

# MUSIC MOOD CLASSIFICATION BY RHYTHM AND BASS-LINE UNIT PATTERN ANALYSIS

*Emiru Tsunoo, Taichi Akase, Nobutaka Ono and Shigeki Sagayama*

Graduate School of Information Science and Technology, The University of Tokyo  
7-3-1 Hongo, Bunkyo-ku, Tokyo, 113-8656, Japan  
{tsunoo,t-akase,onono,sagayama}@hil.t.u-tokyo.ac.jp

## ABSTRACT

This paper discusses an approach for the feature extraction for audio mood classification which is an important and tough problem in the field of music information retrieval (MIR). In this task the timbral information has been widely used, however many musical moods are characterized not only by timbral information but also by musical scale and temporal features such as rhythm patterns and bass-line patterns. In particular, modern music pieces mostly have certain fixed rhythm and bass-line patterns, and these patterns can characterize the impression of songs. We have proposed the extraction of rhythm and bass-line patterns, and these unit pattern analysis are combined with statistical feature extraction for mood classification. Experimental results show that the automatically calculated unit pattern information can be used to effectively classify musical mood.

**Index Terms**— Audio classification, pattern clustering method, feature extraction, dynamic programming,  $k$ -means clustering.

## 1. INTRODUCTION

Due to the increasing size of music collections available on computers and portable music players, researches related to music information retrieval (MIR) including automatic genre classification and mood classification from audio has surged recently. Mood classification systems share the same formulation with genre classification systems which have been provided a established way of evaluating new representations of musical content and generally constructed of two stage processes: feature extraction and classification. These audio classification tasks have been competed in the contest called MIREX [1] for MIR researches for a long time. Mood classification is one of a variation of the audio classification however has even more elusive ground truth [2, 3]. In mood classification, not only instrumental information but also musical scale, rhythm and bass-line informations are thought to be important. In many case especially of modern popular music, some fixed bar-long percussive patterns and bass-line patterns are repeated in whole a song and the unit patterns frequently characterize the music. If such representative bar-long percussive patterns and bass-line patterns in music can be captured automatically as templates, they can potentially be used to characterize different music mood directly from audio signals.

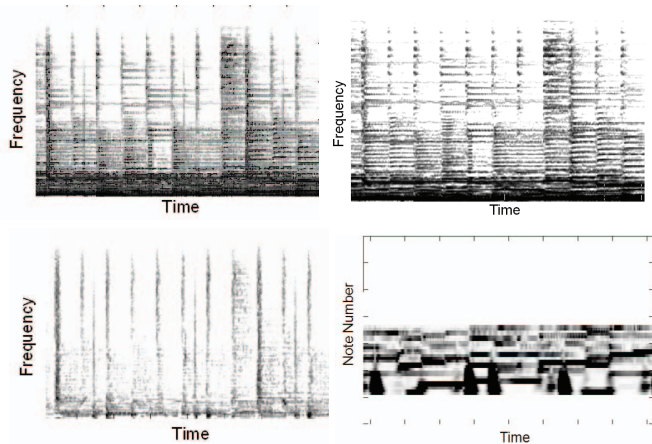
In previous research, timbral features, rhythmic features and pitch features have been used for audio genre classification [4], and this sort of feature extraction is widely used in audio classification including mood classification in the past MIREX contests. In this work, the timbral features were the

most dominant and the other statistical features came out not to be so useful for audio classification. This feature extraction method was tested for mood classification [2] and the effectivity was verified. The statistical features are also classified in various ways such as active learning of support vector machine (SVM) [5]. In comparison to these statistical feature extraction, we have proposed temporal feature extractions including bar-long percussive pattern information [6] and bar-long bass-line pattern information [7], and these two features are tested to the genre classification task.

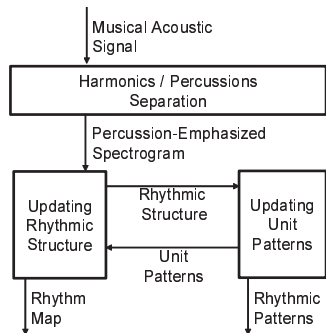
In this paper, we discuss an approach for extracting unit rhythm and bass-line patterns from a number of audio tracks and propose a feature vector for the application to mood classification. First, we separate percussive sound and harmonic sound of the audio tracks. Then we propose a clustering method specialized to rhythm patterns using a combination of dynamic programming and  $k$ -means clustering, and to bass-line patterns based on the  $k$ -means clustering algorithm. For the purpose of an application to audio mood classification, the scheme to extract feature vector based on clustered patterns which contain temporal information is suggested. Finally, the effectiveness of the proposed feature extraction for mood classification is verified experimentally.

## 2. HARMONIC/PERCUSSIVE SOUND SEPARATION

Generally, harmonic and percussive sounds are mixed in the observed spectrograms of audio pieces. Therefore in order to perform rhythm and bass-line pattern analysis it is useful to separate these components as a preprocessing, and percussive components are used for rhythm analysis and harmonic components for bass-line analysis. We utilize the harmonic/percussive sound separation (HPSS) technique proposed by Ono [8] that is based on the difference of general timbral features. By looking at the upper left figure in Fig. 1, a typical instance of spectrogram, one can observe that harmonic components tend to be continuous along the temporal axis in particular frequencies. On the other hand, percussive components tend to be continuous along the frequency axis and temporally short. Mask functions for separating the two components (harmonic and percussive) are calculated following a maximum a priori (MAP) estimation approach using the expectation maximization (EM) algorithm. Applying this approach to the shown spectrogram, harmonic and percussive components are separated (harmonic and percussive components are shown in the upper right and the lower left of Fig. 1 respectively).



**Fig. 1.** The original spectrogram (upper left), the harmonics-emphasized spectrogram (upper right) and the percussion-emphasized spectrogram (lower left) of a popular music piece (RWC-MDB-G-2001 No.6[9]). The low-pass filtered logarithmic spectrogram calculated using wavelet transform from harmonics-emphasized spectrogram is shown in lower right figure.



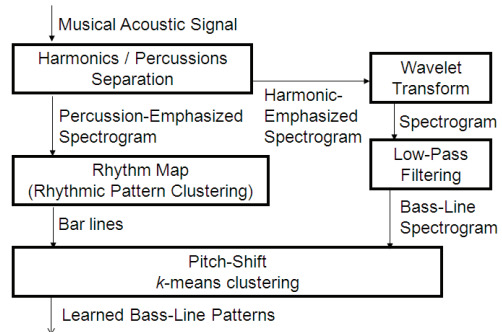
**Fig. 2.** The flow diagram of the extraction of rhythm patterns.

### 3. BAR-LONG UNIT PATTERN CLUSTERING

#### 3.1. Rhythm Pattern Clustering

Bar-long percussive patterns are frequently common and characteristic of a particular mood or style. Automatically detecting these patterns is a “chicken-and-egg” problem in that sets of bar-long unit rhythm patterns may be determined only after their corresponding unit boundaries in the music pieces are given, and vice versa. This is complicated by tempo fluctuations which might cause the unit pattern to stretch or shrink.

In order to solve these problems, the rhythm map which we have proposed [10, 6] is employed. The rhythm map is an approach to estimate representative bar-long percussive patterns and its segmentation (i.e. bar lines) simultaneously. This algorithm is processed to the percussive-emphasized spectrogram. By iterating dynamic programming (DP) matching and updating the templates used for DP matching based on the  $k$ -means-clustering-like update rules, both segmentation and templates themselves are updated. After convergence, the multiple percussive patterns in the input song are learned and the optimal segmentation is obtained. We use this estimated segmentation as bar lines. Fig. 2 illustrates the flow of this algorithm.



**Fig. 3.** The flow diagram of the bass-line pattern extraction.

#### 3.2. Bass-line Pattern Clustering

Bar-long bass-line patterns are also frequently common and characteristic of a particular mood. Bar lines have been already estimated in the process shown in section 3.1. Therefore, only problem to be solved is that unit bass-line patterns are shifted in pitch according to the chord played. For example, a pattern consists of only root notes in a uniform rhythm, all notes in this pattern need to be pitch-shifted by the same amount of notes according to the chord, because the root note changes accompany with the chord changes. Fig. 3 illustrates the flow of the bass-line pattern clustering algorithm.

If there was no pitch shift a simple  $k$ -means clustering approach can be used to estimate the representative bar-long bass-line patterns: Distances between each bar-long spectrogram pattern and centroid spectrogram patterns are calculated, and the centroids are updated by averaging the sample patterns. In order to deal with pitch-shift problem, we have proposed an approach where every possible pitch-shift is compared in  $k$ -means framework [7]. The lower right of Fig. 1 shows the logarithmic spectrogram which is processed low-pass filtering after Gabor wavelet transform whose frequency resolution is semitone (100 cents). The band-pass filtering was done by setting high and low frequency components to be zero. This kind of spectrogram is used for this purpose.

### 4. FEATURE EXTRACTION

#### 4.1. Mood Classification via Pattern Clustering

Ideally percussive patterns for a particular mood would be fixed and would be automatically extracted perfectly. If that was the case then automatic mood classification could be performed simply by looking at which particular rhythm and bass-line patterns are used in a music piece. However in practice there is no guarantee that patterns are fixed for a particular mood or that their automatic extraction will be perfect. Therefore in many cases the percussive and bass-line patterns of a particular music piece will belong to more than one mood. To address these problems simultaneously we utilize a pattern occurrence histogram and distance related representation followed statistical machine learning to automatically classify music mood. Supervised learning classifiers such as Support Vector Machines (SVM)[11] can be used for this purpose.

#### 4.2. Percussive Pattern Feature Extraction

One possible way to extract percussive feature vector is count up which percussive pattern templates are contained in a song

and calculating the mood pattern occurrence histogram, similarly to Latent Semantic Indexing approach [12].

If  $K$  pattern templates are learned from mood  $m$  ( $m = 1, \dots, M$ ), an alignment can be calculated using dynamic programming to get the templates  $T_{k,m}$  that exist in the song  $s$ . Then, the occurrence number of the patterns from mood  $m$  can be simply calculated by summation as follows:

$$c_{s,m} = \sum_{k=1}^K c_{s,k,m} \quad (1)$$

where  $c_{s,k,m}$  is the number of the template  $T_{k,m}$  in the song  $s$ , and the eventual  $M$  dimensional pattern occurrence histogram features vector  $\vec{x}$  of song  $s$  can be written as

$$x_g = \frac{c_{s,m}}{N_s} \quad (2)$$

which is normalized by  $N_s$ , the number of measure in the song  $s$ .

### 4.3. Bass-line Pattern Feature Extraction

For bass-line pattern, one way to extract feature vector is calculating distances between every measure of input spectrogram and every template pattern and averaging them through whole an input piece. Even though there is a possibility for a music piece to belong to more than one mood templates, the distances between spectrogram in the input piece and learned templates are still affected.

The mathematical definition of the feature vector is following. After bass-line pattern templates are learned, we have  $K \cdot M$  templates when  $K$  templates are learned from  $M$  moods. Then the distances between input song which have  $N$  measures and learned templates are calculated. The averaged distances are obtained as follows:

$$d_l = \frac{1}{N} \sum_{n=1}^N D(X_n, B_l) \quad (3)$$

where  $1 \leq l \leq KM$  is the template number and  $D(\cdot, \cdot)$  is a distance with pitch-shift problem proposed in [7]. The feature vector  $\mathbf{x}$  can be written as

$$\mathbf{x} = (\mathbf{d}_1, \mathbf{d}_2, \dots, \mathbf{d}_{KM})^T. \quad (4)$$

### 4.4. Statistical Feature Extraction

#### 4.4.1. Mel-Frequency Cepstral Coefficients

MFCCs are one of the most compact and efficient spectral expression that contain the general frequency characteristics important to human hearing. Originally MFCCs are developed for automatic speech recognition however they have been found to be powerful also for other auditory domains such as MIR [13].

In order to calculate MFCC feature, the log-magnitude spectrogram is wrapped to the Mel frequency scale and the inverse discrete cosine transform is performed. The first 13 coefficients are calculated in each time window of short time Fourier transform and the mean and the standard deviation for each song are calculated on the whole piece.

#### 4.4.2. Musical Scale Feature

Musical scale also characterizes musical mood. For instance, when a note “do” was sounded the note “mi” makes the harmony bright, however “mi flat” makes it dark. In order to extract this kind of information, the correlation of chroma vector is calculated.

**Table 1.** Mood category definition in MIREX and group of annotation in CAL500 dataset used for labeling.

Category	Definition	Annotation
Passionate	Passionate	Emotional/passionate
	Rousing	Happy
	Confident	Arousing/awakening
	Boisterous	Powerful/strong
	Rowdy	Exciting/thrilling Positive/optimistic
Fun	Fun	Happy
	Rollicking	Pleasant/comfortable
	Cheerful	Cheerful/festive
	Sweet	Carefree/lighthearted
	Amiable/good natured	Exciting/thrilling
Sad	Literate	Sad
	Poignant	Tender/soft
	Wistful	Loving/romantic
	Bittersweet	Touching/loving
	Autumnal	
	Brooding	
Humorous	Humorous	Cheerful/Festive
	Silly	Carefree/lighthearted
	Campy	Light/playful
	Quirky	Bizarre/weird
	Whimsical	
	Witty Wry	
Aggressive	Aggressive	Aggressive
	Fiery	
	Tense/anxious	
	Intense	
	Volatile Visceral	

Chroma vector is a 12 dimensional vector whose columns represent overlapped energies of 12 semitones over octaves. Let the  $i$ th element of chroma vector at the time  $t$  be  $c(t, i)$ , a normalized correlation of the vector can be written as

$$S(t, \tau) = \sum_{i=1}^{12} \frac{\left( \frac{e(t)}{12} - c(t, i) \right) \left( \frac{e(t)}{12} - c(t, i + \tau) \right)}{e(t)^2} \quad (5)$$

where  $e(t) = \sum_{i=1}^{12} c(t, i)$ . It means just an energy in the case  $\tau = 0$  and this correlation is symmetrical. Therefore, the coefficients which contain useful information are only 6 of them, and the mean and the standard deviation for each song are calculated on the whole piece.

## 5. EXPERIMENTAL RESULTS

### 5.1. Dataset

Experiments with the proposed algorithm were conducted on the CAL500 dataset [14]. The dataset had annotation from 1 to 5 for each emotion category. The categorization was done based on MIREX contest where 5 categories are defined as Table 1. In order to label moods following this definition, the total score of annotation in the group shown in Table 1 is calculated per song and the song was labeled as the most dominant group i.e. mood category. However because there were only 31 songs that were labeled as “Humorous,” we eliminated this category and arranged 60 songs per each of the rest 4 categories randomly. This number was thought to be very small for classification using machine learning techniques, therefore, each song was split into 3 ten seconds of sequences. In this way, we arranged 180 songs hypothetically.

**Table 2.** Mood classification accuracy results.

Features	Data 1	Data 2
Baseline (Tzanetakis, <i>et. al.</i> )	44.2%	45.8%
Proposed (only statistics)	50.8%	51.7%
Proposed (all features)	56.4%	53.3%

## 5.2. Template Learning and Feature Extraction

First, common percussive and bass-line pattern templates were learned using the proposed algorithm for each mood. The proposed algorithm was implemented using the audio processing framework, *Marsyas*<sup>1</sup> which is open source software with specific emphasis on MIR [15].

To ensure that the template learning didn't become accidentally good for classification, we divided each dataset 50-50 into two parts randomly avoiding the album effect and obtained two sets of templates for each mood. The album effect is an unduly good performance of classification happened when songs from the same album are used for both training and testing. In this experiment, 10 percussive templates and 20 bass-line templates were learned from each particular mood, and the number of iterations was fixed to 30 times for percussion patterns and 15 times for bass-line patterns because it was enough to converge. The proposed features are extracted using templates learned from the other half data, and the other way around on the other half.

## 5.3. Classification Results

To train a classifier in feature space, the "Weka" machine learning toolkit [16] was employed. All the results shown are based on 10-fold cross-validation using a linear SVM as a classifier. To generalize the feature extraction part, we divided the datasets into two parts, though, in 10-fold cross-validation to generalize the classification, the labeled data was split into 10 folds and each fold is used once for testing with the remaining 9 folds used for training the classifier to generalize the classification.

An existing state-of-the-art mood classification system which uses 68 dimensional timbral features such as MFCC and spectrum centroids proposed by Tzanetakis [4] was used for comparison. The system performed well on several audio classification tasks in MIREX 2008 [1]. This system was used in past mood classification experiment in [2] and a part of the features (MFCCs) was used in [5]. This system achieved 44.2% and 45.8% on each half data.

In comparison to this, proposed features was total 122 dimensional vector and performed 56.4% and 53.3% on each half data, improving about 10% of accuracy from the state-of-the-art system. Even only the proposed statistic features exceeded the past accuracy and achieved 50.8% and 51.7%. It means the musical scale features also have a power to capture the characteristics of a mood, and the rhythm and bass-line features improves the classification accuracy over this statistics. Hence, the effectiveness of the proposed features is verified. The results are listed in Table 2.

## 6. CONCLUSIONS

We discussed an approach for clustering bar-long common percussive patterns and bass-line patterns for particular mood and extracting feature vectors for mood classification. We

<sup>1</sup><http://marsyas.sness.net/>

used HPSS technique to separate percussive components and harmonic components from audio signals as a preprocessing. Percussive patterns were clustered using a combination of one-pass DP and  $k$ -means clustering algorithm. Bass-line patterns were clustered using a new clustering method based on  $k$ -means clustering with pitch-shift. For audio mood classification, new feature vectors were defined as pattern occurrence histograms for percussive patterns and as averaged distances from each template for bass-line patterns. In combination with statistical features including MFCCs and musical scale feature, the effectivity of the features was verified experimentally.

Future work includes using  $n$ -gram model approach rather than only looking at the uni-gram histogram for rhythm pattern features. Additionally, other features than pattern distance vector of bass-line pattern can be devised. Combination with other features like chord transition information can be done to improve further audio classification as well.

## 7. ACKNOWLEDGMENT

This research was partly supported by CrestMuse Project under JST.

## 8. REFERENCES

- [1] "Mirex 2008," <http://www.music-ir.org/mirex/2008/>.
- [2] T. Li and M. Ogihara, "Detecting emotion in music," in *Proc. of ISMIR*, 2003, pp. 239 – 240.
- [3] X. Hu, M. Bay, and SJ Downie, "Creating a simplified music mood classification ground-truth set," in *Proc. of ISMIR*, 2007, pp. 309–310.
- [4] G. Tzanetakis and P. Cook, "Musical genre classification of audio signals," *IEEE Transaction on Speech and Audio Processing*, vol. 10, no. 5, pp. 293–302, 2002.
- [5] M. Mandel, G. Poliner, and D. Ellis, "Support vector machine active learning for music retrieval," *ACM Multimedia Systems*, vol. 12, no. 1, pp. 3 – 13, 2006.
- [6] E. Tsunoo, G. Tzanetakis, N. Ono, and S. Sagayama, "Audio genre classification using percussive pattern clustering combined with timbral features," in *Proc. of ICME*, 2009, pp. 382 – 385.
- [7] E. Tsunoo, N. Ono, and S. Sagayama, "Musical bass-line pattern clustering and its application to audio genre classification," in *Proc. of ISMIR*, 2009.
- [8] N. Ono, K. Miyamoto, H. Kameoka, and S. Sagayama, "A real-time equalizer of harmonic and percussive components in music signals," in *Proc. of the 9th Int. Conf. on Music Information Retrieval*, September 2008, pp. 139–144.
- [9] M. Goto, H. Hashiguchi, T. Nishimura, and R. Oka, "Rwc music database: Music genre database and musical instrument sound database," in *Proc. of the 4th Int. Conf. on Music Information Retrieval*, October 2003, pp. 229–230.
- [10] E. Tsunoo, N. Ono, and S. Sagayama, "Rhythm map: Extraction of unit rhythmic patterns and analysis of rhythmic structure from music acoustic signals," in *Proc. of ICASSP*, 2009, pp. 185–188.
- [11] V. Vapnik, *The Nature of Statistical Learning Theory*, Springer-Verlag, 1995.
- [12] S. Deerwester, S. Dumais, G. Furnas, T. Landauer, and R. Harshman, "Indexing by latent semantic analysis," *Journal of the American Society for Information Science*, vol. 41, no. 6, pp. 391–407, 1990.
- [13] B. Logan, "Mel frequency cepstral coefficients for music modeling," in *Proc. of ISMIR*, 2000.
- [14] D. Turnbull, L. Barrington, D. Torres, and G. Lanckriet, "Towards musical query-by-semantic-description using the cal500 data set," in *Proc. of ACM SIGIR*, 2007, pp. 439 – 446.
- [15] G. Tzanetakis, *Marsyas-0.2: A Case Study in Implementing Music Information Retrieval System*, chapter 2, pp. 31 – 49, Idea Group Reference, 2007, Shen, Shepherd, Cui, Liu (eds).
- [16] I. Witten and E. Frank, *Data Mining: Practical Machine Learning Tools and Techniques*, Morgan Kaufmann, 2005.