

SPARSE CHROMA ESTIMATION FOR HARMONIC NON-STATIONARY AUDIO

Maria Juhlin, Ted Kronvall, Johan Swärd, and Andreas Jakobsson

Centre for Mathematical Sciences, Lund University, Sweden.
email: {juhlin, ted, js, aj}@maths.lth.se

ABSTRACT

In this work, we extend on our recently proposed block sparse chroma estimator, such that the method also allows for signals with time-varying envelopes. Using a spline-based amplitude modulation of the chroma dictionary, the refined estimator is able to model longer frames than our earlier approach, as well as to model highly time-localized signals, and signals containing sudden bursts, such as trumpet or trombone signals, thus retaining more signal information than other methods for chroma estimation. The performance of the proposed estimator is evaluated on a recorded trumpet signal, clearly illustrating the improved performance, as compared to other used techniques.

Index Terms— chromagram, amplitude modulation, block sparsity, convex optimization, ADMM

1. INTRODUCTION

Music is an art-form that most enjoy. Even more so today than earlier, as personalized computers and smart telephones have enabled ubiquitous music listening and allow everyone to be their own hobby-DJ. When listening, learning, composing, mixing, and identifying music, there are a number of aspects and approaches one may utilize, such as a composition's timbre, pitch, tempo, beat, rhythm, and chroma (see, e.g. [1]). Many such features involve analyzing the spectral content of the signal. Pitch, as a musical concept, is an ordinal scale of sounds which is related to, but not necessarily as cardinally specific, as the frequency scale. A single pitch is from a spectral point of view a combination of many narrowband spectral peaks, which typically share an integer relationship in terms of their frequencies. In this sense, the pitch is typically defined by the component of lowest frequency, i.e., the fundamental, whereas the other frequencies are referred to as its harmonics. The number of harmonics in a certain pitch, as well as the magnitude power of these, varies greatly between different sounds. Identifying pitches in a way similar to our human perception has proved to be a difficult estimation problem. Partly, this difficulty is due

to octaves; two pitches where one has exactly twice the fundamental frequency as the other are referred to as being octave equivalent as the distance in pitch by a factor of two is called an octave. The octave equivalence is a central part of the Western musicological system. Within each octave, the Western musical system defines twelve so called semi-tones, or chromas. The same chroma is then cyclically defined to each doubling of fundamental frequency, for all twelve chromas [2]. Methods for multi-pitch estimation in audio have been thoroughly examined in the literature (see e.g., [3–5], and the references therein). Typically, trouble arises when the complexity of the audio signal increases, such that there are simultaneously two or more pitches present, played by more than one instrument. Separating these complex combinations of components in the signal often proves difficult, even if the harmonic structure of the signal is taken into account. As introduced in [6], by collecting the pitches in groups in accordance with their respective chroma, we simplify the estimation, only focusing on chroma, while retaining much of the musical information. Chroma features are widely used in applications such as cover song detection, transcription, and recommender systems (see, e.g. [7–9]). Most methods for chroma estimation begin with some pitch estimation, which then maps into its respective chroma. In this approach, some take the harmonic structure into account, and others do not. The commonly used method by Ellis [10] is formed via a time-smoothed version of the short time Fourier transform, whereas the method by Müller and Ewert uses a filterbank approach [11]. Neither of these use the pitches' harmonic structure for estimation. On the other hand, taking this structure into account often requires knowledge of the number of pitches and their respective number of harmonics, which is notoriously difficult to obtain for multi-pitch signals. Instead, we propose to estimate the present chromas using a sparse model reconstruction approach, where explicit model orders are not required. These parameters are instead controlled implicitly using some tuning parameters, which may typically be set using cross-validation, or by using some simple heuristics. Recently, we proposed such a technique [6], generalizing an earlier work exploiting block sparsity for multi-pitch estimation [12]. Herein, we extend on this model by generalizing it in accordance with the methods presented in [13, 14]. The proposed extension allows the signal to have a time vary-

This work was supported in part by the Swedish Research Council, Carl Trygger's foundation, and the Royal Physiographic Society in Lund.

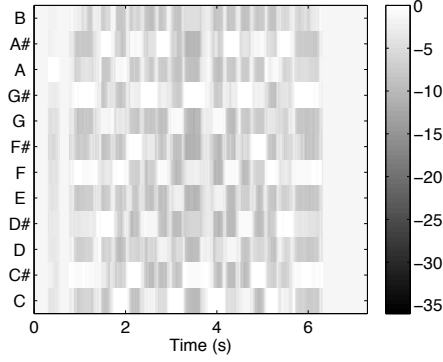


Fig. 1. The normalized log-chromagram for the trumpet scale using the method by Müller and Ewert.

ing amplitude, extending the usability of the method to also allow for highly non-stationary signals, or signals with sudden bursts, like trumpets, whose nature may easily be misinterpreted using ordinary chroma selection techniques. As in [13], the extended model uses a spline basis to detail the time-varying envelope of the signal, thereby enabling the amplitudes to evolve smoothly with time. The time-localization offered by the new method also enables a better signal matching, such that more overall information is retained in the resulting chromagram. The performance of the proposed estimator is illustrated using a recorded trumpet scale, clearly illustrating the improved performance as compared to typical reference methods, and to our earlier proposed estimator.

2. THE SIGNAL MODEL

As shown in [6], a harmonically related audio signal may be well modeled as a sum of K distinct pitch signals, each consisting of L_k harmonically related sinusoids with normalized fundamental frequencies f_k . In this work, we allow the amplitudes of the harmonic components to vary over time, such that

$$y(t) = \sum_{k=1}^K \sum_{\ell=1}^{L_k} \alpha_{k,\ell}(t) e^{i2\pi f_k \ell t}, \quad (1)$$

for $t = 1, \dots, N$, where $\alpha_{k,\ell}(t)$ represents the amplitude of the ℓ th harmonic of the k th pitch, at time instant t . Reminiscent to [13], we model the amplitudes' time-varying nature using a spline basis with uniformly spaced knots, i.e.,

$$\alpha_{k,l} = \sum_{r=1}^R \gamma_r s_{r,k,l} = \mathbf{\Gamma} \mathbf{s}_{k,l}. \quad (2)$$

Here, the amplitude vector $\alpha_{k,l}$ is a linear combination of the $\gamma_r \in \mathbb{R}^N$ spline basis vectors, and $s_{r,k,l}$ denotes the corresponding complex amplitude at spline point r of the l th har-

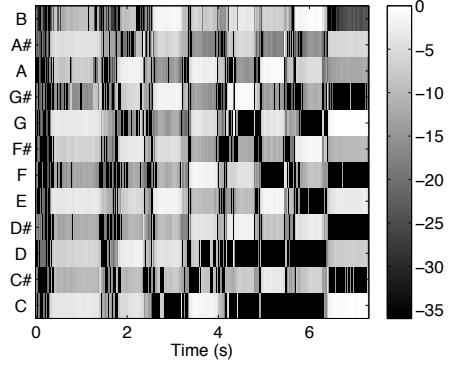


Fig. 2. The normalized log-chromagram for the trumpetscale using the method developed by Ellis.

monic for the k th source, and with

$$\alpha_{k,l} = [\alpha_{k,l}(1) \quad \alpha_{k,l}(2) \quad \cdots \quad \alpha_{k,l}(N)]^T, \quad (3)$$

$$\mathbf{s}_{k,l} = [s_{1,k,l} \quad s_{2,k,l} \quad \cdots \quad s_{R,k,l}]^T, \quad (4)$$

$$\mathbf{\Gamma} = [\gamma_1 \quad \gamma_2 \quad \cdots \quad \gamma_R], \quad (5)$$

where $[\cdot]^T$ denotes the transpose. To mould our algorithm for the use on harmonic audio signals, we, in accordance with [6], make the partition of different pitches into the twelve equivalence classes known as $C, C\#, D, D\#, E, F, F\#, G, G\#, A, A\#,$ and B . Furthermore, we design the range of f_k to have the structure $f_k = f_{base} \cdot 2^{c_k/12+o_k}$ where c_k and o_k denote the equivalence class and the octave belonging of the pitch k , respectively, and f_{base} denotes a normalized tuning parameter. The reason for this special design of the range space is that it conforms with the here examined Western music scale, which uses a cyclic scale partitioned with twelve semitones within an octave, spaced by a relative absolute frequency of $2^{1/12}$ [2]. In this work, we have chosen the tuning parameter $f_{base} = 440/2^{9/12+4}$ Hz, which corresponds to the note $C0$. Reminiscent to [6], we thus propose to extend the signal model to

$$y(t) \approx \sum_{c=0}^{11} \sum_{o=\underline{O}}^{\bar{O}} \sum_{\ell=1}^{L_{max}} \alpha_{c,o,\ell}(t) e^{i2\pi f_{base} 2^{(c/12+o)} \ell t}, \quad (6)$$

with \underline{O}, \bar{O} , and L_{max} denoting the lowest considered octave, the highest considered octave, and the maximum number of overtones, respectively. This may be expressed compactly as

$$y(t) = \sum_{c=1}^{11} \mathbf{W}_c(t) \alpha_c(t), \quad (7)$$

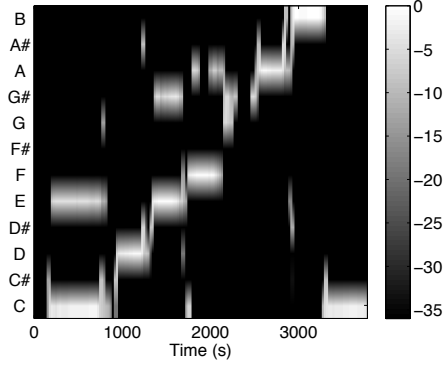


Fig. 3. The normalized log-chromagram for the trumpet scale using the CEBS method.

where

$$\begin{aligned} \mathbf{W}_c &= \begin{bmatrix} \mathbf{w}_c^O & \cdots & \mathbf{w}_c^{\bar{O}} \end{bmatrix}^T, \\ \mathbf{w}_c &= \begin{bmatrix} \mathbf{z}_c^1 & \cdots & \mathbf{z}_c^{L_{max}} \end{bmatrix}^T, \\ \mathbf{z}_c &= \begin{bmatrix} e^{i2\pi 2^{c/12} 1} & \cdots & e^{i2\pi 2^{c/12} N} \end{bmatrix}^T, \\ \boldsymbol{\alpha}_c &= \begin{bmatrix} \alpha_{c,O,1} & \cdots & \alpha_{c,\bar{O},L_{max}} & \cdots & \alpha_{c,\bar{O},L_{max}} \end{bmatrix}^T. \end{aligned}$$

Using (2), one may rewrite (7) as

$$\mathbf{y}(t) = \sum_{c=0}^{11} \text{diag}(\boldsymbol{\Gamma} \mathbf{S}_{c,o} \mathbf{W}_{c,o}^T), \quad (8)$$

where

$$\mathbf{S}_{c,o} = \begin{bmatrix} s_{c,o,1} & \cdots & s_{c,o,L_{max}} \end{bmatrix}, \quad (9)$$

$$\mathbf{s}_{c,o,l} = \begin{bmatrix} s_{1,c,o,l} & \cdots & s_{R,c,o,l} \end{bmatrix}^T. \quad (10)$$

As a result, the sought chroma features of the considered signal frame may be found as the parameters minimizing

$$\underset{S_{0,\bar{O}} \cdots S_{11,\bar{O}}}{\text{minimize}} \frac{1}{2} \left\| \mathbf{y} - \sum_{c=0}^{11} \sum_{o=Q}^{\bar{O}} \text{diag}(\boldsymbol{\Gamma} \mathbf{S}_{c,o} \mathbf{W}_{c,o}^T) \right\|_2^2, \quad (11)$$

where \mathbf{y} denotes the vector containing the measured signal. To promote a sparse solution, one may rewrite and extend (11) as

$$\begin{aligned} \underset{S_p}{\text{minimize}} \frac{1}{2} \left\| \mathbf{y} - \sum_{p=1}^P \text{diag}(\boldsymbol{\Gamma} \mathbf{S}_p \mathbf{W}_p^T) \right\|_2^2 \\ + \lambda \sum_{p=1}^P \sum_{l=1}^{L_{max}} \|\mathbf{s}_{p,l}\|_2 + \gamma \sum_{c=0}^{11} \|\tilde{\mathbf{S}}_c\|_F, \end{aligned} \quad (12)$$

where the reparametrization from c, o to p is

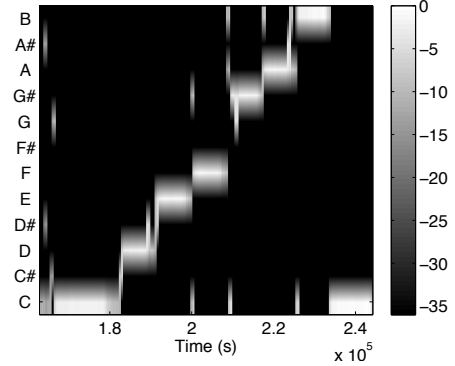


Fig. 4. The normalized log-chromagram for the trumpet scale using the CEAMS method.

$p = 12(o - Q) + c$, and thus P denotes the total number of chroma-octave pairs in the dictionary, and with

$$\tilde{\mathbf{S}}_c = \begin{bmatrix} \mathbf{S}_{c,O} & \cdots & \mathbf{S}_{c,\bar{O}} \end{bmatrix}. \quad (13)$$

The first penalty term in (12) has the effect of forcing columns in $\mathbf{s}_{p,l}$ with small l_2 norm to zero, whereas the second promotes the sparsity of the resulting chroma estimate.

3. IMPLEMENTATION

Since the problem at hand is convex, one may implement the proposed method efficiently using the Alternating Direction Method of Multipliers (ADMM) (see e.g. [15]). Denoting $\mathbf{S} = \begin{bmatrix} \mathbf{S}_1 & \cdots & \mathbf{S}_P \end{bmatrix}$, (12) may be rewritten as

$$\underset{\mathbf{X}, \mathbf{Z}}{\text{minimize}} f(\mathbf{X}) + g(\mathbf{Z}) \quad \text{subject to} \quad \mathbf{X} - \mathbf{Z} = \mathbf{0} \quad (14)$$

where

$$\begin{aligned} f(\mathbf{X}) &= \frac{1}{2} \left\| \mathbf{y} - \sum_{p=1}^P \text{diag}(\boldsymbol{\Gamma} \mathbf{X}_p \mathbf{W}_p) \right\|_2^2 \\ g(\mathbf{Z}) &= \lambda \sum_{p=1}^P \sum_{l=1}^{L_{max}} \|\mathbf{Z}_{p,l}\|_2 + \gamma \sum_{c=0}^{11} \|\mathbf{Z}_c\|_F \end{aligned} \quad (15)$$

with \mathbf{X} and \mathbf{Z} having the same structure as \mathbf{S} . It is worth noting that the ADMM separates the sought variable into two unknown variables, here denoted \mathbf{X} and \mathbf{Z} , enabling the original problem to be decomposed into easier sub-problems. These are in turn solved iteratively until convergence. Introducing the Lagrangian of (14), i.e.,

$$L_\rho(\mathbf{X}, \mathbf{Z}, \mathbf{U}) = f(\mathbf{X}) + g(\mathbf{Z}) + \frac{\rho}{2} \|\mathbf{X} - \mathbf{Z} + \mathbf{U}\|_2^2 \quad (16)$$

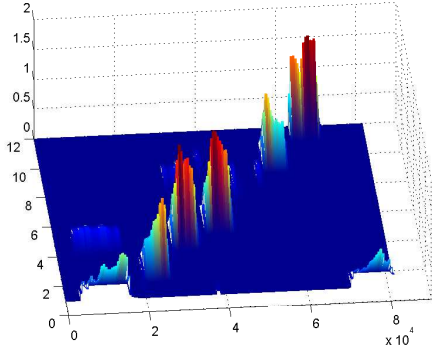


Fig. 5. The normalized 3-D log-chromagram for the trumpet scale using the CEBS method.

where \mathbf{U} represents the scaled dual variable [15], allows (16) to be solved iteratively as

$$\mathbf{X}^{(r+1)} = \arg \min_{\mathbf{X}} L_{\rho}(\mathbf{X}, \mathbf{Z}^{(r)}, \mathbf{U}^{(r)}), \quad (17)$$

$$\mathbf{Z}^{(r+1)} = \arg \min_{\mathbf{Z}} L_{\rho}(\mathbf{X}^{(r+1)}, \mathbf{Z}, \mathbf{U}^{(r)}), \quad (18)$$

$$\mathbf{U}^{(r+1)} = \mathbf{X}^{(r+1)} - \mathbf{Z}^{(r+1)} + \mathbf{U}^{(r)}. \quad (19)$$

To solve (17), one differentiates $f(\mathbf{X}) + \frac{\rho}{2} \|\mathbf{X} - \mathbf{Z} + \mathbf{U}\|_2^2$ with respect to \mathbf{X}_p and sets the result equal to zero, which yields

$$\begin{aligned} & - \sum_{n=1}^N y(n) \mathbf{\Gamma}(n, \cdot)^H \mathbf{W}_p(\cdot, n)^H + \frac{\rho}{2} (\mathbf{X}_p - \mathbf{Z}_p + \mathbf{U}_p) \\ & + \sum_{u=1}^P \sum_{n=1}^N \mathbf{\Gamma}(n, \cdot)^H \mathbf{\Gamma}(n, \cdot) \mathbf{X}_u \mathbf{W}_u(\cdot, n) \mathbf{W}_p(\cdot, n)^H = 0. \end{aligned}$$

By stacking all columns in \mathbf{X} on top of each other, this may be represented as

$$\begin{aligned} & \sum_{n=1}^N \mathbf{a}(p, n)^H y(n) + \frac{\rho}{2} (\mathbf{z}_p - \mathbf{u}_p) \\ & = \sum_{n=1}^N \sum_{u=1}^P \mathbf{a}(p, n)^H \mathbf{a}(u, n) \mathbf{x}_u + \frac{\rho}{2} \mathbf{x}_p, \end{aligned} \quad (20)$$

where

$$\mathbf{a}(u, n) = \mathbf{W}_u(\cdot, n)^T \otimes \mathbf{\Gamma}(n, \cdot), \quad (21)$$

$$\mathbf{x}_u = \text{vec}(\mathbf{X}_u), \quad (22)$$

$$\mathbf{z}_u = \text{vec}(\mathbf{Z}_u), \quad (23)$$

$$\mathbf{u}_u = \text{vec}(\mathbf{U}_u), \quad (24)$$

with \otimes denoting the Kronecker product, and $\mathbf{W}_u(\cdot, n)$ and $\mathbf{\Gamma}(n, \cdot)$ denoting the n th column in \mathbf{W}_u and the n th row $\mathbf{\Gamma}$, respectively. Let

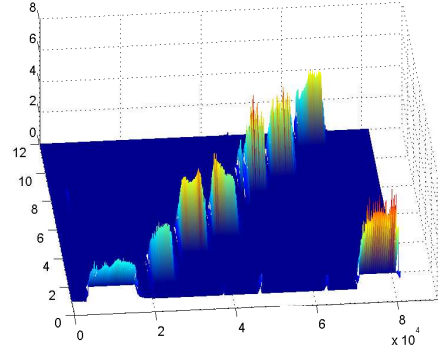


Fig. 6. The normalized 3-D log-chromagram for the trumpet scale using the proposed CEAMS method.

$$\mathbf{A}(p, u) = \sum_{n=1}^N \mathbf{a}(p, n)^H \mathbf{a}(u, n), \quad (25)$$

$$\tilde{\mathbf{y}}(p) = \sum_{n=1}^N \mathbf{a}(p, n)^H y(n), \quad (26)$$

$$\tilde{\mathbf{Y}} = [\tilde{\mathbf{y}}(1) \quad \cdots \quad \tilde{\mathbf{y}}(P)]^T, \quad (27)$$

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}(1, 1) & \cdots & \mathbf{A}(1, P) \\ \vdots & \ddots & \vdots \\ \mathbf{A}(P, 1) & \cdots & \mathbf{A}(P, P) \end{pmatrix}. \quad (28)$$

This yields the proposed algorithm, which is summarized in Algorithm (1). We term this the Chroma Estimation of Amplitude Modulated Signals (CEAMS) method. The soft thresholds \mathcal{T} and \mathcal{T} , used in Algorithm (1), are interpreted column wise, and defined as

$$\mathcal{T}(\mathbf{x}, \kappa_1) = \max \left(\frac{\mathbf{x}}{\|\mathbf{x}\|_2} (\|\mathbf{x}\|_2 - \kappa_1), 0 \right) \quad (29)$$

$$\mathcal{T}(\mathbf{X}, \kappa_2) = \left(\frac{\mathbf{X}}{\|\mathbf{X}\|_F} (\|\mathbf{X}\|_F - \kappa_2), 0 \right) \quad (30)$$

4. NUMERICAL RESULTS

The proposed method was evaluated using a concert C-scale played by a trumpet acquired from [16]. Figures 1-4 illustrate the resulting chromagrams as obtained using the estimators in [11], [10], and [6], respectively, as well as the here proposed CEAMS estimator. For the latter, we use the parameter values $\lambda = 0.3$ and $\gamma = 193$, a window length of 1024 samples, a sampling frequency of 22050 Hz, $L_{max} = 9$ overtones, and 9 spline points. As is clear from Figures 1 and 2, both the estimators in [10, 11] suffer from apparent problems in choosing the correct chroma-bin for the scale. The CEBS estimate, shown in Figure 3, is on the other hand notably cleaner, but does still suffer from some spurious chroma features. As is clear from Figure 4, these peaks are correctly

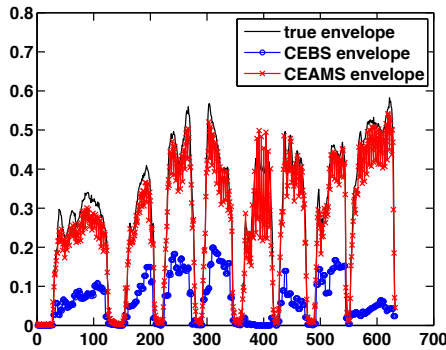


Fig. 7. The envelopes of the raw signal, the estimation using CEBS, and the estimation using the proposed CEAMS.

estimated by CEAMS. Here, we have used the same basic settings for CEBS as for CEAMS, and with $\lambda_2 = 0.05$, $\lambda_3 = 3$ and $\lambda_4 = 0.1$ (in setting these parameters, we have taken care to find the best possible setting for CEBS). Note that the G in the scale is not detected by any method. This is because the fundamental frequency found in those time frames is 808 Hz, which is slightly closer to $G\#5$ than to $G5$, using concert tuning. To illustrate the difference in time-localization between CEBS and CEAMS, Figures 5 and 6 show the 3-D chromagrams, where it once again can be noted that CEBS fails to identify the chroma-bin at $G\#$. Moreover, one notes the spurious peaks produced in CEBS, as they are of significant magnitude, compared to the rest of the chromagram. This is in contrast to CEAMS, where none of the above mentioned behaviour is present. This is also illustrated in Figure 7, showing the envelopes of the measured signal together with the CEBS and CEAMS estimates, clearly indicating the better fit of the latter.

5. REFERENCES

- [1] M. Müller, D. P. W. Ellis, A. Klapuri, and G. Richard, “Signal Processing for Music Analysis,” *IEEE J. Sel. Top. Sign. Proc.*, vol. 5, no. 6, pp. 1088–1110, 2011.
- [2] R. Shepard, “Circularity in Judgements of Relative Pitch,” *Journal of Acoustical Society of America*, vol. 36, no. 12, pp. 2346–2353, Dec. 1964.
- [3] M. Christensen and A. Jakobsson, *Multi-Pitch Estimation*, Morgan & Claypool, 2009.
- [4] A. Klapuri and M. Davy, *Signal Processing Methods for Music Transcription*, Springer, 2006.
- [5] S. I. Adalbjörnsson, A. Jakobsson, and M. G. Christensen, “Multi-Pitch Estimation Exploiting Block Sparsity,” *Elsevier Signal Processing*, vol. 109, pp. 236–247, April 2015.
- [6] T. Kronvall, M. Juhlin, S. I. Adalbjörnsson, and A. Jakobsson, “Sparse Chroma Estimation for Harmonic Audio,” in *IEEE Int. Conf. on Acoustics, Speech, and Sig. Proc.*, Brisbane, Apr. 19–24 2015.

Algorithm 1 The proposed CEAMS algorithm

- 1: Initiate $\mathbf{X} = \mathbf{X}(0)$, $\mathbf{Z} = \mathbf{Z}(0)$, and $r = 0$
 - 2: **repeat**
 - 3: $\mathbf{X}^{(r+1)} = (\mathbf{A}^H \mathbf{A} + \frac{\rho}{2} \mathbf{I})^{-1} \mathbf{A}^H \tilde{\mathbf{Y}}$
 - 4: $\mathbf{Z}^{(r+1)} = \mathcal{T}(\mathbf{T}(X_p^{(r+1)} + U_p^{(r)}), \beta/\rho), \alpha/\rho)$
 - 5: $\mathbf{U}^{(r+1)} = \mathbf{X}^{(r+1)} - \mathbf{Z}^{(r+1)} + \mathbf{U}^{(r)}$
 - 6: $r \leftarrow r + 1$
 - 7: **until** convergence
-

- [7] M. A. Bartsch and G. H. Wakefield, “Audio Thumbnailing of Popular Music Using Chroma-based Representations,” *IEEE Transactions on Multimedia*, vol. 7, no. 1, pp. 96–104, Feb. 2005.
- [8] S. Kim and S. Narayanan, “Dynamic Chroma Feature Vectors with Applications to Cover Song Identification,” in *10th IEEE Workshop on Multimedia Signal Processing*, 2008, pp. 984–987.
- [9] T.-M. Chang, E.-T. Chen, C.-B. Hsieh, and P.-C. Chang, “Cover Song Identification with Direct Chroma Feature Extraction from AAC Files,” in *IEEE 2nd Global Conference on Consumer Electronics*, Oct. 2013, pp. 55–56.
- [10] D. P. W. Ellis, “Chroma Feature Analysis and Synthesis,” <http://www.ee.columbia.edu/~dpwe/resources/matlab/chroma-ansyn/>, accessed Sept. 2014.
- [11] M. Müller and S. Ewert, “Chroma Toolbox: MATLAB Implementations for Extracting Variants of Chroma-based Audio Features,” in *Proceedings of the 12th International Conference on Music Information Retrieval (ISMIR)*, 2011.
- [12] S. I. Adalbjörnsson, A. Jakobsson, and M. G. Christensen, “Estimating Multiple Pitches Using Block Sparsity,” in *38th IEEE Int. Conf. on Acoustics, Speech, and Signal Processing*, Vancouver, May 26–31, 2013.
- [13] S. I. Adalbjörnsson, J. Swärd, T. Kronvall, and A. Jakobsson, “A Sparse Approach for Estimation of Amplitude Modulated Signals,” in *48th Annual Asilomar Conference on Signals, Systems, and Computers*, Pacific Grove, USA, Nov. 2–5 2014.
- [14] M. G. Christensen, A. Jakobsson, S. V. Andersen, and S. H. Jensen, “Amplitude Modulated Sinusoidal Signal Decomposition for Audio Coding,” *IEEE Signal Process. Lett.*, vol. 13, no. 7, pp. 389–392, July 2006.
- [15] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, “Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers,” *Found. Trends Mach. Learn.*, vol. 3, no. 1, pp. 1–122, Jan. 2011.
- [16] Mrs. Thomas, “Sound examples,” <http://www.hffmcsd.org/webpages/arushkoski/nyssma.cfm>, accessed Feb. 2015.