Context-free 2D tree structure model of musical notes for Bayesian modeling of polyphonic spectrograms 東京大学

Hirokazu Kameoka^{1,2}, Kazuki Ochiai¹, Masahiro Nakano¹, Masato Tsuchiya¹, Shigeki Sagayama¹



 $\left. \begin{array}{c} H_{\omega,1} \\ H_{\omega,2} \\ H_{\omega,3} \end{array} \right|_{H_{\omega,3}}$

 $au_2 \quad au_3 au_4^{ ext{Time}}$

 $W_{r,t}$

¹ The University of Tokyo, ² NTT Communication Science Laboratories

Abstract: This paper proposes a Bayesian model for automatic music transcription. Automatic music transcription involves several subproblems that are interdependent of each other: multiple fundamental frequency estimation, onset detection, and rhythm / tempo recognition. We formulate a unified generative model, which could be used to jointly solve the problems of determining the pitch and onset time of each musical note, the rhythm and the overall tempo variation of a piece of music, with a Bayesian approach.

1. Introducti

- Automatic music transcription consists of subproblems that are interdependent of each other:



Multipitch analysis

- Audio source separation
- Multiple *F*_o estimation
- Onset/offset estimation
- Rhythm analysis
 - Note value estimation
 - Tempo/beat estimation
- Ambiguities involved in these subproblems:
- Multipitch analysis

Spectrogram of mixture of two notes (C3 and G3)



 Since each note consists of many overtones, there can be multiple interpretations of which notes are present in a mixture.

Rhythm analysis

Possible interpretations

• Since we cannot define a note value without having a

Subprocess for generating note onset positions

 μ_2

• S_r : Onset position (in units)

Time

Time

Multipitch analysis

Rhythm analysis

- of ticks) of note r
- Music has a 2D hierarchical structure, consisting of
 - Time-spanning structure •*Synchronizing* structure

• We introduce extension of probabilistic context-free grammar (PCFG):

•

00

9

- For each node n in the parse tree: $b_n \sim \text{Bernoulli}(b_n; \phi^T)$ [Choose either "Emission" or "BINARY-PRODUCTION"] Parent If $b_n = \text{Emission}$ Ο \rightarrow $S_r \sim \delta_{S_r,S_n}, \ \ L_r \sim \delta_{L_r,L_n}$ Emit terminal symbol \rightarrow If $b_n = \mathbf{B}$ inary-production \rightarrow $\rho_n \sim \text{Bernoulli}(\rho_n; \phi^N)$ [Choose either "Synchronization" or "Time-spanning"] If $\rho_n = \mathbf{Synchronization}$ $S_{n_1} \sim \delta_{S_{n_1},S_n}, \quad S_{n_2} \sim \delta_{S_{n_2},S_n}$ Produce 2 copies $L_{n_1} \sim \delta_{L_{n_1},L_n}, \quad L_{n_2} \sim \delta_{L_{n_2},L_n}$ of note *n* If $\rho_n = \text{TIME-SPANNING}$ $S_{n_1} \sim \delta_{S_{n_1}, S_n}, \quad S_{n_2} \sim \delta_{S_{n_2}, S_n + L_{n_1}}$ $L_{n_1} \sim \delta_{L_{n_1}, L_n - L_{n_2}}$ Split note *n* into two consecutive $L_{n_2} \sim \operatorname{Discrete}(L_{n_2}; \boldsymbol{\phi}_{L_n}^{\mathrm{B}})$ notes n_1 and n_2
- [A3] Power of each musical note varies Time-span continuously in time. production [A1][A2] -> Spectrogram can be described as Synchronization → Power of note *r* Time-span production $X_{\omega,t} = \sum \underline{H_{\omega,\kappa_r}} \underline{W_{r,t}}$ at time t Spectrum of note *r* $[A3] \rightarrow$ We describe power envelope $W_{r,t}$ using a parametric function: Time-spanning binary $W_{r,t} = \sum G_{r,m,t},$ production rules Child 1 Child 2 $G_{r,m,t} = \frac{w_r u_{r,m}}{\sqrt{2\pi}\varphi} e^{-(t - (m-1)\varphi - \underline{\tau_r})^2/2\varphi^2}$ •Onset position τ_r of note *r* should be placed near the absolute time into which S_r is converted. To ensure this, we assume $\underline{\tau_r} \sim \mathcal{N}(\tau_r; \psi_{S_r}, (\sigma^{\tau})^2)$. Finally, we assume that observed Parent and children nodes spectrogram is generated according to: \mathcal{T}_{1} $Y_{\omega,t} \sim \text{Poisson}(Y_{\omega,t}; X_{\omega,t})$ Note: The maximization of the Poisson likelihood with respect to $X_{\omega,t}$ amounts to optimally fitting $X_{\omega,t}$ to $Y_{\omega,t}$ by using the I-divergence as the fitting criterion.

notion for tempo and vice versa, there can be infinite interpretations regarding what the intended rhythm was and how the tempo varied.

"Chicken-and-egg" interdependencies

Time



 Simultaneous estimation is generally preferable when several estimation problems are interdependent!

2. Motivation

- How do humans recognize music ? Humans can recognize music by combining
 - hypothesis made according to the observation about when and which notes are active, and
 - common sense about how music and music performance should be (e.g., how onset occurrences of notes should be temporally structured, how tempo varies in time)



4. Posterior Inference

Variational inference

- (--)

 $\tau = \{\tau_r\}_r$

 $\kappa = \{\kappa_r\}_r$

 $\mu = \{\mu_d\}_d$

 $S = \{S_r\}_r$

- Goal is to compute $p(\Theta|Y)$
- : total energy of note r $w = \{w_r\}_r$: shape of power envelope of note r $u = \{u_{r,m}\}_{r,m}$
 - : onset time (sec) of note r
 - : pitch index assigned to note r
- : absolute time corresponding to d ticks $\psi = \{\psi_d\}_d$: local tempo between d and d+1 ticks : onset position of note r (in ticks)
- : duration of note r (in ticks) $L = \{L_r\}_r$

To obtain exact posterior, we must derive marginal distribution p(Y), but this involves intractable integrals.

Instead, we consider finding

 $\hat{q}(\Theta) = \operatorname{argmin} \operatorname{KL}(q(\Theta) || p(\Theta | Y))$ $q \in Q$

Kullback-Leibler divergence between $q(\Theta)$ and $p(\Theta|Y)$

- Tractable choices of Q allows us to obtain principled approximate solutions
- Mean-field approximation

5. Music Transcription Experime

Experiment 1

- Test data: piano recording (RWC-MDB-C-2001 No. 26) Morzart: Piano Sonata No. 11 in A major, K. 331/300i
- Monaural
- sampling rate=16kHz
- Constant-Q transform was used to obtain the observed spectrogram.
- Note onsets au_1, \ldots, au_R were given manually.
- Result: (b) and (c)
- Experiment 2
- Next, we applied our method to the same data without providing any information about au.
- Result: (d) and (e)

The result showed that many octave



(b) Detected beat locations along with the estimate of W_{rt}





(d) Detected beat locations along with the estimate of $W_{r,t}$



We propose to model the generative process of an entire spectrogram of a piece of music by formulating the following three sub-processes and combining them into one process:

(1) subprocess by which a tempo curve is generated, 2 subprocess by which a set of note onset positions (in terms of the relative time) is generated, and 3 subprocess by which spectrogram is generated

Parameter inference given the observation would then give a musically likely interpretation of what is actually happening in the spectrogram (i.e., a musical score).

 $q(\Theta) = q(C)q(w)q(u)q(\tau,\psi,\mu)q(S,L)$

Coordinate ascent

• By iteratively optimizing one factor q(.) in the mean-field approximation of the posterior at a time while fixing all the other factors, we can find a local optimum.

Probability distributions

 $x \sim \delta_{x,y}$ means x = y (with probability 1) $\mathcal{N}(x;\mu,\sigma) \propto e^{-rac{(x-\mu)^2}{2\sigma^2}}$ Bernoulli $(x; y) = y^x (1 - y)^{1 - x}$ Beta $(y; z) \propto y^{z_1 - 1} (1 - y)^{z_2 - 1}$ where $x \in \{0,1\}, 0 \le y \le 1$ and $z = (z_1, z_2)$ Discrete $(x; y) = y_x$ Dirichlet $(y; z) \propto \prod y_i^{z_i - 1}$ where $y = (y_1, ..., y_I)$ with $y_1 + \cdots + y_I = 1$ and $z = (z_1, ..., z_I)$ $Poisson(y;x) = x^y e^{-x}/y!$

errors had occurred. This kind of error often occurs when there is a mismatch between a spectrum model and an actual spectrum. The validity of the assumptions we have made about the spectra of musical sounds must be carefully examined in the future.



References

[1] N. Bertin, R. Badeau, and G. Richard, "Blind signal decompositions for automatic transcription of polyphonic music: NMF and K-SVD on the benchmark," in *Proc. ICASSP*, Vol. 1, pp. 65-68, 2007. [2] K. Ochiai, H. Kameoka, and S. Sagayama, "Explicit beat structure modeling for non-negative matrix factorization-based multipitch analysis," in *Proc. ICASSP*, pp. 133-136, 2012. [3] A. T. Cemgil, "Bayesian inference in non-negative matrix factorisation models," Technical Report CUED/F-INFENG/TR.609, University of Cambridge, 2008. [4] M. D. Hoffman, D. M. Blei, and P. R. Cook, "Bayesian nonparametric matrix factorization for recorded music," in Proc. ICML, pp. 439-446, 2010. [5] K. Yoshii, and M. Goto, "A nonparametric Bayesian multipitch analyzer based on infinite latent harmonic allocation," IEEE Trans. Audio, Speech, Language Process., Vol. 20, No. 3, pp. 717-730, 2012. [6] M. Nakano, Y. Ohishi, H. Kameoka, R. Mukai, and K. Kashino, "Bayesian nonparametric music parser," in *Proc. ICASSP*, pp. 461-464, 2012. [7] P. Liang, S. Petrov, M. I. Jordan, and D. Klein, "The infinite PCFG using hierarchical Dirichlet processes," in Proc. EMNLP, pp. 688-697, 2007. [8] H. Kameoka, T. Nishimoto, and S. Sagayama, "A multipitch analyzer based on harmonic temporal structured clustering," IEEE Trans. on Audio, Speech, Language Process., Vol. 15, No. 3, pp. 982-994, 2007.