

## On Automatic Music Genre Recognition by Sparse Representation Classification using Auditory Temporal Modulations

Sturm, Bob L.; Noorzad, Pardis

*Published in:*  
Proceedings of the 9th International Symposium on Computer Music Modeling and Retrieval

*Publication date:*  
2012

*Document Version*  
Early version, also known as pre-print

[Link to publication from Aalborg University](#)

*Citation for published version (APA):*  
Sturm, B. L., & Noorzad, P. (2012). On Automatic Music Genre Recognition by Sparse Representation Classification using Auditory Temporal Modulations. In *Proceedings of the 9th International Symposium on Computer Music Modeling and Retrieval* (pp. 379-394). <http://cmmr2012.eecs.qmul.ac.uk/programme>

### General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal -

### Take down policy

If you believe that this document breaches copyright please contact us at [vbn@aub.aau.dk](mailto:vbn@aub.aau.dk) providing details, and we will remove access to the work immediately and investigate your claim.

# On Automatic Music Genre Recognition by Sparse Representation Classification using Auditory Temporal Modulations

Bob L. Sturm<sup>1</sup> and Pardis Noorzad<sup>2</sup>

<sup>1</sup> Department of Architecture, Design and Media Technology  
Aalborg University Copenhagen

Lautrupvang 15, 2750 Ballerup, Denmark

<sup>2</sup> Department of Computer Engineering and IT  
Amirkabir University of Technology  
424 Hafez Ave., Tehran, Iran

bst@create.aau.dk, pardis@aut.ac.ir

**Abstract.** A recent system combining sparse representation classification (SRC) and a perceptually-based acoustic feature (ATM) [31, 30, 29], outperforms by a significant margin the state of the art in music genre recognition, e.g., [3]. With genre so difficult to define, and seemingly based on factors more broad than acoustics, this remarkable result motivates investigation into, among other things, why it works and what it means for how humans organize music. In this paper, we review the application of SRC and ATM to recognizing genre, and attempt to reproduce the results of [31]. First, we find that classification results are consistent for features extracted from different analyses. Second, we find that SRC accuracy improves when we pose the sparse representation problem with inequality constraints. Finally, we find that only when we reduce the number of classes by half do we see the high accuracies reported in [31].

## 1 Introduction

Simply because we lack clearly definitive examples, and any utilitarian definitions, the automatic recognition of music genre is different from other tasks in music information retrieval. The human categorization of music seems natural, yet appears fluid and often arbitrary by the way it appears motivated by more than measurable characteristics of audible changes in pressure [9, 26, 17]. Extra-musical information, such as artist fashion, rivalries and the fan-base, associated dance styles, lyrical subjects, societal and political factors, religious beliefs, and origins in time and location, can position a particular piece of music into one category or another, not often without debate [38]. With the changing fashions of communities and the needs of companies, new genres are born [17]. And genres become irrelevant and lost, though we might still hear the recorded music and classify it as something entirely different.

It seems daunting then to make a computer recognize genre with any success. Yet, in developments between 2002 and 2006 [23], we have seen the accuracy of such algorithms progress from about 60% — using parametric models created from bags of features [42] — to above 80% — aggregating features over long time scales and employing weak classifiers [3]. The majority of approaches developed so far use features derived only from the waveform, and/or its symbolic form.

Some work has also explored mining user tags [27] and written reviews [1], or analyzing song lyrics [22]. Since humans have been measured to have accuracies around 70% after listening to 3 seconds of music — which surprisingly drops only down to about 60% for only half a second of listening [13] — the results of the past decade show that the human categorization of music appears grounded to a large extent in acoustic features, at least at some coarse granularity.

Recently, we have seen a large leap in genre classification accuracy. In [31, 30, 29], the authors show that with a perceptually-motivated acoustic feature, and a framework of sparse representation classification (SRC) [45], we move from 82.5% accuracy [3] to up to 93.7%. SRC, which has produced very promising results in computer vision [45, 46] and speech recognition [11, 35], can be thought of as a generalization of  $k$ -nearest neighbors ( $k$ NN) for multiclass classification with many important advantages. It is a global method, in the sense that it classifies based on the entire training set; and it does not rely only on local similarity information as does  $k$ NN. SRC can prevent overcounting of neighborhood information by virtue of its emphasis on sparsity in the representation. Additionally, SRC assigns a weight to each training set sample, thus quantifying the degree of its importance. All of these points make SRC a very strong classification method.

The massive improvement in genre recognition that accompanies this approach motivates many questions, not only about what is working and why it is working so well, but also about how we perceive rich acoustic scenes, and the way we think and talk about music. For instance, are these purely acoustic features so discriminative because they are modeled on the auditory system of humans? Since the features are computed from segments of long duration (30 s), how robust is the method to shorter durations? Do the misclassifications make sense, and does multiclass assignment improve performance? Do the features cluster in a way that resembles the genres we use? Are the genre clusters themselves organized in a sensible way, and do subgenres appear as smaller clusters within larger clusters? Can we compute high-level descriptors from these features, such as rhythm, harmony, or tempo? Is music genre more objective than we think?

In this work, we review the approach proposed in [31], and describe our attempt to reproduce the results, making explicit the many decisions we have had to make to produce the features, and to build the classifier. We find evidence that the perceptual nature of the features has little significant impact on the classifier accuracy. We also see improvement in accuracy when we pose the sparse representation problem using inequality constraints rather than the equality constraints in [31, 30, 29]. Finally, we find that only when we reduce by half the number of classes do we see the reported high accuracies. We make available our MATLAB code, both classification and feature extraction, and with which all figures in this article can be reproduced: <http://imi.aau.dk/~bst/software.html>.

## 2 Background

We now review SRC from a general perspective, and then we review modulation analysis for feature extraction, and its application specifically to music genre recognition. Throughout, we work in a real Hilbert space with inner product  $\langle \mathbf{x}, \mathbf{y} \rangle := \mathbf{y}^T \mathbf{x}$ , and  $p$ -norm  $\|\mathbf{x}\|_p^p := \sum_i |\mathbf{x}_i|^p$ , for  $p \geq 1$ , where  $\mathbf{x}_i$  is the  $i$ th component of the column vector  $\mathbf{x}$ .

## 2.1 Classification via sparse representation in labeled features

Define a set of  $N$  labeled features, each belonging to one of  $C$  enumerated classes

$$\mathcal{D} := \{(\mathbf{x}_n, c_n) : \mathbf{x}_n \in \mathbb{R}^m, c_n \in \{1, \dots, C\}\}_{n \in \{1, \dots, N\}}. \quad (1)$$

And define  $\mathcal{I}_c \subset \{1, \dots, N\}$  as the indices of the features in  $\mathcal{D}$  that belong to class  $c$ . Given an unlabeled feature  $\mathbf{y} \in \mathbb{R}^m$ , we want to determine its class using  $\mathcal{D}$ . In  $k$ NN, we assume that the neighborhood of  $\mathbf{y}$  carries class information, and so we classify it by a majority vote of its  $k$ -nearest neighbors in  $\mathcal{D}$ . Instead of iteratively seeking the best reconstruction of  $\mathbf{y}$  by a single training sample (i.e., its  $i$ th nearest neighbor), we find a reconstruction of  $\mathbf{y}$  by all training samples. Then we chose the class whose samples contribute the most to the reconstruction. We have SRC when we enforce a sparse reconstruction of  $\mathbf{y}$ .

SRC essentially entails finding nearest to an unlabeled feature its linear approximation by class-restricted features. To classify an unlabeled feature  $\mathbf{y}$ , we first find the linear combination of features in  $\mathcal{D}$  that constructs  $\mathbf{y}$  with the fewest number of non-zero weights, regardless of class membership, posed as

$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{a}\|_0 \quad \text{subject to} \quad \mathbf{y} = \mathbf{D}\mathbf{a} \quad (2)$$

where we define the  $m \times N$  matrix  $\mathbf{D} := [\mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_N]$ , and the pseudonorm  $\|\mathbf{a}\|_0$  is defined as the number of non-zero weights in  $\mathbf{a} := [a_1, a_2, \dots, a_N]^T$ . We might not want to enforce equality constraints, and so we can instead pose this

$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{a}\|_0 \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{D}\mathbf{a}\|_2^2 \leq \epsilon^2 \quad (3)$$

where  $\epsilon^2 > 0$  is a maximum allowed error in the approximation. All of this, of course, assumes that we are using features that are additive. We can extend this to non-linear combinations of features by adding such combinations to  $\mathcal{D}$  [35], which can substantially increase the size of the dictionary.

We now define the set of class-restricted weights  $\{\mathbf{a}_c\}_{c \in \{1, 2, \dots, C\}}$

$$[\mathbf{a}_c]_n := \begin{cases} a_n, & n \in \mathcal{I}_c \\ 0, & \text{else.} \end{cases} \quad (4)$$

The non-zero weights in  $\mathbf{a}_c$  are thus only those specific to class  $c$ . From these, we construct the set of  $C$  approximations and their labels  $\mathcal{Y}(\mathbf{a}) := \{\hat{\mathbf{y}}_c(\mathbf{a}) := \mathbf{D}\mathbf{a}_c\}_{c \in \{1, 2, \dots, C\}}$ , and we assign a label to  $\mathbf{y}$  simply by a nearest neighbor criterion

$$\hat{c} := \arg \min_{c \in \{1, \dots, C\}} \|\mathbf{y} - \hat{\mathbf{y}}_c(\mathbf{a})\|_2^2. \quad (5)$$

Thus, SRC picks the class of the nearest approximation of  $\mathbf{y}$  in  $\mathcal{Y}(\mathbf{a})$ .

We cannot, in general, efficiently solve the sparse approximation problems above [8], but there exist several strategies to solve them. We briefly review the convex optimization approaches, but [41] provides a good overview of many

more; and [47] is a large study of SRC using many approaches. Basis pursuit (BP) [6] proposes relaxing strict sparsity with the convex  $\ell_1$ -norm

$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{a}\|_1 \quad \text{subject to} \quad \mathbf{y} = \mathbf{D}\mathbf{a}. \quad (6)$$

And without equality constraints, BP denoising (BPDN) [6] poses this as

$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{a}\|_1 \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{D}\mathbf{a}\|_2^2 \leq \epsilon^2. \quad (7)$$

One could also change the  $\ell_2$  error to  $\ell_1$  to promote sparsity in the error [46, 12]. We have the LASSO [40] when we switch the objective and constraint of BPDN

$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{y} - \mathbf{D}\mathbf{a}\|_2^2 \quad \text{subject to} \quad \|\mathbf{a}\|_1 \leq \rho \quad (8)$$

where  $\rho > 0$ . Furthermore, we can pose the problem in a joint fashion

$$\min_{\mathbf{a} \in \mathbb{R}^N} \frac{1}{2} \|\mathbf{y} - \mathbf{D}\mathbf{a}\|_2^2 + \lambda \|\mathbf{a}\|_1 \quad (9)$$

where  $\lambda > 0$  tunes our preference for sparse solutions versus small error.

Along with using the  $\ell_1$  norm, we can reduce the dimensionality of the problem in the feature space [46]. For instance, the BPDN principle (7) becomes

$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{a}\|_1 \quad \text{subject to} \quad \|\Phi\mathbf{y} - \Phi\mathbf{D}\mathbf{a}\|_2^2 \leq \epsilon^2 \quad (10)$$

where  $\Phi$  is a fat full-rank matrix mapping the features into some subspace. To design  $\Phi$  such that the mapping might benefit classification, we can compute it using information from  $\mathbf{D}$ , e.g., by principal component analysis (PCA) or non-negative matrix factorization (NMF), or we can compute it non-adaptively by random projection. With PCA, we obtain an orthonormal basis describing the directions of variation in the features, from which we define  $\Phi^T$  as the  $d \leq m$  significant directions, i.e., those having the  $d$  largest principal components.

Given  $d \leq m$ , NMF finds a positive full rank  $m \times d$  matrix  $\mathbf{U}$  such that

$$\min_{\mathbf{U} \in \mathbb{R}_+^{m \times d}} \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{U}\mathbf{v}_n\|_2^2 \quad \text{subject to} \quad \mathbf{v}_n \succeq 0. \quad (11)$$

The full-rank matrix  $\mathbf{U}$  contains  $d$  templates that approximate each feature in  $\mathbf{D}$  by an additive combination. Thus the range space of  $\mathbf{U}$  provides a good approximation of the features in  $\mathbf{D}$ , with respect to the mean  $\ell_2$ -norm of their errors. In this case, we make  $\Phi^T := (\mathbf{U}^T \mathbf{U})^{-1} \mathbf{U}^T$ .

Finally, we can reduce feature dimensionality by random projection [7, 4, 21], where we form the entries of  $\Phi$  by sampling from a random variable, e.g., Normal, and without regard to  $\mathbf{D}$ . We normalize the columns to have unit  $\ell_2$ -norm, and ensure  $\Phi$  has full rank. While this approach is computationally simple, its non-adaptivity can hurt classifier performance [21].

## 2.2 Modulation Analysis

Modulation representations of acoustic signals describe the variation of spectral power in scale, rate, time and frequency. This approach has been motivated by the human auditory and visual systems [43, 15, 36, 39, 24]. In the literature, we find two types of modulation representations of acoustic signals, which seemingly have been developed independently. One might see these approaches as a form of feature integration, which aggregate a collection of small scale features.

In [43, 24], the authors model the output of the human primary auditory system as a multiscale spectro-temporal modulation analysis, which [43] terms a “reduced cortical representation” (RCR). To generate an RCR, one first produces an “auditory spectrogram” (AS) approximating the time-frequency distribution of power at the output of the early stage of the auditory system [48]. This involves filtering the signal with bandpass filters modeling the frequency responses of the hair cells along the basilar membrane, then calculating activations of the nerve cells in each band, and finally extracting a spectral power estimate from the activation patterns [48, 43, 24]. In the next step, which models the central auditory system, one performs a “ripple analysis” of the AS, giving the local magnitudes and phases of modulations in scale and modulation rate over time and frequency [43, 24]. This procedure uses 2-D time-frequency modulation-selective filters, equivalent to a multiresolution affine wavelet analysis sensitive to fast and slow upward and downward changes in frequency [43]. To obtain spectro-temporal modulation content [24], one integrates this four-dimensional representation over time and/or frequency.

A similar representation is proposed in [15], where the authors extract modulation information by applying a Fourier transform to the output of a set of bandpass filters modeling the basilar membrane. The magnitude output of this gives a time-varying modulation spectrogram. One could instead apply a wavelet transform to each row of a magnitude spectrogram, and then integrate the power at each scale of each band along the time axis. This produces a modulation rate-scale representation [39].

Motivated by its perceptual foundation [48, 43], and success in automatic sound discrimination [44, 24], the work of [28] appears to be the first to use modulation analysis features for music genre recognition, which they further refine in [31, 30, 32]. In [28], the authors use the toolbox by the Neural Systems Laboratory (NSL) [33, 34] to derive an RCR and perform a ripple analysis. They then average this across time to produce tensors of power distributed in modulation rate, scale, and acoustic frequency. While the features in that work are built from the RCR [43, 24], the features used in [31, 30] are a joint scale-frequency analysis [39] of an AS created from the model in [48]. This feature, which they call an “auditory temporal modulation” (ATM), describes power variation over modulation scale in each primary auditory cortex channel.

## 3 Recreating the Features and Classifier

In this section, we first describe how we generate ATM features, which are described in part in [31, 32]; and then we describe the approach to classify an ATM using SRC [31].

### 3.1 Building Auditory Temporal Modulations

The authors take 30 seconds of music, downsample it to 16 kHz, then make it zero mean and unit variance. They then compute an AS following the model of the primary auditory system of [48], except they use a constant-Q transform of 96 bandpass filters covering a 4-octave range (24 filters per octave), whereas [48] uses an affine wavelet transform of 64 scales covering 5 octaves from about 173 Hz to 5.9 kHz. Finally, they pass each channel of the AS through a Gabor filterbank sensitive to particular modulation rates, and form the ATM by integrating the energy output at each filter.

To create ATMs, we have tried to follow as closely as possible the description in [31, 32]. We first generate a constant-Q filter bank with 97 bands spaced over a little more than four octaves, with  $N_f = 24$  filters per octave. We center the first filter at 200 Hz because that is specified in [48]. The last filter is thus centered on 3200 Hz. Since in [48] the model of the final stage of the primary audio cortex computes first-order derivatives across adjacent frequency bands, we end up with a 96 band AS as specified in [31, 32].

We create our constant-Q filter bank as a set of finite impulse response filters designed by the windowing method [25]. Since it is not mentioned in [48, 31, 32], we make all filters independent, and to have the same gain. To generate the impulse responses of our filterbank, we modulate a prototype lowpass window to logarithmically spaced frequencies. Because of its good low passband characteristic, we use a Hamming window, which for the  $k$ th filter ( $k \geq 1$ ) produces the impulse response sampled at  $F_s$  Hz

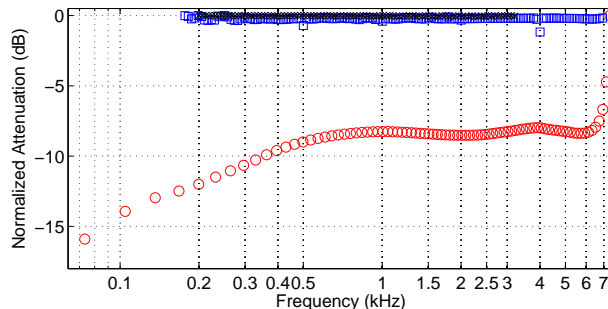
$$h_k(n) := \gamma_k \left[ 0.54 - 0.46 \cos \left( \frac{2\pi n}{l_k} \right) \right] e^{j2\pi\omega_k n/F_s}, \quad 0 \leq n < l_k \quad (12)$$

with a modulation frequency  $\omega_k := f_{\min} 2^{(k-1)/N_f}$  Hz, and length in samples

$$l_k := \left\lceil \frac{q}{2^{k/N_f} - 2^{(k-1)/N_f}} \frac{F_s}{f_{\min}} \right\rceil. \quad (13)$$

We set the gain  $\gamma_k$  such that there is no attenuation at the  $k$ th center frequency, i.e.,  $|\mathcal{F}\{h_k(n)\}(\omega_k)| = 2$ , where  $\mathcal{F}\{x(n)\}(\omega)$  is the Fourier transform of  $x(n)$  evaluated at frequency  $\omega$ . The factor  $q > 0$  tunes the width of the main lobe. We choose  $q \approx 1.316$  such that adjacent filters overlap at their -3 dB stopband.

This model of the basilar membrane is simplified considering its non-adaptive and uniform nature, e.g., it does not take into account masking and equal loudness curves. An alternative model of the cochlea is given by Lyon [20], which involves a filterbank with center frequencies spread uniformly below a certain frequency, and logarithmically above [37]. Figure 1 shows that the Lyon model attenuates single sinusoids at frequencies tuned to the center frequencies of its filterbank. Our filterbank uniformly passes these frequencies, albeit over a smaller four octave range [31, 32] assumed to begin at 200 Hz. Figure 1 also shows that the filterbank of the NSL model [34] by and large has a uniform attenuation.



**Fig. 1.** Attenuations of single sinusoids with the same power, at frequencies identical to center frequencies in the filterbanks. (x) Our constant-Q filter bank. (o) Lyon passive ear model [20, 37]. (□) NSL ear model [34].

We pass through our constant-Q filter bank a sampled, zero-mean and unit-variance acoustic signal  $y(n)$  [31, 32], which produces for the  $k$ th filter the output

$$y_k(n) := \sum_{m=0}^{l_k-1} h_k(m)y(n-m-\Delta_k) \quad (14)$$

where  $\Delta_k > 0$  is the group delay of the  $k$ th filter at  $\omega_k$ . This delay correction is necessary because the filters we use to model the basilar membrane have different lengths. This correction is unnecessary in the implementation of the Lyon [37] or NSL models [34], since they use filter structures having identical delays.

As in [31, 32], we next take the sample wise difference in each band

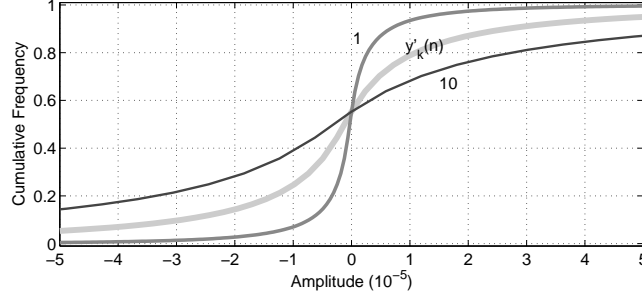
$$y'_k(n) := y_k(n) - y_k(n-1). \quad (15)$$

which models the action potential of the hair cell [48]. This now goes through a non-linear compression, followed by a low pass filter modeling leakage in the hair cell membrane. Referring to [48], we see the compression can be modeled as a sigmoidal function, and that the output of the  $k$ th channel is

$$g_k(n) := \frac{1}{1 + e^{-\gamma y'_k(n)}} - \frac{1}{2} \quad (16)$$

where  $\gamma > 0$  is depends on sound pressure level [48]. Furthermore, “... saturation in a given fiber is limited to 30-40 dB” [48], implying  $\gamma$  is somehow set adaptively. In reality, we cannot equate the values of the digital samples in  $y'_k(n)$  with the physical pressure embodied in this compression. However, working naively, we might absorb into  $\gamma$  such a conversion, and find some value that actually compresses. Figure 2 shows the cumulative distribution of amplitudes input to the compressor (15) with a 30 second music signal having unit energy [31, 32]. For  $\gamma = 1$ , we see that this distribution is compressed, whereas setting  $\gamma = 10$  results in an expansion. Thus, we set  $\gamma = 1$  independent of the input, and assume it compresses  $y'_k(n)$  from any 30 second music signal scaled to have unit energy.





**Fig. 2.** Cumulative distributions of amplitude input to compressor ( $y'_k(n)$ ), and output as a function of  $\gamma$  (labeled).

The compressor output  $g_k(n)$  is then smoothed by the hair cell membrane and attendant leakage [48, 32], which passes frequencies only up to 4-5 kHz [48]. Thus, we pass each  $g_k(n)$  through a 6th-order Butterworth filter having a cutoff frequency of 4 kHz, producing  $f_k(n)$ . This is then processed by a “lateral inhibitory network,” described in [48], which detects discontinuities in the response. This entails a spatial derivative across channels with smoothing, a half-wave rectifier, and then integration; but [31, 32] does not specify smoothing, and states the process can be approximated by a first order derivative across logarithmic frequency. Thus, we compute for channel  $s \in \{1, \dots, 96\}$

$$v_s(n) := [f_{s+1}(n) - f_s(n)]\mu[f_{s+1}(n) - f_s(n)] \quad (17)$$

where  $\mu(u) = 1$  if  $u \geq 0$ , and zero otherwise.

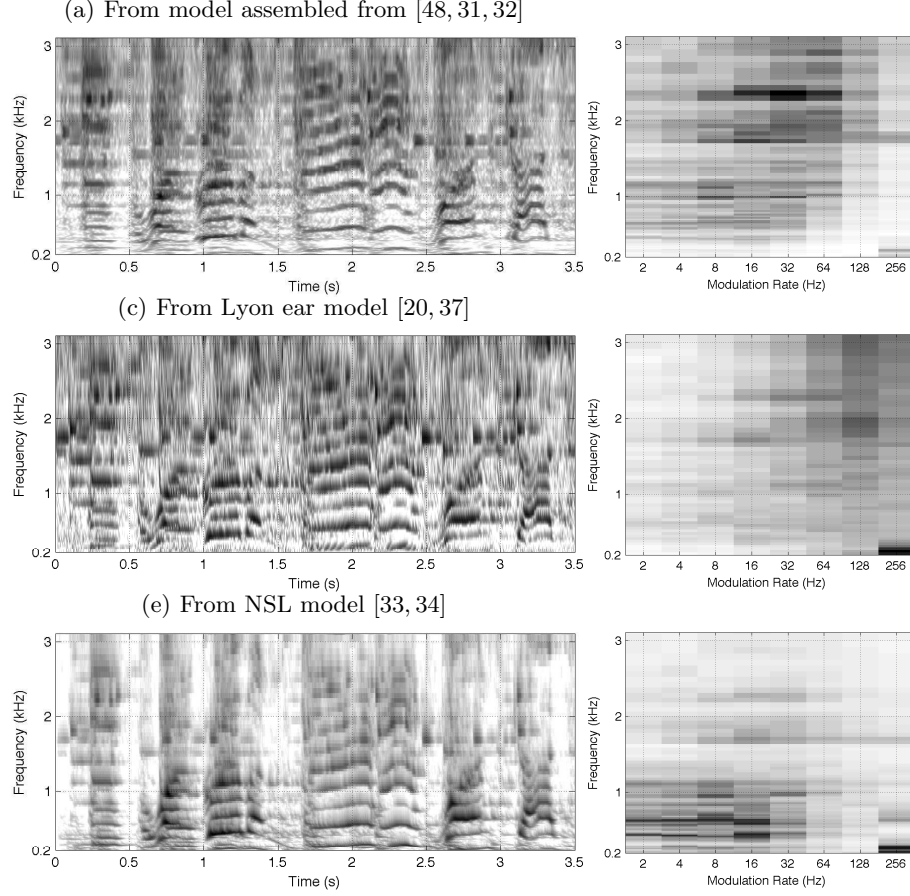
In the final step, we integrate the output with “a [possibly rectangular window with a] long time constant (10-20 ms)” [48], or a 2-8 ms exponential window [31, 32]. Thus, we compute the  $n$ th sample of the  $k$ th row of the AS by

$$A_k(n) := \sum_{m=0}^{\lfloor F_s \tau \rfloor} v_s(n-m)e^{-m/F_s \tau} \quad (18)$$

where we define  $\tau := 8$  ms. This completes the first step of building an ATM.

Figure 3 compares the resulting AS from our model built from interpreting [48, 31, 32], that of the auditory model designed by Lyon [20, 37], and the cortical representation from the NSL model [33, 34]. The Lyon model uses 86 bands non-uniformly spread over a little more than 6.5-octaves in 80 – 7630 Hz [20, 37], whereas the NSL model covers 5.33 octaves with 24 filters per octave logarithmically spread over 180 – 7246 Hz [33, 34]. Though the frequency range of those models are larger, we only show a four-octave frequency range [31, 32].

To generate an ATM, [31, 32] describe first performing a multiresolution wavelet decomposition of each row of an AS, and then integrating the squared output across the translation axis. Based on experimental evidence [36], the authors use a set of Gabor filters sensitive to eight modulation rates  $\{2, 4, 8, \dots, 256\}$



**Fig. 3.** Auditory spectrograms (left) and their auditory temporal modulations (right).

Hz. We assume this Gabor filterbank can be assembled as follows. We define the sampled impulse response truncated to length  $N_l$  of our complex Gabor filter tuned to a modulation rate  $f_0 2^l \geq 0$  Hz, and of scale  $F_s \alpha / f_0 2^l > 0$

$$\psi(n; f_0 2^l) := \frac{f_0 2^l}{F_s \alpha} \left[ e^{-(f_0 2^l / \alpha)^2 ((n - N/2) / F_s)^2} e^{j 2 \pi f_0 2^l n / F_s} - \mu_l \right] \quad (19)$$

for  $n = 0, \dots, N_l - 1$ , where we define  $\mu_l$  such that  $\psi(n; f_0 2^l)$  has zero mean. The normalization constant assures uniform attenuation at each modulation frequency, as used in joint scale-frequency analysis [39]. We set  $\alpha = 256/400$  and  $N_l = 4 F_s \alpha / f_0 2^l$ . Since a Gabor filter tuned to a low frequency has a high DC component, we make each row of the AS zero mean, thus producing  $A'_k(n)$ . Passing the  $k$ th row of this AS through the  $l$ th channel ( $l \in \{0, 1, \dots, 7\}$ ) of the Gabor filterbank produces the convolution  $R_{k,l}(n) := [\psi(m; f_0 2^l) \star A'_k(m)](n)$ . Finally, as in [31, 32], we sum the squared modulus of the output sampled at all

wavelet translations, producing the  $(k, l)$  element of the ATM

$$[\mathbf{A}]_{kl} := \sum_{p \in \mathbb{Z}} |R_{k,l}(p \lfloor F_s \alpha / f_0 2^{l+1} \rfloor)|^2 \quad (20)$$

where  $p$  is an integer multiplying the wavelet translations, which we assume is half the wavelet scale.

To the right of each AS in Fig. 3 we see the resulting ATM. Portions of these ATMs appear similar, with major differences in scaling and feature dimensionality. Within the four octave range specified in [31, 32], the dimensionality of the vectorized features are: 768 for our own ATM [31, 32], 416 for that created from the model by Lyon [20, 37], and 800 using the NSL model [33, 34].

### 3.2 Classifying Genre by Auditory Temporal Modulations

Given a set  $\mathcal{D}$  of vectorized ATM features, each associated with a single music genre, we can use the machinery of SRC to label an unknown vectorized ATM  $\mathbf{y}$ . Following [31], we first make all features of  $\mathcal{D}$  have unit  $\ell_2$ -norm, as well as the test feature  $\mathbf{y}$ . We next solve the BP optimization problem posed in [31]

$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{a}\|_1 \quad \text{subject to} \quad \Phi \mathbf{y} = \Phi \mathbf{D} \mathbf{a} \quad (21)$$

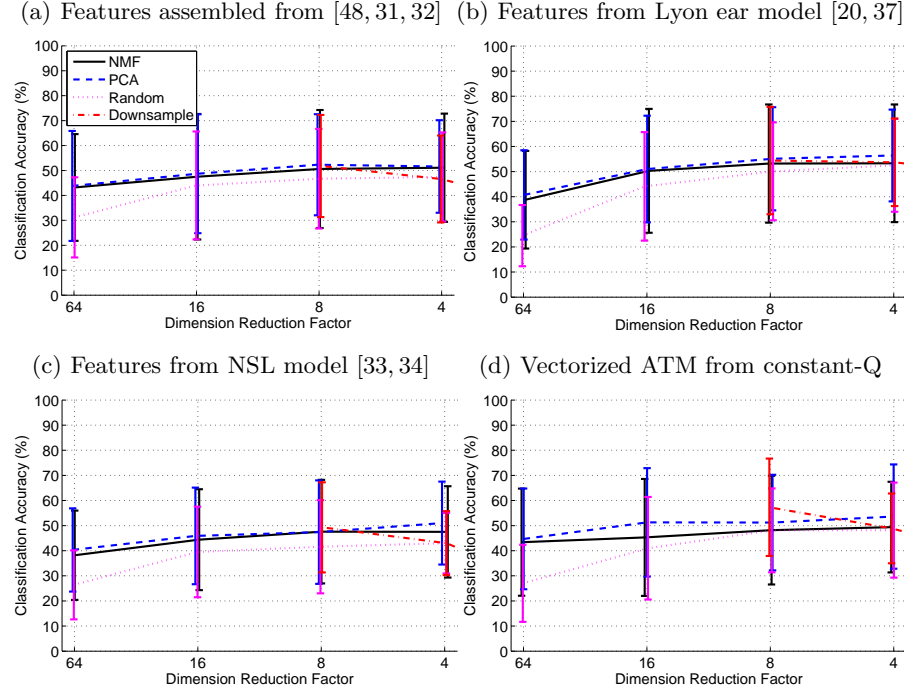
where  $\Phi$  reduces the features by, e.g., PCA. Finally, to classify  $\mathbf{y}$ , we construct the set of weights in (4), and assign a single genre label using the criterion (5).

Since we are working with real vectors, we can solve (21) as a linear program [6], for which numerous solvers have been implemented, e.g., [5, 10, 2, 14]. Because of its speed, we choose as the first step the root-finding method of the SGPL1 solver [2]. If this fails to find a solution, then we use the primal-dual method of  $\ell_1$ -Magic [5], which takes as its starting point the minimum  $\ell_2$ -norm solution  $\mathbf{a}_2 = (\Phi \mathbf{D})^\dagger \mathbf{y}$ . This initial solution satisfies the constraints of (21) as long as  $\Phi \mathbf{D}$  has full rank, but probably is not the optimal solution. If the solution  $\hat{\mathbf{a}}$  does not satisfy  $\|\Phi \mathbf{y} - \Phi \mathbf{D} \hat{\mathbf{a}}\|_2^2 < 10^{-16}$  (numerical precision), we set  $\hat{\mathbf{a}} = \mathbf{a}_2$ .

## 4 Experimental Results

As in [31], we use the music genre dataset of [42], which has 1000 half-minute sound examples drawn from music in 10 broad genres: blues, classical, country, disco, hiphop, jazz, metal, pop, reggae, and rock. We define  $\Phi$  by PCA, NMF, or random sampling; and as in [31], we test dimension reduction by factors of  $\{64, 16, 8, 4\}$ , e.g., a feature vector of 768 dimensions is reduced by a factor of four to 192 dimensions. We also test downsampling the features, but we define it as vectorizing the result of lowpass filtering and decimating each column of the ATM (20). It is not clear how downsampling is done in [31]. In our case, a factor of  $f$  downsampling results in a vectorized feature of dimension  $8 \lceil 96/f \rceil$  when using our 96-channel features. Finally, as done in [3, 31], we use stratified 10-fold cross-validation to test the classifier.

Figure 4 shows our classification results for four different features, including just the vectorized modulation-analysis of the magnitude output of the constant-Q filterbank that precedes (15). For all features, we see no mean accuracies above

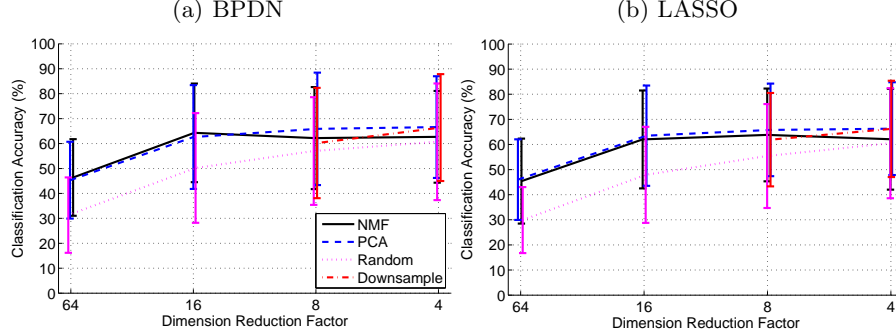


**Fig. 4.** Mean classification accuracy (10 classes) of SRC based on (21) for four different feature design methods, four dimension reduction methods, and several reduction factors. Overlaid is the  $\pm 1$  standard deviation about the mean. (We add a slight x-offset to each bar for readability.)

57.3%, which we achieve with features derived from the constant-Q filterbank, and downsampled by a factor of 8. Though this feature is nothing more than a modulation analysis of the output of a logarithmic filterbank, we see its performance is not significantly different from the other features that are modeling the entire primary auditory cortex. With the exception of random projection, we do not see any significant performance variation between the approaches for dimensionality reduction and the dimension of the features used in SRC. although the trend is that mean accuracy decreases with fewer dimensions.

The large differences between our results and those for the same dataset and experimental protocol reported in [31] — mean accuracy of around 91% for features reduced a factor of 4 by NMF — signify something is not correct. We have verified with synthetic signals with known modulations that our modulation analysis is working [31]. We have tested and confirmed on a handwritten digits dataset [16] that our feature reduction is working, and that our SRC classifier performs comparably to other classifiers. We believe then that these differences in results come from several things, three of which are significant.

First, it is common in classification to preprocess features to remove problems that can come from comparing dimensions with different scales. This entails making the values of each row of  $\mathbf{D}$  be in  $[0, 1]$  by finding and subtracting the minimum, and then dividing by the difference of the maximum and minimum.



**Fig. 5.** Mean classification accuracy (10 classes) of SRC based on the BPDN (22) and LASSO (23), with standardized data and a normalized projected dictionary, for ATM features from the Lyon model [20, 37], four dimension reduction methods, and several reduction factors. Overlaid is the  $\pm 1$  deviation about the mean.

This “standardization” is not mentioned in [31, 32], or the original reference to SRC [46]. When we rerun the experiments above with standardized data, we see the mean accuracy increases, but does not exceed the highest of 64% for the NSL features reduced in dimensionality a factor of 4 by PCA.

The second problem is posing the sparse representation with equality constraints in (21), which forces the sparse representation algorithm to model a feature exactly when instead we just want to find a good model of our feature. We thus pose the problem instead using BPDN [6] (7)

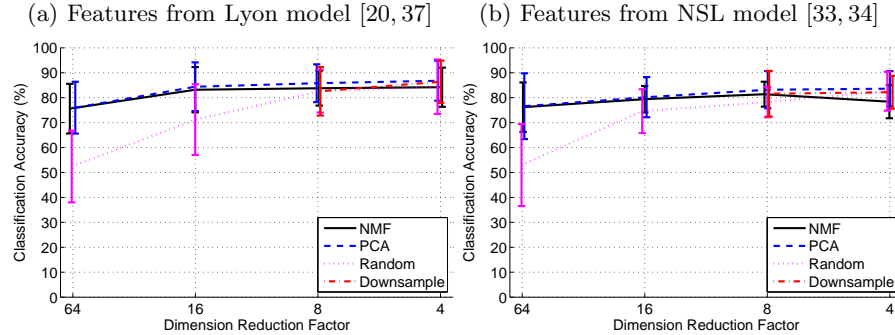
$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\mathbf{a}\|_1 \quad \text{subject to} \quad \|\Phi\mathbf{y} - \Phi\mathbf{D}\mathbf{a}\|_2^2 \leq \epsilon^2. \quad (22)$$

or as the LASSO [40] (8)

$$\min_{\mathbf{a} \in \mathbb{R}^N} \|\Phi\mathbf{y} - \Phi\mathbf{D}\mathbf{a}\|_2^2 \quad \text{subject to} \quad \|\mathbf{a}\|_1 \leq \rho. \quad (23)$$

Solving these can produce an informative representation using few features instead of an exact fit by many.

For standardized features, and normalizing the columns of the mapped dictionary  $\Phi\mathbf{D}$ , Fig. 5(a) shows the results of using BPDN (22) where we define  $\epsilon^2 = 10^{-4}$ ; and Fig. 5(b) shows the results when we pose the problem as the LASSO (23) where we define  $\rho = 1$ . (We show only the results from the Lyon model since the other features did not give significantly different results.) In both cases, we use SGPL1 [2] with at most 100 iterations, and use the result whether it is in the feasible set or not. This is different from our approach to solving (21), where we run  $\ell_1$ -Magic [5] if SPGL1 fails, and then use the minimum  $\ell_2$ -norm solution if this too fails. In our experiments, we see (23) is solved nearly all the time for  $\rho = 1$ , and (22) is solved only about 5% of the time; yet we see no significant difference between the accuracies of both cases. With these changes, we see a slight increase in mean accuracies to about 68% for the features derived from the Lyon model [20, 37].



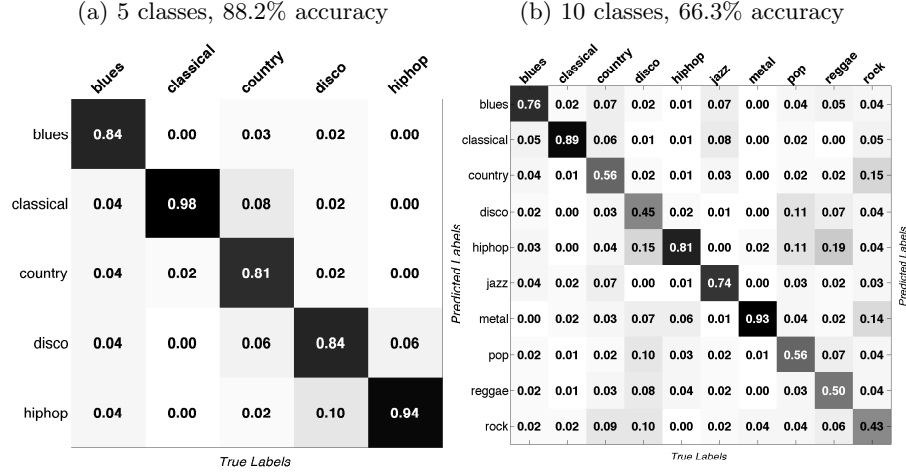
**Fig. 6.** Mean classification accuracy (5 classes) of SRC based on the LASSO (23), with standardized data and a normalized projected dictionary, for ATM features derived from a larger frequency range, four dimension reduction methods, and several reduction factors. Overlaid is the  $\pm 1$  deviation about the mean.

The third significant problem comes from the definition of the features. We find that accuracy improves slightly if we use features from a wider frequency range than four octaves, e.g., all 86 bands of the AS from the Lyon model, covering 80 – 7630 Hz [20, 37], or all 128 bands of the AS from the NSL model [33, 34], logarithmically spread over 180 – 7246 Hz. With these changes, however, our mean accuracies do not exceed 70%, still well below the 91% reported [31].

The only way we have obtained something close to such a high accuracy is by, in addition to the changes above, limit the classification problem to these five genres: blues, classical, country, disco, and hiphop. Figure 6 shows our results using the expanded features derived from the Lyon and NSL models; and Fig. 7 shows the confusion matrices of the two best performing features and classifiers we have found for the 5- and 10-class problems. These do show sensible behaviors: classical is rarely confused for another genre; rock is often confused for country and metal; rock and metal are confused; reggae, hiphop, pop and disco are confused. Listening to the examples that have been misclassified in Fig. 7(a) show some of them make sense, such as the disco-confused but disco-esque “Let Your Love Flow” (country, no. 13 in dataset [42]). Another is the blues-confused but country-labeled “Running Bear Little White Dove” by Johnny Preston (no. 48). However, the two classical segments misclassified as country are a portion of Gershwin’s “Rhapsody in Blue” (no. 48), and an operatic duo in Italian (no. 54). More confusing is that we find that the pre-autotune George McCrae disco hit “I can’t leave you alone” (no. 39), and Willie Nelson’s “Oh they tell me (unclouded day)” (country, no. 77) are both misclassified as classical.

## 5 Conclusion

Were the difficult problem of music similarity solved, it would present a wonderful tool for exploring many interesting questions; and were it solved using solely acoustic features, it would say something significant about a process that appears influenced by much more than sound. Though the approach of [31] appears extremely promising in light of state of the art — it is based on a perceptually-informed acoustic feature and a classification method built upon sparse repre-



**Fig. 7.** Confusion matrices (frequency of classification) found by stratified 10-fold cross-validation using features derived from the Lyon Ear Model [20, 37] reduced in dimensionality a factor of 4 by PCA, and classification using the LASSO (23).

sentations in exemplars, which has its own biological motivations, e.g., [19, 18] — we have not been able to reproduce their results without reducing the number of classes from 10 to 5. We have shown in as much detail possible the variety of decisions we have had to make in our work, and provide all our code for discussion. Though our results point to a negative conclusion, we have shown that the modeling the primary auditory cortex may not provide additional discriminative information for genre. We have also shown that relaxing the constraints in the sparse representation component of SRC, and standardizing the features of the dictionary, improves classification accuracy.

### Acknowledgments

B. L. Sturm is supported in part by Independent Postdoc Grant 11-105218 from Det Frie Forskningsrd. The authors would like to thank Dr. Costas Kotropoulos and Yannis Panagakakis for helpful discussion about their work.

### References

1. Baumann, S., Pohle, T., Vembu, S.: Towards a socio-cultural compatibility of MIR systems. In: Proc. ISMIR. pp. 460–465. Barcelona, Spain (Oct 2004)
2. van den Berg, E., Friedlander, M.P.: Probing the pareto frontier for basis pursuit solutions. SIAM J. on Scientific Computing 31(2), 890–912 (Nov 2008)
3. Bergstra, J., Casagrande, N., Erhan, D., Eck, D., Kégl, B.: Aggregate features and adaboost for music classification. Mach. Learn. 65(2-3), 473–484 (2006)
4. Bingham, E., Mannila, H.: Random projection in dimensionality reduction: Application to image and text data. In: Proc. Int. Conf. Knowledge Discovery Data Mining. pp. 245–250. San Francisco, CA (Aug 2001)
5. Candès, E., Romberg, J.:  $\ell_1$ -magic: Recovery of sparse signals via convex programming. Tech. rep., Caltech, Pasadena, CA, USA (2005)

6. Chen, S.S., Donoho, D.L., Saunders, M.A.: Atomic decomposition by basis pursuit. *SIAM J. Sci. Comput.* 20(1), 33–61 (Aug 1998)
7. Dasgupta, S.: Experiments with random projection. In: *Proc. Conf. Uncertainty in Artificial Intelligence*. pp. 143–151. Stanford, CA, USA (June 2000)
8. Davis, G., Mallat, S., Avellaneda, M.: Adaptive greedy approximations. *J. Constr. Approx.* 13(1), 57–98 (Jan 1997)
9. Fabbri, F.: A theory of musical genres: Two applications. In: *Proc. First International Conference on Popular Music Studies*. Amsterdam, The Netherlands (1980)
10. Figueiredo, M., Nowak, R., Wright, S.J.: Gradient projection for sparse reconstruction: Application to compressed sensing and other inverse problems. *IEEE J. Sel. Topics Signal Process.* 1(4), 586–597 (Dec 2007)
11. Gemmeke, J., ten Bosch, L., L.Boves, Cranen, B.: Using sparse representations for exemplar based continuous digit recognition. In: *Proc. EUSIPCO*. pp. 1755–1759. Glasgow, Scotland (Aug 2009)
12. Giacobello, D., Christensen, M., Murthi, M.N., Jensen, S.H., Moonen, M.: Enhancing sparsity in linear prediction of speech by iteratively reweighted  $\ell_1$ -norm minimization. In: *Proc. IEEE Int. Conf. Acoustics, Speech, Signal Process.* Dallas, TX (Mar 2010)
13. Gjerdingen, R.O., Perrott, D.: Scanning the dial: The rapid recognition of music genres. *J. New Music Research* 37(2), 93–100 (Spring 2008)
14. Grant, M., Boyd, S.: CVX: Matlab software for disciplined convex programming, version 1.21. <http://cvxr.com/cvx> (Apr 2011)
15. Greenberg, S., Kingsbury, B.E.D.: The modulation spectrogram: in pursuit of an invariant representation of speech. In: *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.* pp. 1647–1650. Munich, Germany (Apr 1997)
16. LeCun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proc. IEEE* 86(11), 2278–2324 (Nov 1998)
17. Lena, J.C., Peterson, R.A.: Classification as culture: Types and trajectories of music genres. *American Sociological Review* 73, 697–718 (Oct 2008)
18. Lewicki, M.S.: Efficient coding of natural sounds. *Nature Neuroscience* 5(4), 356–363 (Mar 2002)
19. Lewicki, M.S., Sejnowski, T.J.: Learning overcomplete representations. *Neural Computation* 12, 337–365 (Feb 2000)
20. Lyon, R.F.: A computational model of filtering, detection, and compression in the cochlea. In: *Proc. ICASSP*. pp. 1282–1285 (1982)
21. Majumdar, A., Ward, R.K.: Robust classifiers for data reduced via random projections. *IEEE Trans. Systems, Man, Cybernetics* 40(5), 1359–1371 (Oct 2010)
22. Mayer, R., Neumayer, R., Rauber, A.: Rhyme and style features for musical genre classification by song lyrics. In: *Proc. Int. Symp. Music Info. Retrieval* (2008)
23. McKay, C., Fujinaga, I.: Music genre classification: Is it work pursuing and how can it be improved? In: *Proc. Int. Symp. Music Info. Retrieval* (2006)
24. Mesgarani, N., Slaney, M., Shamma, S.A.: Discrimination of speech from nonspeech based on multiscame spectro-temporal modulations. *IEEE Trans. Audio, Speech, Lang. Process.* 14(3), 920–930 (May 2006)
25. Mitra, S.K.: *Digital Signal Processing: A Computer Based Approach*. McGraw Hill, 3 edn. (2006)
26. Pachet, F., Cazaly, D.: A taxonomy of musical genres. In: *Proc. Content-based Multimedia Information Access Conference*. Paris, France (Apr 2000)
27. Pampalk, E., Flexer, A., Widmer, G.: Hierarchical organization and description of music collections at the artist level. In: *Research and Advanced Technology for Digital Libraries*. pp. 37–48 (2005)



28. Panagakis, Y., Benetos, E., Kotropoulos, C.: Music genre classification: A multi-linear approach. In: Proc. ISMIR. pp. 583–588. Philadelphia, PA (Sep 2008)
29. Panagakis, Y., Kotropoulos, C.: Music genre classification via topology preserving non-negative tensor factorization and sparse representations. In: Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. pp. 249–252. Dallas, TX (Mar 2010)
30. Panagakis, Y., Kotropoulos, C., Arce, G.R.: Music genre classification using locality preserving non-negative tensor factorization and sparse representations. In: Proc. Int. Symp. Music Info. Retrieval. pp. 249–254. Kobe, Japan (Oct 2009)
31. Panagakis, Y., Kotropoulos, C., Arce, G.R.: Music genre classification via sparse representations of auditory temporal modulations. In: Proc. European Signal Process. Conf. pp. 1–5. Glasgow, Scotland (Aug 2009)
32. Panagakis, Y., Kotropoulos, C., Arce, G.R.: Non-negative multilinear principal component analysis of auditory temporal modulations for music genre classification. *IEEE Trans. Acoustics, Speech, Lang. Process.* 18(3), 576–588 (Mar 2010)
33. Ru, P.: Cortical Representations and Speech Recognition. Ph.D. thesis, University of Maryland, College Park, MD, USA (Dec 1999)
34. Ru, P.: Multiscale multirate spectro-temporal auditory model. Tech. rep., Neural Systems Laboratory, University of Maryland College Park (2001), <http://www.isr.umd.edu/Labs/NSL/Software.htm>
35. Sainath, T.N., Carmi, A., Kanevsky, D., Ramabhadran, B.: Bayesian compressive sensing for phonetic classification. In: Proc. ICASSP (2010)
36. Shamma, S.A.: Encoding sound timbre in the auditory system. *IETE J. Research* 49(2), 145–156 (Mar-Apr 2003)
37. Slaney, M.: Auditory toolbox. Tech. rep., Interval Research Corporation (1998)
38. Sordo, M., Celma, O., Blech, M., Guaus, E.: The quest for musical genres: Do the experts and the wisdom of crowds agree? In: Proc. ISMIR (2008)
39. Sukittanon, S., Atlas, L.E., Pitton, J.W.: Modulation-scale analysis for content identification. *IEEE Trans. Signal Process.* 52(10), 3023–3035 (Oct 2004)
40. Tibshirani, R.: Regression shrinkage and selection via the lasso. *J. Royal Statist. Soc. B* 58(1), 267–288 (Jan 1996)
41. Tropp, J.A., Wright, S.J.: Computational methods for sparse solution of linear inverse problems. *Proc. IEEE* 98(6), 948–958 (June 2010)
42. Tzanetakis, G., Cook, P.: Musical genre classification of audio signals. *IEEE Trans. Speech Audio Process.* 10(5), 293–302 (July 2002)
43. Wang, K., Shamma, S.A.: Spectral shape analysis in the central auditory system. *IEEE Trans. Speech Audio Process.* 3(5), 382–395 (Sep 1995)
44. Woolley, S.M.N., Fremouw, T.E., Hsu, A., Theunissen, F.E.: Tuning for spectro-temporal modulations as a mechanism for auditory discrimination of natural sounds. *Nature Neuroscience* 8(10), 1371–1379 (Oct 2005)
45. Wright, J., Ma, Y., Mairal, J., Sapiro, G., Huang, T., Yan, S.: Sparse representation for computer vision and pattern recognition. *Proc. IEEE* 98(6), 1031–1044 (June 2009)
46. Wright, J., Yang, A.Y., Ganesh, A., Sastry, S.S., Ma, Y.: Robust face recognition via sparse representation. *IEEE Trans. Pattern Anal. Machine Intell.* 31(2), 210–227 (Feb 2009)
47. Yang, A.Y., Ganesh, A., Zhou, Z., Sastry, S.S., Ma, Y.: A review of fast  $\ell_1$ -minimization algorithms for robust face recognition. (preprint) (2010), <http://arxiv.org/abs/1007.3753>
48. Yang, X., Wang, K., Shamma, S.A.: Auditory representations of acoustic signals. *IEEE Trans. Info. Theory* 38(2), 824–839 (Mar 1992)