

変調スペクトルを考慮した HMM 音声合成*

© 高道 慎之介, 戸田 智基, Graham Neubig, Sakriani Sakti, 中村 哲 (奈良先端大)

1 はじめに

HMM 音声合成 [1] において, 生成されるパラメータ系列の過剰な平滑化は, 音質劣化の一因となる. これに対して本稿では, パラメータ系列の変調スペクトル (MS: Modulation Spectrum) に基づくポストフィルタを提案する. 提案法では, 生成パラメータ系列の MS が, 自然音声のパラメータ系列の MS と近くなるように, 生成パラメータ系列に対してフィルタ処理を施す. 実験的評価によりその有効性を示す.

2 HMM 音声合成のパラメータ生成法

2.1 HMM 尤度最大化基準 [2]

HMM 音声合成では, 自然音声のパラメータ系列からコンテキスト依存 HMM を学習する. 生成時には, 合成対象のテキストに対応する文 HMM を形成し, 静的・動的特徴量間の明示的な制約条件の下で HMM 尤度を最大化することで, パラメータ系列を生成する.

$$\hat{\mathbf{c}} = \underset{\mathbf{c}}{\operatorname{argmax}} P(\mathbf{W}\mathbf{c}|\lambda) \quad (1)$$

ただし, $\mathbf{c} = [\mathbf{c}_1^\top, \dots, \mathbf{c}_T^\top]^\top$ は T フレームの音声パラメータ系列, $\mathbf{c}_t = [c_t(1), \dots, c_t(D)]^\top$ は時刻 t における D 次元の音声パラメータ, \mathbf{W} は動的特徴量の計算に用いる重み係数によって決定される行列 [2], λ は HMM のパラメータセットを表す.

式 (1) により生成されるパラメータは, 過剰に平滑化される傾向にある.

2.2 HMM 尤度・GV 尤度最大化基準 [3]

系列内変動 (GV: Global Variance) は, パラメータ系列全体の変動成分を表し, 次式で定義される.

$$\mathbf{v}(\mathbf{c}) = [v(1), \dots, v(D)]^\top \quad (2)$$

$$v(d) = \frac{1}{T} \sum_{t=1}^T \left(c_t(d) - \frac{1}{T} \sum_{\tau=1}^T c_\tau(d) \right)^2 \quad (3)$$

生成時には, HMM 尤度及び GV 尤度を最大化してパラメータ系列を生成する.

$$\hat{\mathbf{c}} = \underset{\mathbf{c}}{\operatorname{argmax}} P(\mathbf{W}\mathbf{c}|\lambda) P(\mathbf{v}(\mathbf{c})|\lambda_v)^w \quad (4)$$

ただし, λ_v は GV の確率密度関数のパラメータセット, w は GV 尤度の重みを表す.

生成パラメータ系列の GV は, GV 尤度を考慮することで補償され, 合成音声の音質が改善する.

3 変調スペクトルを考慮した HMM 音声合成

3.1 パラメータ系列の変調スペクトル分析

MS は, 本来, パラメータ系列をフーリエ変換した値を表す [4] が, 本稿では, その対数振幅スペクトルを MS と呼ぶ. パラメータ系列 \mathbf{c} に対する MS を次式で定義する.

$$\mathbf{s}(\mathbf{c}) = [\mathbf{s}(1)^\top, \dots, \mathbf{s}(D)^\top]^\top \quad (5)$$

$$\mathbf{s}(d) = [s_d(0), \dots, s_d(M)]^\top \quad (6)$$

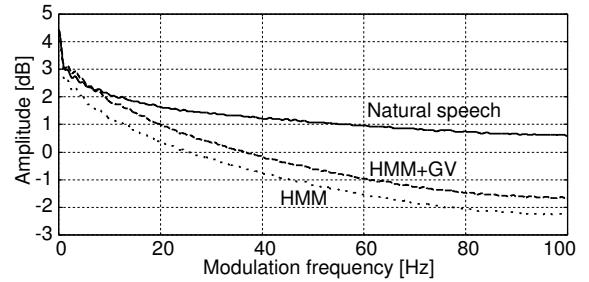


Fig. 1 メルケプストラム系列の変調スペクトル

ただし, $s_d(m)$ は, d 次元目のパラメータ系列 $[c_1(d), \dots, c_T(d)]^\top$ に対する, 周波数インデックス m の MS, M は離散フーリエ変換のサンプル数の半分を表す. 本稿では発話毎に MS を計算する.

Fig. 1 に, 式 (1) (HMM) と式 (4) (HMM+GV) で生成された第 10 次メルケプストラム系列の MS を示す. 比較のため, 自然音声 (Natural speech) の同系列の MS も示す. “HMM” の MS は, 自然音声のパラメータ系列の MS と比較して, 大きく減衰していることが確認できる. また, “HMM+GV” の MS は, GV の導入により比較的補償されるものの, 未だに大きく減衰している. 故に, MS の直接的な補償により, 合成音声の音質改善がもたらされると期待される.

3.2 変調スペクトルに基づくポストフィルタ

MS を補償するポストフィルタ処理を提案する. ポストフィルタは, 学習データを用いて事前に設計する.

3.2.1 学習部

自然音声のパラメータ系列から, 次式に示す確率密度関数を学習する.

$$P(\mathbf{s}(\mathbf{c})|\lambda_s) = \mathcal{N}(\mathbf{s}(\mathbf{c}); \boldsymbol{\mu}^{(N)}, \boldsymbol{\Sigma}^{(N)}) \quad (7)$$

ただし, $\mathcal{N}(\cdot; \boldsymbol{\mu}^{(N)}, \boldsymbol{\Sigma}^{(N)})$ は平均 $\boldsymbol{\mu}^{(N)} = [\mu_{1,0}^{(N)}, \dots, \mu_{D,M}^{(N)}]^\top$ と対角共分散行列 $\boldsymbol{\Sigma}^{(N)} = \operatorname{diag} \left[\left(\sigma_{1,0}^{(N)} \right)^2, \dots, \left(\sigma_{D,M}^{(N)} \right)^2 \right]$ の正規分布, $\mu_{d,m}^{(N)}$ と

$\left(\sigma_{d,m}^{(N)} \right)^2$ は $s_d(m)$ の平均と分散, λ_s は MS の確率密度関数のパラメータセットを表す. 同様に, HMM 音声合成で生成されたパラメータ系列から正規分布 $\mathcal{N}(\cdot; \boldsymbol{\mu}^{(G)}, \boldsymbol{\Sigma}^{(G)})$ を学習する. なお, 自然音声のパラメータ系列と生成パラメータ系列間の継続長の違いが MS に影響することを避けるために, 正規分布の学習に用いる生成パラメータ系列は, 自然音声の継続長において生成する.

3.2.2 生成部

生成されたパラメータ系列 \mathbf{c} の MS に対して次式のポストフィルタを適用する.

$$s'_d(m) = (1-k)s_d(m) + k \left[\frac{\sigma_{d,m}^{(N)}}{\sigma_{d,m}^{(G)}} \left(s_d(m) - \mu_{d,m}^{(G)} \right) + \mu_{d,m}^{(N)} \right] \quad (8)$$

*HMM-based speech synthesis considering modulation spectrum, by TAKAMICHI, Shinnosuke, TODA, Tomoki, SAKTI, Sakriani, NEUBIG, Graham, NAKAMURA, Satoshi (NAIST)

ただし、 k はポストフィルタ強度係数 ($0 \leq k \leq 1$) を表す。フィルタ後のMSは、 $k = 1$ の際には自然音声のパラメータ系列のMSに近い値となり、 $k = 0$ の際にはポストフィルタ処理前と等価となる。ポストフィルタ後のパラメータ系列は、式(8)で計算されたMSと、フィルタ処理前のパラメータ系列の周波数位相特性から計算する。

4 実験的評価

4.1 実験条件

学習データは女性話者によるATR音素バランス文[5] A-Iセット450文とする。評価データは同Jセット53文を使用する。学習データのサンプリング周波数は16 kHz、フレームシフトは5 msとする。スペクトル特徴量は、STRAIGHT分析[6]による0次から24次のメルケプストラム係数、音源特徴量は、対数 F_0 、5周波数帯域における平均非周期成分を使用する。5状態left-to-right型の隠れセミマルコフモデル (Hidden Semi-Markov Model: HSMM) の学習を行う。変調スペクトルにおける離散フーリエ変換のサンプル数は4096点とする。これは、学習・評価データのパラメータ系列のフレーム数を十分に超える値である。提案法はスペクトルパラメータに対してのみ適用し、音源パラメータは式(1)で生成する。

以下に示す手法を用いて評価を行う。

HMM: 式(1)で生成

HMM+MS: 式(1)で生成したパラメータ系列に対して提案法を適用

HMM+GV: 式(4)で生成

HMM+GV+MS: 式(4)で生成したパラメータ系列に対して提案法を適用

まず、ポストフィルタ強度係数を決定するための評価を行う。ポストフィルタ強度係数を0から1まで0.05刻みで変化させ、ポストフィルタ処理後のパラメータ系列に対するHMM尤度、GV尤度及びMS尤度を計算する。同時に、自然音声 (Natural speech) のパラメータ系列に対する尤度も計算する。

次に、提案法による音質改善効果を対比較実験により評価する。評価は8人の受聴者に対するプリファレンススコアとする。

4.2 客観評価結果

ポストフィルタ強度係数を変化させた時の、ポストフィルタ後のパラメータ系列に対するHMM対数尤度をFig. 2に、GV対数尤度をFig. 3に、MS対数尤度をFig. 4に示す。Fig. 2から、ポストフィルタ強度係数を大きくするに従い、生成パラメータ系列に対するHMM尤度は大きく減少することがわかる。しかしながら、その尤度は自然音声のパラメータ系列に対するHMM尤度よりも依然として大きい。一方、Fig. 3から、ポストフィルタ強度係数を大きくするに従いGV尤度は変化し、ポストフィルタ強度係数を0.85に設定した場合に、“HMM+MS”と“HMM+GV+MS”の両方の尤度が自然音声に接近していることがわかる。対してFig. 4から、生成パラメータ系列に対するMS尤度は、自然音声のパラメータ系列に対する尤度よりも常に小さいことがわかる。以上の結果から、ポストフィルタ強度係数を0.85に設定する。

4.3 主観評価結果

音質の主観評価結果をFig. 5に示す。“HMM”における生成パラメータ系列に対して提案法を適用することで、スコアが著しく上昇し、“HMM+GV”と同等の音質が得られることが分かる。また、“HMM+GV”

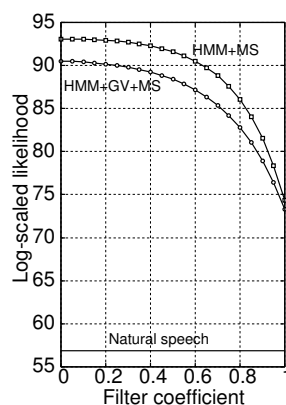


Fig. 2 HMM 尤度

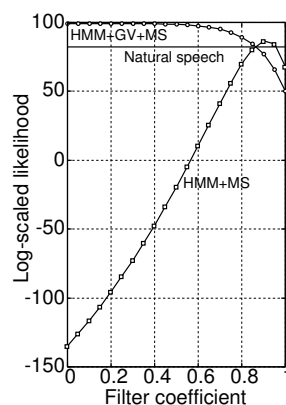


Fig. 3 GV 尤度

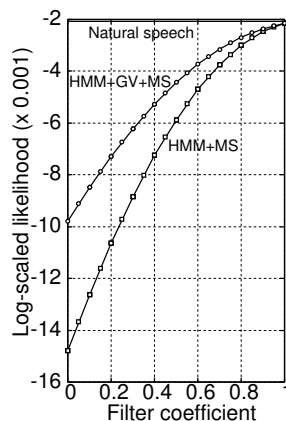


Fig. 4 MS 尤度

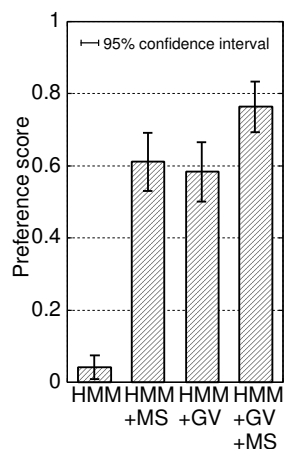


Fig. 5 主観評価結果

におけるパラメータ系列に対する提案法の適用により、スコアは更に上昇することがわかる。以上の結果から、提案法による音質の改善が確認できる。

5 まとめ

本稿では、HMM音声合成の音質改善を目的として、従来のHMM音声合成の生成パラメータ系列に対して、変調スペクトルを補償するポストフィルタ処理を提案し、その有効性を実験的評価により示した。今後は、変調スペクトルを考慮したパラメータ生成法の検討を行う。

謝辞 本研究の一部は、JSPS 科研費 22680016 の助成を受け実施した。また、本研究の一部は、(独)情報通信研究機構の委託研究「知識・言語グリッドに基づくアジア医療交流支援システムの研究開発」の一環として実施した。

参考文献

- [1] H. Zen *et al.*, *Speech Commun.*, 51(11), pp. 1039–1064, 2009.
- [2] K. Tokuda *et al.*, *Proc. ICASSP*, Vol. 3, pp. 1315–1318, 2000.
- [3] T. Toda *et al.*, *IEICE Transactions*, Vol. E90–D, No. 5, pp. 816–824, 2007.
- [4] L. Atlas *et al.*, *EURASIP Journal*, Vol. 7, pp. 668–675, 2003.
- [5] 阿部 他, ATR テクニカルレポート, TR-I-0166, 1990.
- [6] H. Kawahara *et al.*, *Speech Communication*, Vol. 27, No. 3–4, pp. 187–207, 1999.