

A Stochastic Model of Singing Voice F_0 Contours for Characterizing Expressive Dynamic Components

Yasunori Ohishi[†], Hirokazu Kameoka[†], Daichi Mochihashi[‡], Kunio Kashino[†]

[†]NTT Communication Science Laboratories, NTT Corporation,

[‡]The Institute of Statistical Mathematics

{ohishi.yasunori,kameoka.hirokazu,kashino.kunio}@lab.ntt.co.jp, daichi@ism.ac.jp

Abstract

We present a novel stochastic model of singing voice fundamental frequency (F_0) contours for characterizing expressive dynamic components, such as *vibrato* and *portamento*. Although dynamic components can be important features for any singing voice applications, modeling and extracting these components from a raw F_0 contour have yet to be accomplished. Therefore, we describe a process for generating dynamic components explicitly and represent the process as a stochastic model. Then we develop an algorithm for estimating the model parameters based on statistical techniques. Experimental results show that our method successfully extracts the expressive components from raw F_0 contours.

Index Terms: Singing voice, Fundamental frequency, Second-order linear system, Stochastic model

1. Introduction

The fundamental frequency (F_0) contour in a singing voice contains two main types of dynamic components. One is those generated by the physical constraints of the vocal folds, such as *overshoot*, *preparation* and *fine fluctuations* [1, 2]. The other type is dynamic components generated by singer's musical expressive intentions, such as *vibrato* and *portamento* [3]. Most previous papers have reported that these dynamic components strongly affect singing-voice perception, and that the former relates to the naturalness and individuality of a singing voice while the latter relates to singing styles and skills [4]. Accordingly, extracting these components from a raw F_0 contour automatically can be potentially very beneficial for any singing voice applications, such as singer identification, singing skill evaluation and the synthesis of more natural and varied singing voices [5, 6, 7].

Previous studies have represented the dynamic components generated by the physical constraints as the output of a second-order linear system [8, 9]. The input is a stepwise signal representing a musical-note sequence. The transfer function of the system is described as

$$\mathcal{H}(s) = \frac{\Omega^2}{s^2 + 2\zeta\Omega s + \Omega^2}, \quad (1)$$

where ζ and Ω denote the damping ratio and the natural angular frequency, respectively. The dynamic components were controlled by adjusting these parameters

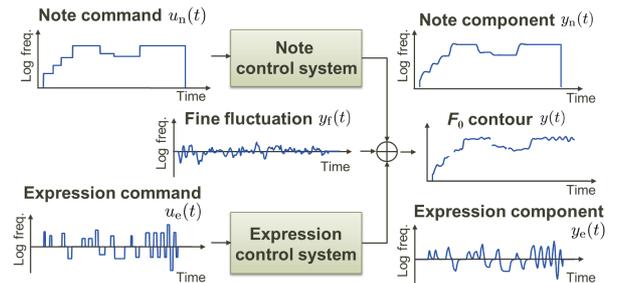


Figure 1: Process of generating a singing voice F_0 contour

manually [9]. In [10], we proposed a method for solving the inverse problem of estimating the parameters from a raw F_0 contour. However, modeling and extracting the dynamic components generated by expressive intentions from the raw F_0 contour have not yet been accomplished. For example, although the rate and amplitude of *vibrato* are time-varying, *vibrato* has been conventionally modeled as a sinusoidal signal [9]. We consider that *vibrato* is controlled variedly by the singer's expressive intentions.

In this paper we propose a new stochastic model of singing voice F_0 contours to describe the process of generating various dynamic components explicitly. This model is based on an analogy with the Fujisaki model [11], which describes the process of generating speech F_0 contours, and assumes that an F_0 contour on a logarithmic scale, $y(t)$, where t is time, is the superposition of three components (Fig. 1). The note and expression components are the outputs of second-order linear systems driven by the note and expression commands that correspond to the musical note sequence and the musical expressive intentions, respectively. The note component contains the note transition and *overshoot*, and the expression component contains *vibrato* and *portamento*. The fine fluctuation component consists of an irregular fluctuation higher than 10 Hz [9]. Then, we formulate a discrete-time stochastic version of this process and develop an algorithm for estimating the model parameters based on statistical techniques. Experimental results show that our method successfully extracts the expressive intentions for *vibrato* and *portamento* from a raw F_0 contour. Furthermore, we verify that the extracted expression components are perceived as *vibrato* and *portamento* through a psychoacoustic experiment.

describe the parameters as $\{\varphi^{(i)}, \psi^{(i)}\}_{i=1}^I$. Hence, it follows from Eq. (5) that

$$\mathbf{y}_n \sim \mathcal{N}(\mathbf{A}^{-1}\boldsymbol{\mu}_n, \sigma_n^2 \mathbf{A}^{-1}(\mathbf{A}^{-1})^T), \quad (6)$$

$$\mathbf{y}_e \sim \mathcal{N}(\mathbf{B}^{-1}\boldsymbol{\mu}_e, \sigma_e^2 \mathbf{B}^{-1}(\mathbf{B}^{-1})^T). \quad (7)$$

As for the fine fluctuation component $\mathbf{y}_f := (y_f[1], \dots, y_f[K])^T$, we assume that it is white Gaussian noise such that $\mathbf{y}_f \sim \mathcal{N}(\mathbf{0}, \sigma_f^2 \mathbf{I}_K)$.

Overall, the likelihood function of the model parameters Θ given \mathbf{y} can be written as

$$P(\mathbf{y}|\Theta) = \frac{|\boldsymbol{\Sigma}^{-1}|^{1/2}}{(2\pi)^{K/2}} \exp\left\{-\frac{1}{2}(\mathbf{y} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{y} - \boldsymbol{\mu})\right\},$$

$$\boldsymbol{\mu} = \mathbf{A}^{-1}\boldsymbol{\mu}_n + \mathbf{B}^{-1}\boldsymbol{\mu}_e, \quad (8)$$

$$\boldsymbol{\Sigma} = \sigma_n^2 \mathbf{A}^{-1}(\mathbf{A}^{-1})^T + \sigma_e^2 \mathbf{B}^{-1}(\mathbf{B}^{-1})^T + \sigma_f^2 \mathbf{I}_K,$$

where $\Theta := \{\theta_n, \{\varphi^{(i)}, \psi^{(i)}\}_{i=1}^I, \xi, \sigma_f^2\}$. As for the prior probability of Θ , we assume that the parameters are independent of each other and distributed according to the following probability density functions: $P(s_1, \dots, s_K) = P(s_1) \prod_{k=2}^K P(s_k | s_{k-1})$, $B_e^{(i,j)} \sim \mathcal{N}(\mu_{B^{(i,j)}}, \sigma_B^2)$, $\varphi^{(i)} \sim \mathcal{N}(\mu_\varphi, \sigma_\varphi^2)$, $\psi^{(i)} \sim \mathcal{N}(\mu_\psi, \sigma_\psi^2)$, $\xi \sim \mathcal{N}(\mu_\xi, \sigma_\xi^2)$, $d_i \sim \mathcal{N}(0, \sigma_d^2)$. σ_n^2, σ_e^2 and σ_f^2 are uniformly distributed.

4. Parameter estimation algorithm

We describe an iterative algorithm, which locally maximizes the posterior density $P(\Theta|\mathbf{y}) \propto P(\mathbf{y}|\Theta)P(\Theta)$. By regarding a set consisting of the note, expression and fine fluctuation components, $\mathbf{x} := (\mathbf{y}_n^T, \mathbf{y}_e^T, \mathbf{y}_f^T)^T$, as the complete data, this problem can be viewed as an incomplete data problem, which can be dealt with using the EM algorithm. Taking the conditional expectation of the log-likelihood with respect to \mathbf{x} given \mathbf{y} and $\Theta = \Theta'$, and then adding $\log P(\Theta)$, we obtain the Q function

$$Q(\Theta, \Theta') \stackrel{c}{=} \frac{1}{2} [\log |\boldsymbol{\Lambda}^{-1}| - \text{tr}(\boldsymbol{\Lambda}^{-1} \mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta'])]$$

$$+ 2m^T \boldsymbol{\Lambda}^{-1} \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] - m^T \boldsymbol{\Lambda}^{-1} \mathbf{m} + \log P(\Theta),$$

$$\mathbf{m} := \begin{bmatrix} \mathbf{A}^{-1}\boldsymbol{\mu}_n \\ \mathbf{B}^{-1}\boldsymbol{\mu}_e \\ \mathbf{0} \end{bmatrix}, \boldsymbol{\Lambda}^{-1} := \begin{bmatrix} \mathbf{A}^T \mathbf{A} / \sigma_n^2 & \mathbf{O} & \mathbf{O} \\ \mathbf{O} & \mathbf{B}^T \mathbf{B} / \sigma_e^2 & \mathbf{O} \\ \mathbf{O} & \mathbf{O} & \mathbf{I}_K / \sigma_f^2 \end{bmatrix}.$$

Because the relationship between the incomplete data and the complete data can be written as $\mathbf{y} = \mathbf{H}\mathbf{x}$ where $\mathbf{H} := [\mathbf{I}_K, \mathbf{I}_K, \mathbf{I}_K]$, $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']$ and $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']$ are given explicitly as

$$\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] = \mathbf{m} + \boldsymbol{\Lambda} \mathbf{H}^T (\mathbf{H} \boldsymbol{\Lambda} \mathbf{H}^T)^{-1} (\mathbf{y} - \mathbf{H} \mathbf{m}),$$

$$\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta'] = \boldsymbol{\Lambda} - \boldsymbol{\Lambda} \mathbf{H}^T (\mathbf{H} \boldsymbol{\Lambda} \mathbf{H}^T)^{-1} \boldsymbol{\Lambda} \boldsymbol{\Lambda}$$

$$+ \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta'] \mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']^T.$$

The E-step computes $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']$ and $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']$. In the M-step, we maximize $Q(\Theta, \Theta')$ with respect to each parameter of Θ in which we treat $\mathbb{E}[\mathbf{x} | \mathbf{y}; \Theta']$ and $\mathbb{E}[\mathbf{x}\mathbf{x}^T | \mathbf{y}; \Theta']$ as constants. In particular, the state sequence $\{s_k\}_{k=1}^K$ can be solved efficiently using the

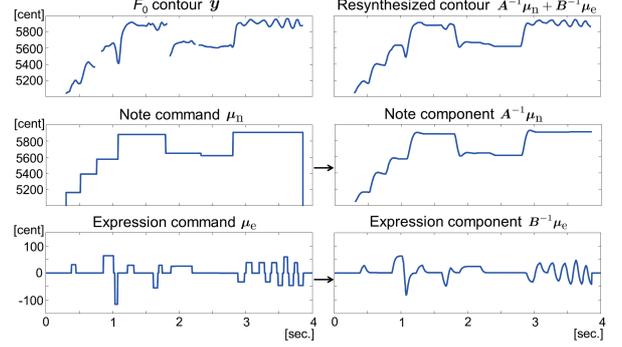


Figure 3: Estimation result from a raw F_0 contour

Viterbi algorithm [12]. Owing to space limitations, we omit all the mathematical details.

Looking back at Eq. (8), we have thus far implicitly assumed that we are given a set of F_0 observations on the whole sample period. However, the F_0 data in unvoiced regions are missing. This can simply be viewed as a missing data imputation problem, which can be effectively dealt with using EM algorithm [12].

5. Experimental evaluations

We tested our method in two experiments. Experiment A evaluated the expression commands estimated from raw F_0 contours, experiment B evaluated singing voices synthesized using resynthesized contours subjectively.

In these experiments, we used an F_0 contour annotated manually in the AIST annotation [13]. The title of the song is "PROLOGUE" (RWC-MDB-P-2001 [14], Song No.07). Although the F_0 contour should essentially be estimated from the acoustic signal, we used these data to evaluate the upper limit of the performance of our method. The F_0 values were annotated every 5 ms ($t_0 = 5$ ms) and were represented in cents so that one equal-tempered semitone corresponded to 100 cents. We used the MIDI data to determine I and $\{A_n^{(i)}\}_{i=1}^I$. The initial conditions are shown in the footnote¹.

Experiment A: Fig. 3 shows note command $\boldsymbol{\mu}_n$, expression command $\boldsymbol{\mu}_e$, note component $\mathbf{A}^{-1}\boldsymbol{\mu}_n$ and expression component $\mathbf{B}^{-1}\boldsymbol{\mu}_e$ estimated from a raw F_0 contour. Resynthesized contour $\mathbf{A}^{-1}\boldsymbol{\mu}_n + \mathbf{B}^{-1}\boldsymbol{\mu}_e$ is represented as the sum of two components, and is similar to the F_0 contour. Since σ_f is 13.6 cent, we confirm that \mathbf{y}_f can be estimated as a fine fluctuation component.

In the F_0 contour of Fig. 4(a), *portamento*, which is a vocal slide between two notes, is observed. The result shows that two rectangular pulses for changing pitch gradually are estimated in expression command $\boldsymbol{\mu}_e$. On

¹The iteration was run for 1000 iterations. J was set at 5. The state transition probabilities were set at $\phi_{S_{i-1}, S_{i-1}} = \log(0.9999 \times (J - 1) / J)$, $\phi_{S_{i-1}, S_{i,j}} = \log(0.9999 / J)$, $\phi_{S_{i-1}, S_{i+1,1}} = \log(0.0001)$, $\phi_{S_{i,j}, S_{i,j}} = \log(0.9999)$, $\phi_{S_{i,j}, S_{i,1}} = \log(0.0001)$, with $1 \leq i \leq I$ and $2 \leq j \leq J$. The parameters of the prior distributions were set at $\mu_\varphi = 6$, $\sigma_\varphi^2 = 0.1$, $\mu_\psi = 0.6$, $\sigma_\psi^2 = 0.02$, $\mu_\xi = 3$, $\sigma_\xi^2 = 0.1$, $\sigma_d^2 = 2500$, $\sigma_B^2 = 100$, $\mu_{B^{(i,1)}} = 0$, $\mu_{B^{(i,2)}} = 30$, $\mu_{B^{(i,3)}} = -30$, $\mu_{B^{(i,4)}} = 60$, $\mu_{B^{(i,5)}} = -60$, with $1 \leq i \leq I$.

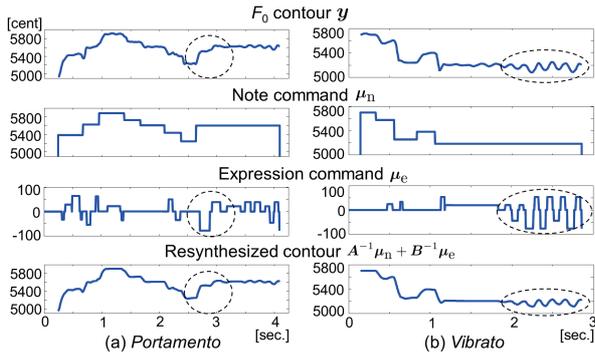


Figure 4: Estimation of commands for portamento and vibrato

the other hand, in the F_0 contour in Fig. 4(b), *vibrato* is observed in the latter part of the note. The result shows that the expression command consisted of quasi-periodic rectangular pulses. We found that our method extracted these commands from the raw F_0 contours and obtained resynthesized contours similar to the F_0 contours.

Experiment B: We conducted a psychoacoustic experiment to evaluate the estimated expression components subjectively. First, we prepared the following contours: (1) Note command μ_n , (2) Note component $A^{-1}\mu_n$, (3) Resynthesized contour $A^{-1}\mu_n + B^{-1}\mu_e$, (4) F_0 contour y . Then, we synthesized four kinds of singing voices using singing synthesis software based on Yamaha’s Vocaloid3 technology [15]. Specifically, notes and lyrics were manually entered using the estimated note command (all the syllables of the lyrics were /na/), *pitch bend (PIT)* and *pitch bend sensitivity (PBS)* were adjusted using the difference between the note command and each contour. Finally, we divided synthesized singing voice signals into short signals consisting of four bars and used these signals as stimuli. The total number of stimuli was 96 signals (24 signals for each contour).

Scheffe’s method of paired comparison was used to evaluate the expressiveness of the stimuli. As shown at the top of Fig. 5, the subjects decided which stimulus was a more expressive singing voice according to an eleven-grade evaluation measure. The expressiveness of a singing voice has a multidimensional meaning. However, in this experiment, we defined an expressive singing voice as a singing voice with vocal ornamentations such as *vibrato* and *portamento*. The pair-wise stimuli were presented through binaural headphones at a comfortable sound pressure level. Each paired stimulus was randomly presented to each subject. The subjects were five males.

Fig. 5 shows the results. The evaluation values are normalized for each subject and are averaged over all subjects. Subjects perceived singing voices based on (3) and (4) contours to be more expressive than those based on (1) and (2) contours. Additionally, it was difficult for subjects to determine the difference between singing voices based on (3) and (4) contours. Therefore, we found that the singer’s expressiveness is represented as the expression components estimated from the F_0 contour. Future work is to compare the expression components with *vibrato* modeled by a sinusoidal signal subjectively.

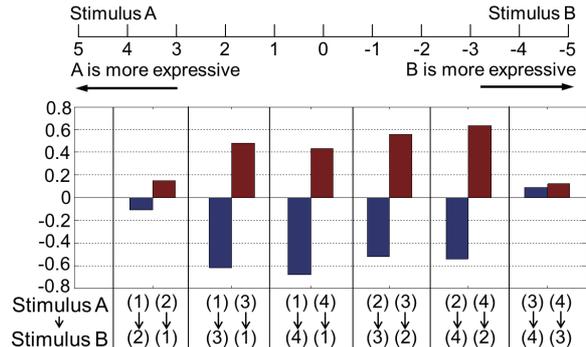


Figure 5: Subjective evaluation results: (1) Note command, (2) Note component, (3) Resynthesized contour, and (4) F_0 contour

6. Conclusions

We proposed a stochastic model of singing voice F_0 contours for characterizing and extracting various dynamic components and develop a model parameter estimation algorithm based on the EM approach. Experimental results show that our method decomposes a raw F_0 contour into the note and expression components explicitly and that the singer’s expressiveness is contained in the expression component. In the future, we plan to evaluate model validity using the large data sets and learn the singing styles using the expression commands.

7. References

- [1] J. Sundberg, *The Science of the Singing*, Northern Illinois University Press, 1987.
- [2] G de Krom and G Bloothoof, “Timing and accuracy of fundamental frequency changes in singing,” in *Proc. ICPhS95*.
- [3] C. E. Seashore, “A musical ornament, the vibrato,” in *Psychology of Music*, 1938, pp. 33–52.
- [4] T. Saitou et al., “Acoustic and perceptual effects of vocal training in amateur male singing,” in *Proc. EUROSPEECH 2009*.
- [5] W. Tsai et al., “An automated singing evaluation method for karaoke systems,” in *Proc. ICASSP 2011*.
- [6] K. Oura et al., “Pitch adaptive training for HMM-based singing voice synthesis,” in *Proc. ICASSP 2012*.
- [7] L. Regnier et al., “Singer verification: singer model .vs. song model,” in *Proc. ICASSP 2012*.
- [8] N. Minematsu et al., “Prosodic modeling of Nagauta singing and its evaluation,” in *Proc. SpeechProsody 2004*.
- [9] T. Saitou et al., “Development of an F0 control model based on F0 dynamic characteristics for singing-voice synthesis,” *Speech Communication*, vol. 46, pp. 405–417, 2005.
- [10] Y. Ohishi et al., “Statistical modeling of F0 dynamics in singing voices based on Gaussian processes with multiple oscillation bases,” in *Proc. INTERSPEECH 2010*.
- [11] H. Fujisaki, “A note on the physiological and physical basis for the phrase and accent components in the voice fundamental frequency contour,” in *Vocal Physiology: Voice Production, Mechanisms and Functions*, pp. 347–355. Raven Press, 1988.
- [12] H. Kameoka et al., “A statistical model of speech F0 contours,” in *Proc. SAPA 2010*.
- [13] M. Goto, “AIST annotation for the RWC music database,” in *Proc. ISMIR 2006*.
- [14] M. Goto et al., “RWC music database: Popular, classical, and jazz music databases,” in *Proc. ISMIR 2002*.
- [15] H. Kenmochi et al., “Singing synthesis as a new musical instrument,” in *Proc. ICASSP 2012*.