

Unsupervised Music Genre Classification with a Model-Based Approach

Luís Barreira, Sofia Cavaco, and Joaquim Ferreira da Silva

CITI, Departamento de Informática
Faculdade de Ciências e Tecnologia
Universidade Nova de Lisboa
2829-516 Caparica, Portugal
lfbarreira@gmail.com, {sc,jfs}@di.fct.unl.pt

Abstract. New music genres emerge constantly resulting from the influence of existing genres and other factors. In this paper we propose a data-driven approach which is able to cluster and classify music samples according to their type/category. The clustering method uses no previous knowledge on the genre of the individual samples or on the number of genres present in the dataset. This way, music *tagging* is not imposed by the users' subjective knowledge about music genres, which may also be outdated. This method follows a model-based approach to group music samples into different clusters only based on their audio features, achieving a perfect clustering accuracy (100%) when tested with 4 music genres. Once the clusters are learned, the classification method can categorize new music samples according to the previously learned created groups. By using Mahalanobis distance, this method is not restricted to spherical clusters, achieving promising classification rates: 82%.

Keywords: automatic music genre classification, audio indexing, unsupervised classification.

1 Introduction

Since today's digital content development triggered the massive use of digital music, an indexing process is very important to guarantee a correct organization of huge databases. While a music genre categorization would be a solution, this may be hard to achieve (manually): on the one hand, music can be associated to one or more musical genres, and on the other hand, cultural differences and human interpretations, make it difficult the attainment of common music genre taxonomy. For example, an expert could label the Gustav Mahler's 2nd symphony as *Erudite - late Romantic*, while a non-expert could label it as *Classical*. Alternatively, an automatic classification based on good audio features may prevent the occurrence of incoherencies related to manual labeling.

While many supervised automatic music genre classifiers have been proposed, these will always be dependent on a previous manual labeling of the

data [4; 5; 6; 8; 9; 12; 14]. As a consequence, these will be unable to evolve with the data and automatically build new clusters driven by new values in the features. Alternatively, an unsupervised approach would not have this dependency and would be able to determine the genre of the music samples only based on their audio features. Nonetheless, only a few unsupervised methods have been proposed. Rauber et al. [10] proposed a growing hierarchical self-organizing map (which is a popular unsupervised artificial neural network) with psycho-acoustic features (loudness and rhythm) to obtain a hierarchical structuring music tree. Shao et al. [11] proposed an unsupervised clustering method that fed rhythmic content, Mel-frequency cepstral coefficients (MFCCs), linear prediction coefficients and delta and acceleration values (improvements in feature extraction) to a hidden Markov model.

Here we propose not only a methodology for unsupervised clustering but also for automatic music genre classification. The clustering method consists of a learning process that is able to cluster music samples based only on their audio properties and uses no previous knowledge on the genre of the training music samples. In addition a Model-Based approach is followed to generate clusters as we do not provide any information about the number of genres in the data set. The features used are related with rhythm analysis, timbre, and melody, among others. As these features represent a large number of dimensions, a feature reduction technique is necessary to reduce the dimensionality of the data. This clustering method achieves 100% accuracy results with classical, fado, metal and reggae music samples. After the clustering process is complete, the classification method can associate new test music samples to the previously created clusters. For that, the classifier uses Mahalanobis distance so that it can consider clusters with different shapes, volumes and orientations. An accuracy of 82% is achieved when classifying new music samples.

In the next section (Feature Extraction) we describe the features we use. Section 3 explains the Clustering Method while section 4 explains the Classification Method. The Results and Conclusions and Future Work are discussed in sections 5 and 6.

2 Feature Extraction

Feature extraction is the first step to be achieved in both automatic music genre clustering and classification. In this section, we describe the features we used, which can be grouped into two distinct groups: computational features (which do not represent any musical meaning and only describe a mathematical analysis over a signal) and perceptual features (which mathematically represent music properties based on the human hearing system).

Since some of the features we used have a very high dimensionality and it is more efficient to describe them with less dimensions, we used a set of statistical spectrum descriptors (SSD) proposed by Lidy and Rauber [7]. This set of descriptors includes: the mean, median, variance, skewness, kurtosis, min and

max-values. (Whenever this property is calculated, we mention it in the text below.)

Computational features are very popular and have been used in many automatic music genre classification studies [3; 4; 5; 6; 8; 9; 14]. To start with we use a set of *timbral texture features* proposed by Tzanetakis and Cook [14]. These include: spectral centroid (which is a measure of the centre of gravity of the magnitude spectrum), spectral roll-off (which corresponds to the frequency below which there is 85% of the energy of the magnitude spectrum), spectral flux (which accounts for the energy difference between successive frames of the spectrogram), zero-crossing rate (ZCR) (which is a measure of the number of times the audio waveform crosses the x -axis per time unit), and low energy (which is the percentage of frames that have lower energy than the average energy over the whole signal).

We also use the SSD of the MFCCs. The MFCCs are a very popular set of features based on the auditive human system that uses a Mel-frequency scale to group the frequency bins. In addition, three other features were also calculated: the root mean square of the spectrograms, which is an approximation of the volume (i.e., loudness) of the signal, the bandwidth, an energy-weighted standard deviation which measures the frequency range of the signal, and the uniformity, which measures the similarity of the energy levels in the frequency bands [3].

The spectral properties mentioned above can follow two different approaches: their values can be calculated over each window of a spectrogram or they can be calculated directly over the spectrum of the whole sound. Usually, these values are calculated over each window of the spectrogram, and that is the approach we used here. In addition, whenever we obtain a set of values with significant dimension (and we do not use their SSDs), we also use means and variances as features.

The perceptual features we used include rhythmic content, rhythm patterns and pitch content. The rhythmic content contains information such as the beat, the tempo, the regularity of the rhythm and time signature. In particular, the beat has been used in several studies on genre classification [4; 6] and it can be extracted from the beat histogram [14]. On the other hand, rhythm patterns represent the loudness sensation for several frequency bands in a time-invariant frequency representation [7]. We use both the SSDs and the rhythm histogram of the rhythm patterns. Finally, the pitch content is used to describe melody and harmony of a music signal. This feature is used quite often in genre classification leading to good accuracy results [4; 6; 13; 14; 15]. The pitch content can be extracted from the pitch histogram [15], and it includes the amplitudes and periods of the highest peaks in the histogram, pitch intervals between the two most prominent peaks and the overall sums of the histograms.

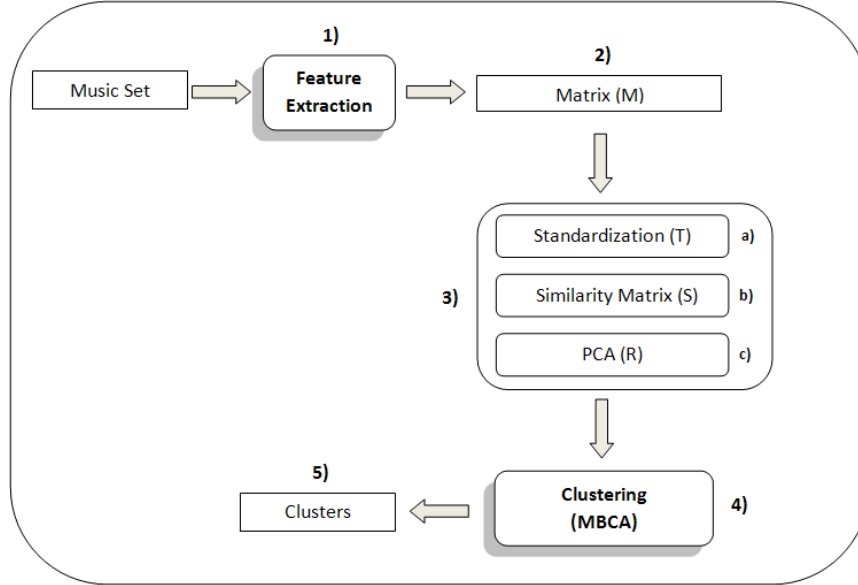


Fig. 1. The clustering method

3 The Clustering Method

The clustering method aims to organize several music samples into clusters without any initial information besides the feature set values of these samples. This method consists of several steps as illustrated in Fig. 1, which we describe below.

3.1 The Feature Reduction Stage

After the audio features have been extracted, they have to be analysed to find clusters of points with similar values. For that, the method starts by representing the features by data matrix \mathbf{M} , whose lines correspond to music samples in the training set and whose columns correspond to features. So, the $m_{s,f}$ cell of \mathbf{M} contains the value of the f th feature for music sample s .

Once this matrix is built, the method performs some transformations as illustrated by box 3 in Fig. 1. In order to set equal importance (scale) to all columns (features) of the data set matrix, in step 3a the method performs a standardization of matrix \mathbf{M} and creates a new matrix \mathbf{T} with the same dimension as matrix \mathbf{M} , that is, both matrices are $(N \times F)$, where N is the number of samples in the training set and F is the number of features. Now, the $t_{s,f}$ cell of \mathbf{T} contains the standardized value of feature f for music s , which is given by

$$t_{s,f} = \frac{m_{s,f} - m_{.,f}}{\sqrt{\text{var}(M_f)}}, \quad (1)$$

where $m_{.,f}$ is the mean value of the f th column of matrix \mathbf{M} , that is,

$$m_{.,f} = \frac{1}{N} \sum_{i=1}^N m_{i,f} , \quad (2)$$

and the variance of feature f , $var(M_f)$, is obtained from

$$var(M_f) = \frac{1}{N-1} \sum_{i=1}^N (m_{i,f} - m_{.,f})^2. \quad (3)$$

As we will show in Sect. 5, depending on the combination of the initial groups of features, the number of columns of \mathbf{M} and \mathbf{T} may be more than 800. Thus, a strong feature reduction has to be made. So, At step 3b, a sample similarity matrix \mathbf{S} is calculated:

$$\mathbf{S} = \begin{bmatrix} Sim(s_1, s_1) & Sim(s_1, s_2) & \dots & Sim(s_1, s_N) \\ Sim(s_2, s_1) & Sim(s_2, s_2) & \dots & Sim(s_2, s_N) \\ \vdots & \vdots & \ddots & \vdots \\ Sim(s_N, s_1) & Sim(s_N, s_2) & \dots & Sim(s_N, s_N) \end{bmatrix} \quad (4)$$

Each cell of the symmetric matrix \mathbf{S} represents the similarity between two music samples and it is calculated by the following correlation from values of matrix \mathbf{T} :

$$Sim(s_i, s_j) = \frac{cov(s_i, s_j)}{\sqrt{cov(s_i, s_i)} \cdot \sqrt{cov(s_j, s_j)}} , \quad (5)$$

where the covariance between music samples s_i and s_j is given by

$$cov(s_i, s_j) = \frac{1}{F-1} \sum_{f=1}^F (t_{s_i,f} - t_{s_i,.})(t_{s_j,f} - t_{s_j,.}) , \quad (6)$$

where $t_{s_i,.}$ is the mean value of the i th line of \mathbf{T} .

Each line of matrix \mathbf{S} , corresponds to a music sample, now characterized by its similarity (within a range from -1 to +1) to all the other samples in the training set. On the other hand, each column of \mathbf{S} may be seen as a new feature reflecting the similarity between a music sample and all the other samples. Clearly, there are as many columns as the number of music samples in the training set. Thus, with \mathbf{S} , the number of features is reduced from the number of initial attributes, usually very high, to a number which is equal to the size of the training set, which may be a much smaller number. As we will show in Sect. 5 we obtained good results using a training set of 60 samples.

Since samples of the same genre tend to show high similarities and, thus, there are strong correlations between the features in \mathbf{S} , another reduction in dimensionality can be obtained by a technique based on Principal Component Analysis (PCA) [1].

Since \mathbf{S} is symmetric, it can be described as $\mathbf{S} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$, where $\mathbf{P} = [e_1, \dots, e_N]$ is the orthogonal matrix of normalized eigenvectors of \mathbf{S} , and $\mathbf{\Lambda}$ is the diagonal matrix of its eigenvalues, $\lambda_1, \dots, \lambda_N$, such that $\lambda_1 \geq \dots \geq \lambda_N \geq 0$. Since $\mathbf{\Lambda}$ is symmetric, $\mathbf{\Lambda} = \mathbf{\Lambda}^{\frac{1}{2}}\mathbf{\Lambda}^{\frac{1}{2}}$ and $\mathbf{\Lambda}^{\frac{1}{2}} = (\mathbf{\Lambda}^{\frac{1}{2}})^T$. Thus,

$$\mathbf{S} = \mathbf{P}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{\Lambda}^{\frac{1}{2}}\mathbf{P}^T = \mathbf{P}\mathbf{\Lambda}^{\frac{1}{2}}(\mathbf{\Lambda}^{\frac{1}{2}})^T\mathbf{P}^T = \mathbf{P}\mathbf{\Lambda}^{\frac{1}{2}}(\mathbf{P}\mathbf{\Lambda}^{\frac{1}{2}})^T = \mathbf{Q}\mathbf{Q}^T, \quad (7)$$

with

$$\mathbf{Q} = \mathbf{P}\mathbf{\Lambda}^{\frac{1}{2}}. \quad (8)$$

The lines in matrix \mathbf{Q} represent the music samples while the columns represent new uncorrelated features; see [1] for more details about this PCA-based technique. The leftmost columns of \mathbf{Q} correspond to the most informative features. Thus, in order to reduce the number of features, we can discard the least informative ones, by ignoring the columns of \mathbf{Q} having a variance (given by cells in $\mathbf{\Lambda}$) lower than a threshold which we set to 1. We call \mathbf{R} to this new reduced matrix (step 3c) of Fig. 1, which is a copy of the k leftmost columns of \mathbf{Q} . This way we build a k -dimensional space in which music samples are represented: the k features correspond to the k axis in this new space, and matrix \mathbf{R} contains the values for these features for each music in the training set. We tried other criteria associated with other threshold values, but this one provided a more reduced number of columns keeping good results. With this technique, we were able to drastically reduce the number of initial dimensions, that is, features (in \mathbf{S}) from 60 to 7 final dimensions (in \mathbf{R}) when we used the training set described in Sect. 5. Now, we are able to submit the resulting matrix \mathbf{R} to the clustering stage.

3.2 The Clustering Stage

In the clustering stage (box 4 in Fig. 1) we use the Model-Based Clustering Analysis (MBCA) as proposed by Fraley and Raftery [2]. This approach uses no initial information about the number of clusters nor their shape or orientation. It represents the data by several possible models, which are characterized by different geometric properties. With this approach, data is represented by a mixture model where each element corresponds to a different cluster. Models with varying geometric properties are obtained through different Gaussian parameterizations and cross-cluster constraints. Partitions (clusters) are determined by the EM (expectation-maximization) algorithm for maximum likelihood, with initial agglomerative hierarchical clustering (see [2] for details). This clustering methodology is based on multivariate normal mixtures. So, the density function associated to cluster c has the form:

$$f_c(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) = \frac{e^{(-\frac{1}{2}(\mathbf{x}_i - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_c))}}{(2\pi)^{\frac{p}{2}} |\boldsymbol{\Sigma}_c|^{\frac{1}{2}}}, \quad (9)$$

where vector \mathbf{x}_i represents an element that belongs to cluster c . Clusters are ellipsoidal and centered at the means $\boldsymbol{\mu}_c$. The covariance matrix $\boldsymbol{\Sigma}_c$ determines

the geometric characteristics of the cluster. This clustering methodology is based on the parameterization of the covariance matrix in terms of the eigenvalue decomposition in the form $\Sigma_{\mathbf{c}} = \lambda_{\mathbf{c}} \mathbf{D}_{\mathbf{c}} \mathbf{A}_{\mathbf{c}} \mathbf{D}_{\mathbf{c}}^T$, where $\mathbf{D}_{\mathbf{c}}$ is the orthogonal matrix of eigenvalues, which determines the orientation of the axes. $\mathbf{A}_{\mathbf{c}}$ is the diagonal matrix whose elements are proportional to the eigenvalues of $\Sigma_{\mathbf{c}}$ and which determines the shape of the ellipsoid. The volume of the ellipsoid is specified by scalar $\lambda_{\mathbf{c}}$. Characteristics (orientation, shape and volume) of distributions are estimated from the input data, and can be allowed to vary between clusters, or constrained to be the same for all clusters. Once all models are created, MBCA uses the Bayesian Information Criterion (BIC) to measure the evidence of clustering for each pair (*model, number of clusters*), and the larger the value of BIC, the stronger the evidence for the pair. So, by choosing the pair having the larger BIC, the most reliable model is automatically obtained, and then a vector indicating which music samples belong to which cluster is returned by this clustering approach. In other words, clusters are automatically formed in a k -dimensional space, according to data in matrix \mathbf{R} .

4 The Classification Method

Once the clusters are learned, the classification method can be used to classify new music samples (not included in the training set). Fig. 2 shows the steps of this method, which we describe below.

4.1 Representing New Music Samples in the k -dimensional Space Built in the Clustering Phase

Given a new (test) music sample s_t , the classification method starts by representing it with the same initial feature set as that used in the clustering method. As a result, the music sample is represented by an F -dimensional vector \mathbf{m}_{s_t} that contains the feature values for music s_t , that is $\mathbf{m}_{s_t}^T = [m_{s_t, f_1}, \dots, m_{s_t, f_F}]$ (step 2 in Fig. 2). Recall that matrix \mathbf{M} (from Sect. 3.1) is a matrix whose lines are vectors of this form for the music samples in the training set.

Afterwards, vector \mathbf{m}_{s_t} needs to be transformed into a new vector that represents music s_t in the k -dimensional space built in the clustering process. Firstly, \mathbf{m}_{s_t} needs to be standardized (box 3a in Fig. 2). This transformation aims to set equal importance (scale) to each feature in vector \mathbf{m}_{s_t} . Despite \mathbf{m}_{s_t} has only one value for each feature, this standardization will take into account the feature values of the music samples in the training set. Thus, the means and standard deviations calculated by the clustering method are used such that each cell of the new vector $\mathbf{t}_{s_t}^T = [t_{s_t, f_1}, \dots, t_{s_t, f_F}]$ is given by an equation similar to (1):

$$t_{s_t, f_i} = \frac{m_{s_t, f_i} - m_{., f_i}}{\sqrt{\text{var}(M_{f_i})}}, \quad (10)$$

where $m_{., f_i}$ and $\text{var}(M_{f_i})$ result from (2) and (3) respectively.

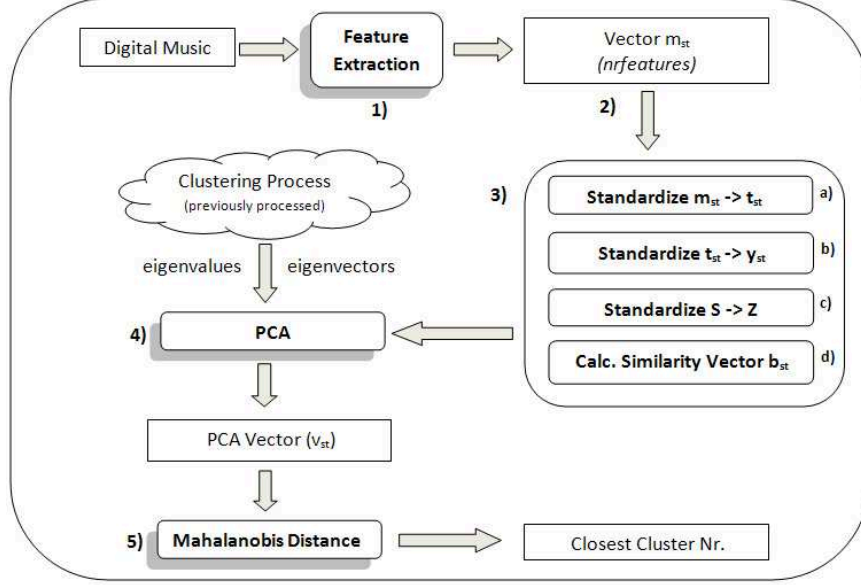


Fig. 2. The classification method

Now, the method aims to calculate a similarity vector between music s_t and all the samples in the training set. For that, we could use the correlation given by (5), used to calculate the similarity between the training set samples. However, since a correlation between non-standardized variables is equivalent to a covariance between the standardization of those variables, for reasons of computational weight, we followed this last option to get the same results. So, once we have vector \mathbf{t}_{s_t} , the referred standardization corresponds to a new vector $\mathbf{y}_{s_t}^T = [y_{s_t, f_1}, \dots, y_{s_t, f_F}]$ (box 3b in Fig. 2), where y_{s_t, f_i} is the standardized value of feature f_i for music s_t , which is given by

$$y_{s_t, f_i} = \frac{t_{s_t, f_i} - t_{s_t, \cdot}}{\sqrt{\text{var}(T_{s_t})}} . \quad (11)$$

$\text{var}(T_{s_t})$ stands for the variance associated to the t_{s_t, f_i} values for sound s_t along all features. So, $\text{var}(T_{s_t}) = \frac{1}{F-1} \sum_{f_i=1}^F (t_{s_t, f_i} - t_{s_t, \cdot})^2$ and $t_{s_t, \cdot} = \frac{1}{F} \sum_{f_i=1}^F t_{s_t, f_i}$.

At this step we need to relate our music s_t , now represented by \mathbf{y}_{s_t} , with the samples used in the learning process. In order to obtain a similarity vector \mathbf{b}_{s_t} (box 3d in Fig.2), we need the information given by the similarity matrix \mathbf{S} , which may also be given by another matrix \mathbf{Z} – see Appendix for details concerning matrices \mathbf{S} and \mathbf{Z} – such that each column of \mathbf{Z} is a vector $\mathbf{z}_s^T = [z_{s, f_1}, \dots, z_{s, f_F}]$ that represents the training set sample s using standardized values. In other words, each of these standardized values z_{s, f_i} is calculated

by

$$z_{s,f_i} = \frac{t_{s,f_i} - t_{s,\cdot}}{\sqrt{\text{var}(T_s)}} . \quad (12)$$

Thus, vector \mathbf{b}_{s_t} represents the similarity vector between \mathbf{y}_{s_t} and each sample of the training set:

$$\mathbf{b}_{s_t}^T = \frac{1}{F-1} \mathbf{y}_{s_t}^T \mathbf{Z} . \quad (13)$$

Now, by using the information obtained by the PCA-based technique from Sect. 3.1, that is, with \mathbf{A} and \mathbf{P} , we can transform \mathbf{b}_{s_t} into a vector \mathbf{u}_{s_t} , such that

$$\mathbf{u}_{s_t}^T = [u_{s_t,1}, \dots, u_{s_t,N}] = \mathbf{b}_{s_t}^T \mathbf{P} \mathbf{A}^{-\frac{1}{2}} , \quad (14)$$

where N is still the number of samples of the training set.

Similarly to what was mentioned in Sect. 3.1 about the most informative columns of matrix \mathbf{Q} , only the k leftmost cells of $\mathbf{u}_{s_t}^T$ are used to obtain a final vector \mathbf{v}_{s_t} that represents the music sample s_t in the k -dimensional space learned by the clustering method. In other words, $\mathbf{v}_{s_t} = [u_{s_t,1}, \dots, u_{s_t,k}]$. In Appendix, the reader may see a detailed proof that \mathbf{v}_{s_t} is the representation of music s_t in the k -dimensional space learned by the clustering method.

4.2 The Classification Stage

Now that music s_t is represented in the k -dimensional space learned by the clustering method, we need to relate \mathbf{v}_{s_t} to the learned clusters (box 5 in Fig. 2). Mahalanobis distance was adopted for this purpose since it takes into account the geometric properties of each cluster, which is important since distances take different impact depending on the data dispersion along each axis. (This characteristic is not achieved when using other metrics such as Euclidean or Manhattan distances.)

The method calculates the Mahalanobis distance between each cluster centroid and \mathbf{v}_{s_t} , and proposes the class represented by the cluster having a smaller distance as the most likely class for music s_t . In other words, class c will be associated to \mathbf{v}_{s_t} if $d(\mathbf{v}_{s_t}, \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c^{-1}) = \min_i d(\mathbf{v}_{s_t}, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i^{-1})$ where:

$$d(\mathbf{v}_{s_t}, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i^{-1}) = (\mathbf{v}_{s_t} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{v}_{s_t} - \boldsymbol{\mu}_i) , \quad (15)$$

where $\boldsymbol{\mu}_i = [\mu_{i,\cdot,1}, \dots, \mu_{i,\cdot,k}]$ is the centroid of cluster i , with $\mu_{i,\cdot,f} = \frac{1}{\|\mathcal{C}_i\|} \sum_{s \in \mathcal{C}_i} r_{s,f}$. \mathcal{C}_i is cluster i , that is, the set containing all the samples in this cluster, $\|\mathcal{C}_i\|$ is its size, and $r_{s,f}$ is the value of the f th axis (i.e., final feature) for music s (this is the value corresponding to the line associated to sample s and f th column of matrix \mathbf{R} , see Sect. 3.1). So, $\boldsymbol{\mu}_i$ represents an *average* music sample of cluster i . Finally, $\boldsymbol{\Sigma}_i$ reflects the geometric properties of cluster i in the k -dimensional space:

$$\boldsymbol{\Sigma}_i = \begin{bmatrix} E_{i1,1} & E_{i1,2} & \dots & E_{i1,k} \\ E_{i2,1} & E_{i2,2} & \dots & E_{i2,k} \\ \vdots & \vdots & \ddots & \vdots \\ E_{ik,1} & E_{ik,2} & \dots & E_{ik,k} \end{bmatrix} \quad (16)$$

and

$$E_{i,p} = \frac{1}{\|\mathcal{C}_i\| - 1} \sum_{s \in \mathcal{C}_i} (r_{s,l} - \mu_{i,l})(r_{s,p} - \mu_{i,p}) . \quad (17)$$

The most heavy calculations needed in the classification phase, such as the matrix Σ_i^{-1} for each cluster, can actually be made at the end of the clustering phase, as all needed data is available for that. This way, the classification of new music samples is a fast computation.

5 Results

In order to validate our approach, we used classical, metal, and reggae music samples from Tzanetakis' GTZAN ¹ data collection [14], which has 100 music samples from several different genres. In addition, we added samples from a new genre, *fado*, and therefore, our data collection has 400 music samples (all represented with a sampling frequency of 22050 Hz, 16 bits, and single channel) representing 4 different music genres: classical, fado, metal, and reggae.

Table 1. Clustering error percentages for several feature combinations. The features are: timbral texture features (ttf), rhythm patterns (rssi_rh), beat, root mean square of the spectrogram frames (rmsFrame), MFCCs, spectral centroid + bandwidth + uniformity (centBandUnif), SSD over spectrogram (specStat), low-energy over sample spectrum (lener), spectral centroid (scentroid), and ZCR

	Feature Combination	Error (%)
1	'tff', 'rssi_rh', 'beat', 'rmsFrame', 'mfccs', 'centBandUnif'	0
2	'tff', 'rssi_rh', 'beat', 'rmsFrame', 'mfccs', 'specStat'	0
3	'tff', 'rssi_rh', 'centBandUnif'	0
4	'tff', 'rssi_rh', 'specStat'	2
5	'tff', 'rssi_rh', 'beat', 'rmsFrame', 'mfccs', 'lener', 'scentroid', 'specStat'	3
6	'tff', 'rssi_rh', 'beat', 'rmsFrame', 'mfccs', 'lener', 'scentroid', 'zcr', 'specStat'	3

Even though our clustering methodology does not use any information about the number of genres nor the genre of the samples, we used this labelling information to validate the results. Thus, once the clustering process is complete, we assume that each learned cluster c corresponds to the mostly represented genre in the cluster, and count the number of samples, o_c , in the cluster that have a different labeling. The overall error percentage is given by $e = (100 \sum_c o_c) / N$, where N is the number of samples.

¹ <http://marsyas.info/download/>

Table 2. Accuracy of the classification results for 3 different feature combinations

Feature Combination	Accuracy rate (%)
1	76.5
2	81.8
3	73.8

In order to validate the clustering methodology (described in Sect. 3), we used a training set composed of 60 elements (15 music samples from each of the four music genres mentioned above) and we tested many combinations of features (from Sect. 2). Table 1 shows the clustering results for the best feature combinations. The first three combinations have a 0% error rate, which shows that this approach is able to achieve perfect clustering results (assuming the initial labelling is correct).

Based on Table 1, it is clear that the third combination uses less features than the top two combinations. On the other hand, if we look at the clusters created (Fig. 3), the second combination achieves clustering results that perfectly match the initial labelling of the data. Nonetheless, this does not mean that combinations 1 and 3 achieve worse or incorrect results. It may actually be the case that these two combinations are learning sub-genres within classical, fado and metal. Each feature in the second combination actually represents a group of features as the whole number of real audio sub-features this combination represents is equal to 873.

Once the clusters were learned, we proceeded and classified new music samples. In order to evaluate the performance of the classification method (described in Sect. 4), we used a test set with the remaining 85 music samples (not used for clustering) for each of the four genres, making a total of 340 samples. Tests were made for clusters learned from each of the first three combinations from Table 1. As can be seen in Table 2, combination 2 achieves the best accuracy (precision) results with 81.8% correctly classified samples, which is a very satisfactory result, given this is an unsupervised approach.

6 Conclusions and Future Work

We proposed an unsupervised clustering and classification methodology for automatic genre classification. This kind of approach has the advantage of being totally independent of any influence from a human taxonomy. Since music genres do not present clear boundaries between them, and human genre taxonomy is hard to be achieved, we believe that an unsupervised approach is more suitable for music genre classification, while a supervised approach based on previously labelled data tends to be subjective. Besides, by learning directly from data in features, an unsupervised approach may automatically detect new genres, which is not possible for the more static nature of the supervised approaches.

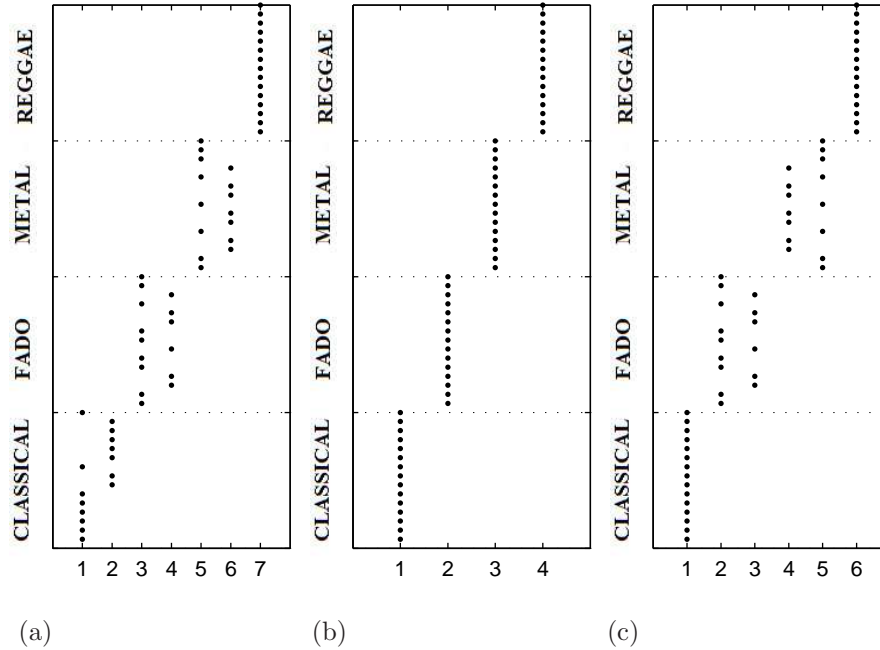


Fig. 3. Clustering results for combinations (a) 1, (b) 2 and (c) 3 from table 1. The x -axis shows the clusters learned, while the y -axis shows the initial labeling of the data. For instance, figure (a) shows that 8 classical music samples fell in cluster number 1 while 7 classical samples fell in cluster number 2

In order to learn the clusters, the clustering method uses only the audio features of the training samples, and no previous knowledge on the genre of the individual samples. In addition, no information on the number of clusters is given a priori. This method achieves a perfect clustering accuracy (100%) when tested with four music genres (even though the genre labeling was not used in the training process, the results agreed with the manual labelling of the data), which shows that it is possible to achieve good accuracy results using an unsupervised method. In addition, as discussed in section 5, depending on the audio features used, the method is also able to find sub-genres within the data.

Once the clusters are learned, the classification method can categorize new music samples according to the previously learned clusters. This method uses Mahalanobis distance so that it is able to deal with clusters of different shapes, volumes and orientations. An 82% classification rate was obtained with four music genres.

We noticed that some misclassified samples were almost equidistant (in terms of Mahalanobis distance) to the chosen cluster and their actual cluster. This suggests that, as future work, a further analysis must be done to detect the existence of possible patterns of the Mahalanobis distances to every cluster,

for both cases (correct classifications and incorrect ones). Such analysis may be important, for instance, to decide if a sample that is too distant from all clusters must be rejected; or to suggest that new clusters ought to be learned (which should be done by running the clustering method again) because several samples are approximately equidistant to two given clusters.

Even though we only reported the clustering and classification results for four music genres, we are currently investigating how the system behaves with more genres. Although this work is still not finished, we were already able to confirm that this clustering method can achieve good results with five and six music genres (at least around 90% clustering accuracy). Working with more music genres may require the use of more (or different) audio features. There are other audio features that we did not explore yet but could be important to discriminate other genres.

Finally, in order to test different feature combinations we simply used a brute force method, that is, with no prior selection. Instead, a possible *filtering* over the extracted features should also be explored in future work, as to only process those features that present higher variances between music samples.

Appendix

Here we prove that the test samples are represented in the k -dimensional space learned by the clustering method, that is, that \mathbf{v}_{s_t} (see Sect. 4.1) is the translation of test music sample s_t in this k -dimensional space.

Proof. Let us suppose we want to classify a sound, say the first music of the training set, which is available in \mathbf{z}_1 , the first column of matrix \mathbf{Z} . So, by (13) $\mathbf{b}_1^T = \frac{1}{F-1} \mathbf{z}_1^T \mathbf{Z}$ since now \mathbf{y}_{s_t} is substituted by \mathbf{z}_1 ; $\mathbf{z}_1^T = [z_{1,1}, \dots, z_{1,F}]$. Then $\mathbf{b}_1^T = [b_{1,1}, \dots, b_{1,N}]$ where

$$b_{1,j} = \frac{1}{F-1} \sum_{i=1}^F z_{1,i} \cdot z_{j,i} . \quad (18)$$

Notice that, by statistics theory, (18) and (5) give the same result since $Sim(s_i, s_j)$ in (5) is a correlation using non-standardized values, and $b_{i,j}$ in (18) (generalizing from 1 to i) is a covariance using the standardization of those values. Then, in order to simplify this proof, let us suppose that we want to classify not just one music from the training set, but the whole training set. Then it is easy to conclude that

$$\mathbf{B} = \frac{1}{F-1} \mathbf{Z}^T \mathbf{Z} . \quad (19)$$

\mathbf{B} would be obtained instead of \mathbf{b}_1 . Note that $\mathbf{B} = \mathbf{S}$, being \mathbf{S} the similarity matrix given by (4), because it contains the similarity vectors between each training set music and all other music samples.

Now let us work with the entire \mathbf{S} as if we wanted to translate all training sounds into vectors in the k -dimensional space. Then, from (14) we would obtain

$\mathbf{G} = \mathbf{S}\mathbf{P}\mathbf{\Lambda}^{-\frac{1}{2}} = \mathbf{S}\mathbf{P}\mathbf{\Lambda}^{-1}\mathbf{\Lambda}^{\frac{1}{2}}$, but since $\mathbf{P}^T\mathbf{P} = \mathbf{I}$ (where \mathbf{I} is the identity matrix) and $\mathbf{S} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T$, then $\mathbf{G} = \mathbf{P}\mathbf{\Lambda}\mathbf{P}^T\mathbf{P}\mathbf{\Lambda}^{-1}\mathbf{\Lambda}^{\frac{1}{2}} = \mathbf{P}\mathbf{\Lambda}\mathbf{\Lambda}^{-1}\mathbf{\Lambda}^{\frac{1}{2}} = \mathbf{P}\mathbf{\Lambda}^{\frac{1}{2}}$.

But $\mathbf{P}\mathbf{\Lambda}^{\frac{1}{2}} = \mathbf{Q}$, the matrix characterizing all sounds by the PCA-based method presented before (see (8)). Since we used a *copy* of the whole training set for classification instead of just one music, we obtained a matrix (\mathbf{Q}) instead of a vector \mathbf{u}_1 . Then, choosing the k leftmost columns of this matrix we would obtain matrix \mathbf{R} referred in Sect. 3.1, which contains the representation of the whole training set in the k -dimensional space. With this, we proved that \mathbf{v}_{s_t} is the representation of the test music sample in the k -dimensional space learned by the clustering method.

Acknowledgments

This work was part of the Videoflow project and partially funded by *Quadro de Referência Estratégica Nacional* (QREN) and *Fundo Europeu para o Desenvolvimento Regional* and *Programa POR Lisboa*.

Bibliography

- [1] Y. Escoufier and H. L'Hermier. A propos de la comparaison graphique des matrices de variance. *Biometrical Journal*, 20(5):477–483, 1978.
- [2] C. Fraley and A. E. Raftery. How many clusters? which clustering method? answers via Model-Based cluster analysis. *The Computer Journal*, 41(8):578–588, 1998.
- [3] S. Golub. Classifying recorded music. Master's thesis, University of Edinburgh - Division of Informatics, 2000.
- [4] A.L. Koerich and C. Poitevin. Combination of homogeneous classifiers for musical genre classification. In *Systems, Man and Cybernetics, 2005 IEEE International Conference on*, volume 1, pages 554–559, 2005.
- [5] C. Lee, J. Shih, K. Yu, and H. Lin. Automatic music genre classification based on modulation spectral analysis of spectral and cepstral features. *IEEE Transactions on Multimedia*, 11(4):670–682, 2009.
- [6] T. Li and M. Ogihara. Toward intelligent music information retrieval. *Multimedia, IEEE Transactions on*, 8(3):564–574, 2006.
- [7] T. Lidy and A. Rauber. Evaluation of feature extractors and psycho-acoustic transformations for music genre classification. In *ISMIR*, pages 34–41, 2005.
- [8] M.F. Mckinney, J. Breebaart, and Prof Holstlaan. Features for audio and music classification. In *ISMIR*, 2003.
- [9] D. Pye. Content-based methods for the management of digital music. In *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 6, pages 2437–2440 vol.4, 2000.
- [10] A. Rauber, E. Pampalk, and D. Merkl. Using Psycho-Acoustic models and Self-Organizing maps to create hierarchical structuring of music by sound similarity. In *ISMIR*, 2002.
- [11] X. Shao, C. Xu, and M.S. Kankanhalli. Unsupervised classification of music genre using hidden markov model. In *IEEE International Conference on Multimedia and Expo (ICME)*, volume 3, pages 2023–2026, 2004.
- [12] H. Soltan, T. Schultz, M. Westphal, and A. Waibel. Recognition of music types. In *Acoustics, Speech and Signal Processing, 1998. Proceedings of the 1998 IEEE International Conference on*, volume 2, pages 1137–1140, 1998.
- [13] T. Tolonen and M. Karjalainen. A computationally efficient multipitch analysis model. *IEEE transactions on speech and audio processing*, 8(6):708–716, 2000.
- [14] G. Tzanetakis and P. Cook. Musical genre classification of audio signals. *Speech and Audio Processing, IEEE Transactions on*, 10(5):293–302, 2002.
- [15] G. Tzanetakis, A. Ermolinskyi, and P. Cook. Pitch histograms in audio and symbolic music information retrieval. *Proceedings of the third international conference on music information retrieval (ISMIR)*, pages 31–38, 2002.