roll-off frequency, Mel-frequency cepstral coefficients (MFCC)
- Rhythmic content features are the most widely used mid-level features in audio-based music classification [11]. These features are used to characterize the temporal regularity of a musical piece. A review of automatic rhythm description systems may be found in [12]. Tzanetakis and Cook proposed the calculation of rhythmic content features for genre classification by extracting periodic changes from the beat histogram [13]. The beat histogram models the distributions of the regularities exhibited in the envelop signal, where rhythmic features can be obtained such as magnitudes and locations of dominant peaks and BPM(beat-per-minute) [11].
- Pitch content features are another set of important mid-level music features. Pitch is determined by what the ear judges to be the most fundamental frequency of the sound [14]. However, the perception of pitch is completely subjective which makes pitch not equal to the fundamental frequency. Tzanetakis and Cook proposed the calculation of pitch content features by extracting periodic changes from the pitch histogram [13]. The pitch histogram models the distribution of candidate pitches extracted by all frames [11].

## Analysis and Texture Window

Most of the low-level features are computed at regular time intervals, over short windows of length 10-100ms. These segments, called analysis windows, have to be small enough so that the frequency characteristics of the magnitude spectrum are relatively stable [13]. However, in order to capture the long term nature of sound "texture", means and variances of the extracted features can be computed over a number of analysis windows. This larger window is referred to as a texture window [13]. In our system, an analysis window of 100ms and a texture window of 2s is used.

### 2.3 Extracted features

This paper uses a variety of low-level features to extract from the audio representation of musical pieces. As documented in two surveys for automatic music classification [1, 11], the low-level features (especially MFCCs) give the best results for the problem of genre classification. We provide below a description of each of the features used in our experiment.

- **MFCC** - This paper uses Mel-frequency cepstral coefficients (MFCCs) as the most important feature to extract from the audio representation of musical pieces. MFCCs are a feature set popular in audio processing. MFCCs fall in the category of timbral texture features that describe the shape of the power spectrum of a signal frame. MFCCs features are based on the short time Fourier transform (STFT). This feature set is obtained as follows: We first frame the audio signal in multiple windows. For each frame the logarithm of the amplitude spectrum based on STFT is computed. Usually the frequencies are divided in 13 bins using Mel-frequency scaling giving rise to 13 coefficients. However, the usage of 20 coefficients has shown to give better results [15]. Therefore we are using 20 MFCCs in this experiment. In the end we apply Discrete Cosine Transform to obtain the feature vectors.
- **Short-time Energy** - The short-time energy is the signal energy over a certain audio frame. It is computed as:

$$E = \sum_{m=1}^N x^2(m) \quad (1)$$

where  $E$  represents the energy of the signal  $x(m)$  [16].

- **Short-time Energy Entropy** - This feature can be interpreted as a measure of abrupt changes in the signal energy. It is computed as:

$$H = - \sum_{t=1}^N E_n[t] \cdot \log_2(E_n[t]) \quad (2)$$

where  $E_n[t]$  is the normalized energy over a sub-frame  $t$ .

- **Zero Crossing Rate** - This feature describes the rate at which the signal changes its sign. It is computed as:

$$zcr = \frac{1}{N} \sum_{t=1}^N \text{sign}(x_t x_{t-1}) \quad (3)$$

where  $x$  is a signal of length  $N$  and  $\text{sign}$  is an indicator function.

- **Spectral Centroid** - The spectral centroid is defined as the center of gravity of the magnitude spectrum of STFT [13]. It is computed as:

$$C_t = \frac{\sum_{n=1}^N M_t[n] \cdot n}{\sum_{n=1}^N M_t[n]} \quad (4)$$

where  $M_t[n]$  is the magnitude of STFT at frame  $t$  and frequency bin  $n$ .

- **Spectral Spread** - The spectral spread is defined as the variance of the magnitude spectrum of STFT. It is computed as:

$$S_t = \sqrt{\sum_{n=1}^N \frac{(C_t - n)^2 \cdot M_t[n]}{M_t[n]}} \quad (5)$$

where  $C_t$  and  $M_t[n]$  are defined as in equation 4.

- **Spectral Entropy** - This feature is similar to Energy Entropy. In contrast the entropy is calculated over sub-frames of the spectrum, not the signal itself.
- **Spectral Flux** - The spectral flux is defined as the squared difference between the normalized magnitudes of successive spectral distributions [13]. It is computed as:

$$F_t = \sum_{n=1}^N (N_t[n] - N_{t-1}[n])^2 \quad (6)$$

where  $N_t[n]$  denotes the normalized magnitude of the Fourier transform at frame  $t$  and frequency bin  $n$ .

- **Spectral Rolloff** - The spectral rolloff is defined as the frequency  $R_t$  below which 85% of the magnitude distribution is concentrated [13]. It is computed as:

$$\sum_{n=1}^{R_t} M_t[n] = 0.85 \cdot \sum_{n=1}^N M_t[n] \quad (7)$$

where  $M_t[n]$  is defined as in equation 4.

- **Delta-Cepstral Coefficients** - These features serve as the first order derivative of MFCCs and are appended to MFCCs as documented below.  
Given MFCCs  $C[n]$  over an analysis window, we calculate the delta features

$$D[n] = C[n] - C[n - 1] \quad (8)$$

where  $n$  is the index of the analysis frame. It has been shown that the usage of delta-cepstral coefficients produces better results in the field of speech recognition [17].

The extraction of features for the experiment was done using pyAudioAnalysis [18].

## 2.4 Research objectives

I state that this guided research does not aim to achieve state of the art solution to the problem of music genre classification. The main objective of the research is to apply a novel ensemble learning approach to this problem, by combining different classifiers. In the end, we provide a comparison of the performances of the different classifiers.



### 3 Ensemble of Classifiers

#### 3.1 Echo State Networks

The following section documents a summary of ESNs as introduced in [3, 19, 20].

Echo state networks provide an architecture and supervised learning principle for recurrent neural networks. The main idea is to drive a random, large, fixed recurrent neural network with the input signal, thereby inducing in each neuron within this reservoir network a nonlinear response signal, and combine a desired output signal by a trainable linear combination of all of these response signals [3].

We consider discrete-time neural networks with  $K$  input units,  $N$  internal network units and  $L$  output units. Activations of input units at time step  $n$  are  $K$ -dimensional vectors  $\mathbf{u}(\mathbf{n}) = (u_1(n) \dots u_K(n))^T$ , of internal units  $N$ -dimensional vectors  $\mathbf{x}(\mathbf{n}) = (x_1(n) \dots x_N(n))^T$ , and of output units  $L$ -dimensional vectors  $\mathbf{y}(\mathbf{n}) = (y_1(n) \dots y_L(n))^T$ . Real-valued connection weights are collected in a  $N \times K$  weight matrix  $\mathbf{W}^{\text{in}}$  for the input weights, in a  $N \times N$  matrix  $\mathbf{W}$  for the internal connections, and in a  $L \times N$  matrix  $\mathbf{W}^{\text{out}}$  for the connections to the output units. The basic architecture of the network is shown in Figure 1.

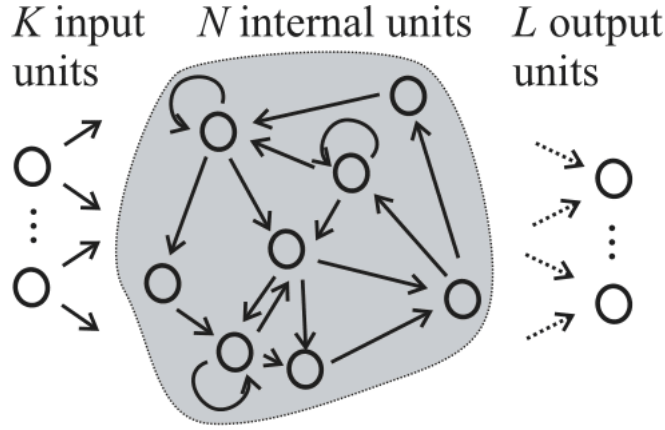


Figure 1: Architecture of a network with  $K$  input units,  $N$  internal units and  $L$  output units (graphics taken from [20]).

In this project we are using the version of ESNs with leaky integration as described in [21]. Using this particular type of the network the activation of internal units follows the equation:

$$\mathbf{x}(n+1) = (1 - \alpha)\mathbf{x}(n) + \alpha\mathbf{f}(\mathbf{W}\mathbf{x}(n) + \mathbf{W}^{\text{in}}\mathbf{u}(n+1) + \mathbf{b}_N) \quad (9)$$

where  $\mathbf{b}_N$  is a fixed random  $N$  dimensional bias vector,  $\alpha$  is the leaking rate of the network and  $\mathbf{f} = (f_1 \dots f_N)$  are the internal unit's output functions. In our case we will be using  $\tanh$  as this function. The output is updated according to the equation:

$$\mathbf{y}(n+1) = g(\mathbf{W}^{\text{out}}\mathbf{x}(n+1)) \quad (10)$$

where  $g$  is an output activation function. In our case we will be using the identity function so  $g$  will be ignored.

## Training ESNs

In the stage of training, the network is driven by a given teacher input signal  $\mathbf{u}(1), \dots, \mathbf{u}(n_{train})$  which yields a sequence  $\mathbf{x}(1), \dots, \mathbf{x}(n_{train})$  of network states. The obtained network states are collected into a matrix  $\mathbf{X}$  of size  $N \times (n_{train} - n_0)$ , where  $n_{train}$  is the training length and  $n_0$  is the washout time. Along the teacher input signal, a teacher output signal is provided. The desired outputs  $\mathbf{y}(n)$  are sorted row-wise into a teacher output collection matrix  $\mathbf{Y}_t$  of size  $(n_{train} - n_0) \times L$  [3]. Then the output weights  $\mathbf{W}^{out}$  are trained using linear regression following the equation:

$$\mathbf{W}^{out} = \mathbf{Y}_t \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \gamma^2 \mathbf{I})^{-1} \quad (11)$$

where  $\gamma^2$  is the regularization coefficient and  $\mathbf{I}$  is the identity matrix. After training the output weight by linear regression the network is ready for testing.

## Achieving Echo States

Under certain conditions, the activations state  $\mathbf{x}(n)$  of a recurrent neural network are a function of the input history presented to the network [19]. More precisely there exists an echo function  $\mathbf{E}$ , such that the activation state of the network is calculated as:

$$\mathbf{x}(n) = \mathbf{E}(\dots, \mathbf{u}(n-1), \mathbf{u}(n)) \quad (12)$$

We give here a proposition from [19] which gives sufficient conditions for the existence and the non-existence of echo states.

**Proposition 1** *Assume a sigmoid network with unit output functions  $f_i = \tanh$ . Let the weight matrix  $\mathbf{W}$  have its largest singular value  $\sigma_{max} < 1$ . Then the network has echo states for all admissible inputs. Let the weight matrix have a spectral radius  $|\lambda_{max}| > 1$ , where the spectral radius denotes the eigenvalue of the weight matrix with the largest absolute value. Then this network has no echo states if  $\mathbf{u}(n) = 0$  is an admissible input sequence.*

Note that scaling the weight matrix  $\mathbf{W}$  with a scalar  $\alpha$  scales the spectral radius and the singular values accordingly. In practice it is enough to achieve echo states if we have a spectral radius marginally smaller than 1. This would be achieved by scaling the matrix  $\mathbf{W}$  with  $\beta/|\lambda_{max}|$  where  $\beta$  is marginally smaller than 1.

## 3.2 Support Vector Machines

Support Vector Machines (SVMs) are a popular machine learning method for classification, regression and other learning tasks. The original SVM algorithm was invented by Vladimir N. Vapnik and Alexey Ya. Chervonenkis in 1963. The current standard version was initially proposed by Cortes and Vapnik in 1993 and published in 1995 [22]. The following sections explain the SVM algorithm as presented in [23, 24, 25].

## Motivation behind SVMs

The basic idea of the SVM algorithm is to find an optimal hyperplane for linearly separable data points. Given linearly separable training data points  $\mathbf{x}_i \in R^n, i = 1, \dots, n$ , where each  $\mathbf{x}_i$  belongs in one of the two classes  $y_i = -1$  or  $+1$ , the SVM algorithm tries to find an optimal hyperplane separating the two classes [23]. This hyperplane can be defined formally as:

$$\mathbf{w} \cdot \mathbf{x} + \mathbf{b} = 0 \quad (13)$$

where  $\mathbf{w}$  is a weight vector normal to the hyperplane and  $\mathbf{b}$  is a bias vector. Support vectors are the data points closest to this separating hyperplane and the aim of the algorithm is to maximize the distance of the hyperplane to the support vectors [24]. This distance is called the margin of the separating hyperplane.

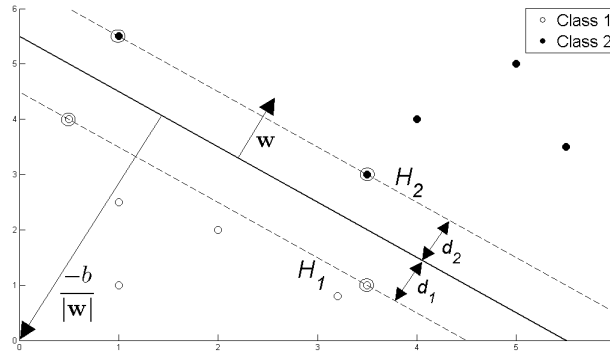


Figure 2: Hyperplane through two linearly separable classes (graphics taken from [23]).

The figure above shows the optimal hyperplane through two linearly separable classes.  $H_1$  and  $H_2$  denote the planes on which the support vectors lie on. The margin of the hyperplane is  $m = |d_1| + |d_2|$ . Referring to this figure, implementing the SVM boils down to selecting variables  $\mathbf{w}$  and  $\mathbf{b}$  such that:

$$\mathbf{x}_i \cdot \mathbf{w} + \mathbf{b} \geq +1 \quad \forall y_i = +1 \quad (14)$$

$$\mathbf{x}_i \cdot \mathbf{w} + \mathbf{b} \leq -1 \quad \forall y_i = -1 \quad (15)$$

The equations can be combined into:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + \mathbf{b}) - 1 \geq 0 \quad \forall i \quad (16)$$

As we said earlier the SVM algorithm aims to maximize the margin  $m$ . Simple vector geometry shows that  $m = \frac{2}{\|\mathbf{w}\|}$  [23]. Maximizing  $m$  is equivalent to minimizing  $\frac{1}{2}\|\mathbf{w}\|^2$ , which leads to solving the following constrained optimization problem:

$$\min \frac{1}{2}\|\mathbf{w}\|^2 \quad s.t. \quad y_i(\mathbf{x}_i \cdot \mathbf{w} + \mathbf{b}) - 1 \geq 0 \quad \forall i \quad (17)$$

The solution to this problem and the mathematical formalism behind it are documented in [24]. We ignore the step by step solution and present just the result. The original problem [17], known as the primal problem, is equivalent to the following dual problem [24]:

$$\max W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \quad s.t. \quad \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0 \quad (18)$$

This is a quadratic programming problem and a solution  $\alpha = (\alpha_i)$  can be found using different approaches. For SVM, sequential minimal optimization seems to be the most popular [24]. A nice characteristic of the solution is that many of the  $\alpha_i$  are 0 and the  $\mathbf{x}_i$  with non-zero  $\alpha_i$  are the support vectors. Therefore we can construct  $\mathbf{w}$  as a linear combination of a small number of data points [24]:

$$\mathbf{w} = \sum_{j=1}^s \alpha_{t_j} \mathbf{y}_{t_j} \mathbf{x}_{t_j} \quad (19)$$

where  $t_j (j = 1, \dots, s)$  are the indices of the  $s$  support vectors. After the optimal hyperplane is computed, the SVM classifier can categorize new data  $\mathbf{z}$  by looking at the sign of the expression:

$$f = \mathbf{w} \cdot \mathbf{z} + \mathbf{b} \quad (20)$$

and categorize it as class 1 if the  $f$  is positive and class  $-1$  otherwise.

### Extension to non-linearly separable data

In order to extend the SVM methodology to handle data that is not fully linearly separable, we relax the constraints 14 and 15 to allow for misclassified points [23]. This is done by introducing positive slack variables  $\xi_i$ ,  $i = 1, \dots, n$  and the new equations can be combined into:

$$\mathbf{y}_i(\mathbf{x}_i \cdot \mathbf{w} + \mathbf{b}) - 1 + \xi_i \geq 0 \quad \text{where} \quad \xi_i \geq 0 \quad \forall i \quad (21)$$

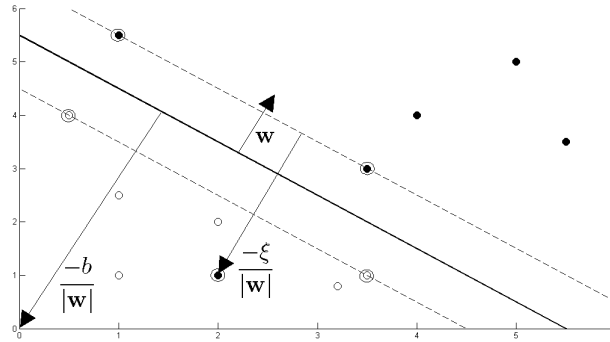


Figure 3: Hyperplane through two non-linearly separable classes (graphics taken from [23]).

In this SVM version, data points on the incorrect side of the margin boundary have a penalty that increases with the distance from it. The correctly classified data points have  $\xi_i = 0$ . As we are trying to reduce the number of misclassifications, we adapt the previous objective function 17 in the following way:

$$\min \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \quad \text{s.t.} \quad \mathbf{y}_i(\mathbf{x}_i \cdot \mathbf{w} + \mathbf{b}) - 1 + \xi_i \geq 0 \quad \forall i \quad (22)$$

where the parameter  $C$  controls the trade-off between the slack variable penalty and the size of the margin. The new constrained optimization problem can be solved in the same

way as in the linear separable case, except there is an upper bound  $C$  on  $\alpha_i$  now [24]. The dual problem in this case is:

$$\max W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j \mathbf{y}_i \mathbf{y}_j \mathbf{x}_i \mathbf{x}_j \quad s.t \quad C \geq \alpha_i \geq 0, \sum_{i=1}^n \alpha_i \mathbf{y}_i = 0 \quad (23)$$

The weight vector  $\mathbf{w}$  can be constructed in the same way as in the linearly separable case.

### Extension to non-linear decision boundary

Many datasets cannot be separated by a linear hyperplane. The idea is to transform these data points into a higher dimensional feature space which is linearly separable. After the transformation, the SVM algorithm can be applied to the transformed data points. Computation in this new space can be quite costly but this is not a problem for the SVM, since we only calculate the inner product of the data points. As long as we can calculate the inner product in the feature space, we do not need the mapping explicitly [24]. This is the so called kernel trick. We define the kernel function  $K$ :

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i) \phi(\mathbf{x}_j) \quad (24)$$

where  $\phi(\mathbf{x})$  is a map which transforms the data points into the linearly separable feature space. For the linear decision boundary case,  $\phi$  is the identity function. The dual problem in the case of non-linearly decision boundary transforms to:

$$\max W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j \mathbf{y}_i \mathbf{y}_j K(\mathbf{x}_i, \mathbf{x}_j) \quad s.t \quad C \geq \alpha_i \geq 0, \sum_{i=1}^n \alpha_i \mathbf{y}_i = 0 \quad (25)$$

The new data  $\mathbf{z}$  can be categorized by looking at the sign of the expression:

$$f = \mathbf{w} \phi(\mathbf{z}) + \mathbf{b} = \sum_{j=1}^s \alpha_{t_j} \mathbf{y}_{t_j} K(\mathbf{x}_{t_j}, \mathbf{x}) + \mathbf{b} \quad (26)$$

where  $t_j$  is defined as in 19.

### 3.3 Logistic Regression

### 3.4 Ensemble Learning

## 4 Documentation of Methods

### 4.1 Problem and Dataset Description

The experiments are based on the famous standard dataset used by Tzanetakis and Cook [13]. This dataset consists of 1000 30s long audio files which have been annotated a genre. The audio files cover 10 different genres with 100 files per each genre. The

ten represented genres are: Blues, Classical, Country, Disco, Hip-hop, Jazz, Metal, Pop, Reggae and Rock.

Given this dataset, we aim to train three different classifiers and apply a 10-fold cross validation scheme to classify the dataset into the given genres. Performance of each classifier is analyzed. In the end we combine the three classifiers by a soft-voting strategy and analyze the performance of the ensemble.

## 4.2 State of the Art

In our context, state of the art refers to the performance of the classifiers applied on the dataset that we are using. The following table, taken from [11], compares the performance of different algorithms using different type of features. The first results come from the creators of the dataset who achieved up to 60% classification accuracy. The most used classifier are support vector machines. The MFCC feature set is used in almost every occasion. Most of the documented classifiers achieve mean classification accuracy in the range of 70-80%.

Reference	Features	Classifier	Accuracy (%)
[1] <sup>a</sup>	$\{\text{STFT+MFCC}\} \times \text{MuVar+beat+pitch}$	K-NN	60
[1] <sup>b</sup>	$\{\text{STFT+MFCC}\} \times \text{MuVar+beat+pitch}$	GMM	61
*[79]	$\{\text{MFCC}\} \times \text{FP}$	SVM	$77.7 \pm 2.8$
*[113] <sup>a</sup>	$\{\text{MFCC}\} \times \text{GMM}$	K-NN	$70.6 \pm 3.0$
*[113] <sup>b</sup>	$\{\text{MFCC}\} \times \text{GMM}$	SVM	$70.4 \pm 3.1$
[12] <sup>a</sup>	$\{\text{STFT+MFCC}\} \times \text{MuVar+beat+pitch}$	SVM	$72 \pm 5.1$
[12] <sup>b</sup>	$\{\text{STFT+MFCC}\} \times \text{MuVar}$	SVM	$71.8 \pm 4.8$
[12] <sup>c</sup>	$\text{DWCH+STFT+MFCC} \times \text{MuVar}$	SVM	$78.5 \pm 4.1$
*[35]	$\text{MFCC} \times \text{MuCov}$	SVM	$78.6 \pm 2.4$
[84]	$\text{STFT+MFCC} \times \text{MuVar}^2$	SVM	79.8
[9]	$\text{STFT+FFT+MFCC+LPC}$	AdaBoost.DT	82.4
[16]	$\text{CR} \times \text{NTF}$	SVM	$78.2 \pm 3.8$
[18]	$\{\text{MFCC+ASE+OSC}\} \times \text{FP} \times \text{LDA}$	NC	$90.6 \pm 3.1$
[19]	$\text{CR} \times \text{NTF}$	SRC	$92.4 \pm 2.0$
[123]	$\{\text{MFCC+ASE+OSC}\} \times \{\text{MuCov,FP}\} + \text{beat+chord}$	SVM (MKL)	90.4
[123]	$\{\text{MFCC+ASE+OSC}\} \times \{\text{MuCov,FP}\} + \text{beat+chord}$	SVM (SG)	90.9

Figure 4: Performance comparison on the standard dataset

## 4.3 Pre-processing

The procedure of feature extraction differs for ESNs and the other two classifiers. We feed continuous live data to the network while the features fed to the SVM and logistic regression algorithm are 'global' features computed for the whole song. We document below the two different strategies used in the pre-processing phase.

- ESNs - For each audio file, the features described in 2.3 are computed for each analysis window of 100ms long. We are calculating 20 MFCCs, 20 delta-cepstral coefficients and 8 single features. After the feature extraction procedure, each audio file is represented as a matrix  $\mathbf{A}$  of size  $48 \times f$ , where  $f$  is the number of frames of length 100ms for each audio file.

- **SVMs & Logistic Regression** - The extraction of 'global' features for these two classifiers is based upon the above procedure for ESNs. We compute the mean and the variance of the 48 aforementioned features (calculated for each analysis window) over each texture window of 2s. Therefore, 96 features are computed for each texture window. In the next step, the mean of these features is computed over the whole song segment, providing a global feature set of size 96. In the end, we append to this feature set the covariance matrix of the 20 MFCCs calculated over each analysis window. The usage of the covariance matrix of MFCCs led to better results using SVMs as a classifier [To-Add]. Finally, we have computed for each song a feature set of size 496 which is fed to both the SVM and logistic regression algorithms.

#### 4.4 Cross Validation of Parameters

A 10-fold cross validation is used to determine the parameters of the classifiers that achieve the best performance. The dataset is divided in 10 random smaller sets of 100 audio files each. Each of these subsets is used iteratively as a validation set, while the remaining data is used for training each of the classifiers. For each classifier configuration, the mean classification accuracy over the 10 validation sets is computed and the configurations with the highest accuracy are chosen as the optimal ones. After the parameters of each classifier are decided, the weights of each classifier in the ensemble system are decided in the same fashion.

#### 4.5 ESN Procedure

The feature set fed to the network is of size 48 so our network consists of  $K = 48$  input units. We are using the network to classify each musical piece into one of the 10 genres. Therefore, the network will  $L = 10$  output units.

##### Constructing the network

Weight matrices  $\mathbf{W}^{\text{in}}$ ,  $\mathbf{W}$  and the bias vector  $\mathbf{b}_N$  are randomly initialized with real numbers in the interval  $[-0.5, 0.5]$ . The spectral radius  $\lambda_{\max}$  of the weight matrix  $\mathbf{W}$  is computed and we scale the matrix with the inverse of this factor. After this operation,  $\mathbf{W}$  will have a spectral radius of 1.

##### Training the Network

For each run, the network is trained using the state update equation 9. For each musical piece, we ignore the first 30 frames from its matrix representation  $\mathbf{A}$ . Upon starting training a new musical piece the network state is initialized to 0. The system states are collected in a matrix  $\mathbf{X}$  of size  $N \times (n_{\text{train}} - n_0)$ . In our case  $n_{\text{train}} = 900 \cdot f$  and  $n_0 = 900 \cdot 30$ , where  $f$  is the number of frames for each musical piece. The output teacher matrix  $\mathbf{Yt}$  will be of size  $L \times (n_{\text{train}} - n_0)$ . Each column vector of  $\mathbf{Yt}$  consists of zeros and a 1 on the right genre. By using the equation 11 the output weights matrix  $\mathbf{W}^{\text{out}}$  is computed.

## Testing the network

Using state update equations 9 and 10, the network is tested on each validation set. For each musical piece in the validation set, the computed output vectors are saved in a matrix of size  $L \times (f - 30)$ . The mean of this matrix over each row is computed resulting in a  $L$  dimensional vector. The index of the highest value is chosen as the predicted genre. For each run we compute the classification accuracy by dividing the number of correct predictions by 100.



## References

- [1] N. Scaringella, G. Zoia, and D. Mlynek. Automatic genre classification of music content: a survey. *IEEE Signal Processing Magazine*, 23(2):133–141, March 2006.
- [2] Mantas Lukoeviius and Herbert Jaeger. Reservoir computing approaches to recurrent neural network training. *Computer Science Review*, 3(3):127 – 149, 2009.
- [3] H. Jaeger. Echo state network. *Scholarpedia*, 2(9):2330, 2007. revision #183563.
- [4] Herbert Jaeger and Harald Haas. Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless communication. *Science*, 304(5667):78–80, 2004.
- [5] John Platt et al. Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods. *Advances in large margin classifiers*, 10(3):61–74, 1999.
- [6] Ting-Fan Wu, Chih-Jen Lin, and Ruby C Weng. Probability estimates for multi-class classification by pairwise coupling. *Journal of Machine Learning Research*, 5(Aug):975–1005, 2004.
- [7] Hyeoun Park. An introduction to logistic regression: from basic concepts to interpretation with particular attention to nursing domain. *Journal of Korean Academy of Nursing*, 43(2):154–164, 2013.
- [8] T. Li and M. Ogihara. Toward intelligent music information retrieval. *IEEE Transactions on Multimedia*, 8(3):564–574, June 2006.
- [9] L. Rabiner and H. Juang. *Fundamentals of Speech Recognition*. Prentice Hall, 1993.
- [10] Geoffroy Peeters. A large set of audio features for sound description (similarity and classification) in the cuidado project. 2004.
- [11] Z. Fu, G. Lu, K. M. Ting, and D. Zhang. A survey of audio-based music classification and annotation. *IEEE Transactions on Multimedia*, 13(2):303–319, April 2011.
- [12] Fabien Gouyon and Simon Dixon. A review of automatic rhythm description systems. *Computer music journal*, 29(1):34–54, 2005.
- [13] G. Tzanetakis and P. Cook. Musical genre classification of audio signals. *IEEE Transactions on Speech and Audio Processing*, 10(5):293–302, Jul 2002.
- [14] Laura Williams Macy. *Grove music online*. Macmillan Reference, 2001.
- [15] M. Mandel and D. Ellis. Song-level features and svms for music classification. *Proc. Int. Conf. Music Information Retrieval*, 2005.
- [16] Madiha Jalil, Faran Awais Butt, and Ahmed Malik. Short-time energy, magnitude, zero crossing rate and autocorrelation measurement for discriminating voiced and unvoiced segments of speech signals. In *Technological Advances in Electrical, Electronics and Computer Engineering (TAECE)*, 2013 International Conference on, pages 208–212. IEEE, 2013.
- [17] Kshitiz Kumar, Chanwoo Kim, and Richard M Stern. Delta-spectral cepstral coefficients for robust speech recognition. In *Acoustics, Speech and Signal Processing (ICASSP)*, 2011 IEEE International Conference on, pages 4784–4787. IEEE, 2011.

- [18] Theodoros Giannakopoulos. pyaudioanalysis: An open-source python library for audio signal analysis. *PloS one*, 10(12), 2015.
- [19] H. Jaeger. The" echo state" approach to analysing and training recurrent neural networks-with an erratum note'. *German National Research Center for Information Technology GMD Technical Report*, 148, Jan 2001.
- [20] H. Jaeger. Short term memory in echo state networks. *GMD - German National Research Institute for Computer Science*, 152, Jan 2002.
- [21] H. Jaeger, M. Lukosevicius, D. Popovici, and U. Siewert. Optimization and applications of echo state networks with leaky- integrator neurons. *Neural Networks*, 20:335–352, April 2007.
- [22] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- [23] Tristan Fletcher. Support vector machines explained. *Tutorial paper*, 2009.
- [24] Martin Law. A simple introduction to support vector machines. *Lecture for CSE*, 802, 2006.
- [25] Robert Berwick. An idiots guide to support vector machines (svms). *Retrieved on October*, 21:2011, 2003.