

# 統計的パラメトリック音声合成における 変調スペクトル制約付きトラジェクトリ学習\*

○ 高道 慎之介 (奈良先端大/CMU), 戸田 智基 (奈良先端大),  
Alan W. Black (CMU), 中村 哲 (奈良先端大)

## 1 はじめに

統計的パラメトリック音声合成の音質改善には、音声パラメータ生成部における変調スペクトル (MS: Modulation Spectrum) 補償が有効である [1, 2]. しかしながら、生成部における補償処理は、短遅延合成処理を必要とするシステムに不向きである. そこで本稿では、HMM 音声合成 [3] と GMM 声質変換 [4] において、静的・動的特徴量の制約下で系列をモデル化するトラジェクトリモデル [5] を、変調スペクトル制約下で学習する手法を提案する. 実験的評価により、従来の短遅延合成処理能力を維持しつつ、音質を改善できることを示す.

## 2 従来の枠組み

HMM 音声合成では入力テキストのコンテキストと出力音声のパラメータ系列から、GMM 声質変換では入力・出力音声のパラメータから、各統計モデル  $\lambda$  を学習する. 生成時には、入力コンテキスト又は入力音声特徴量系列  $\mathbf{X}$  から決定される次式の目的関数  $L(\mathbf{y})$  を、静的・動的特徴量間の制約条件下で最大化することで、パラメータ系列  $\hat{\mathbf{y}}_{\hat{q}}$  を生成する.

$$L(\mathbf{y}) = P(\mathbf{W}\mathbf{y}|\mathbf{X}, \hat{q}, \lambda) = \mathcal{N}(\mathbf{W}\mathbf{y}; \mathbf{E}_{\hat{q}}, \mathbf{D}_{\hat{q}}) \quad (1)$$

$$\hat{\mathbf{y}}_{\hat{q}} = \mathbf{R}_{\hat{q}}^{-1} \mathbf{r}_{\hat{q}} = (\mathbf{W}^T \mathbf{D}_{\hat{q}}^{-1} \mathbf{W})^{-1} \mathbf{W}^T \mathbf{D}_{\hat{q}}^{-1} \mathbf{E}_{\hat{q}} \quad (2)$$

ただし、 $\mathbf{y} = [\mathbf{y}_1^T, \dots, \mathbf{y}_T^T]^T$  は  $T$  フレームの音声パラメータ系列、 $\mathbf{y}_t = [y_t(1), \dots, y_t(D)]^T$  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 [6] である. 準最適な HMM の状態系列又は GMM の分布系列  $\hat{q} = [\hat{q}_1, \dots, \hat{q}_T]$  を用いた近似により、目的関数は、平均ベクトル  $\mathbf{E}_{\hat{q}} = [\mathbf{E}_{\hat{q}_1,1}, \dots, \mathbf{E}_{\hat{q}_T,T}]$  と精度行列  $\mathbf{D}_{\hat{q}}^{-1} = \text{diag}[\Sigma_{\hat{q}_1}^{-1}, \dots, \Sigma_{\hat{q}_T}^{-1}]$  から成る正規分布  $\mathcal{N}(\cdot; \mathbf{E}_{\hat{q}}, \mathbf{D}_{\hat{q}})$  で与えられる.  $q_t$  は時刻  $t$  における HMM 状態もしくは GMM 分布インデックス  $q$  であり、 $1 \leq q \leq Q$  とする.  $Q$  は HMM 状態もしくは GMM 分布の総数である.  $\mu_{q,t} = \{\mu_q \text{ (HMM)}, \mathbf{A}_q \mathbf{X}_t + \mathbf{b}_q \text{ (GMM)}\}$  と  $\Sigma_q^{-1}$  はそれぞれ、時刻  $t$  における  $q$  に対応する平均ベクトル及び精度行列である.  $\mu_q$ ,  $\mathbf{A}_q$ ,  $\mathbf{b}_q$  はそれぞれ、 $q$  に対応する平均ベクトル、変換行列、及び、バイアスペクトルであり、 $\mathbf{X}_t$  は、 $\mathbf{X}$  における時刻  $t$  のベクトルである. コレスキー分解 [6] や短遅延生成アルゴリズム [7] により、 $\hat{\mathbf{y}}_{\hat{q}}$  は解析的に生成される.

## 3 提案する学習法

### 3.1 トラジェクトリモデル

トラジェクトリ HMM の学習部 [5] では、次式の目的関数  $L_{\text{trj}}$  を最大化するように  $\lambda$  を推定する.

$$L_{\text{trj}} = P(\mathbf{y}|\mathbf{X}, \hat{q}, \lambda) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}_{\hat{q}}, \mathbf{R}_{\hat{q}}^{-1}) \quad (3)$$

一方で GMM 声質変換において、結合確率密度関数のトラジェクトリ GMM が提案されている [8] が、本

稿では式 (3) と同様に、条件付き確率密度関数としてトラジェクトリ GMM を定式化する.

### 3.2 変調スペクトル (MS)

MS は、系列内変動 (GV: Global Variance) [9, 4] の拡張であり、パラメータ系列のパワースペクトルとして定義される [1]. パラメータ系列  $\mathbf{y}$  の MS  $\mathbf{s}(\mathbf{y}) = [\mathbf{s}(1), \dots, \mathbf{s}(D)]^T$  は次式で表される.

$$\mathbf{s}(d) = [s_d(0), \dots, s_d(D'_s - 1)]^T \quad (4)$$

$$s_d(f) = R_{d,f}^2 + I_{d,f}^2 \quad (5)$$

$$= \left( \sum_{t=1}^T y_t(d) \cos kt \right)^2 + \left( \sum_{t=1}^T y_t(d) \sin kt \right)^2 \quad (6)$$

ただし、 $2D'_s$  は離散フーリエ変換 (DFT) のタップ長、 $k = -\pi f/D'_s$  は変調周波数を表す.  $D'_s$  は MS の次元数を表し、 $D'_s \leq D_s$  である. 本稿では、音声知覚に強く影響する低変調周波数帯域の MS のみを取り扱う. MS の確率密度関数を、平均ベクトル  $\mu_s$  と共分散行列  $\Sigma_s$  から成る正規分布  $\mathcal{N}(\mathbf{s}(\mathbf{y}); \mu_s, \Sigma_s)$  で表し、更に、 $\Sigma_s^{-1} = [\mathbf{p}_s^{(1)}, \dots, \mathbf{p}_s^{(D)}]$  とする.  $\mathbf{p}_s^{(d)}$  は、 $\mathbf{s}(d)$  の精度に対応する行列である.

### 3.3 変調スペクトル制約付きトラジェクトリ学習

提案法では、次式の目的関数  $L_{\text{ms}}$  を最大化するように  $\lambda$  を推定する.

$$L_{\text{ms}} = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}_{\hat{q}}, \mathbf{R}_{\hat{q}}^{-1}) \mathcal{N}(\mathbf{s}(\mathbf{y}); \mathbf{s}(\hat{\mathbf{y}}_{\hat{q}}), \Sigma_s) \frac{\omega_s^T}{D'_s} \quad (7)$$

$\omega_s$  は MS 尤度の重みであり、本稿では 1.0 に設定する. ここで、 $\xi_A = [\mathbf{A}_1^T, \dots, \mathbf{A}_Q^T]^T$ ,  $\xi_b = [\mathbf{b}_1^T, \dots, \mathbf{b}_Q^T]^T$ ,  $\mu = [\mu_1^T, \dots, \mu_Q^T]^T$ ,  $\Sigma^{-1} = [\Sigma_1^{-1}, \dots, \Sigma_Q^{-1}]^T$  を、各々全ての  $\mathbf{A}_q$ ,  $\mathbf{b}_q$ ,  $\mu_q$ ,  $\Sigma_q^{-1}$  を結合したパラメータとする.  $\lambda$  は、最急降下法  $\xi_A^{(i+1)} = \xi_A^{(i)} + \alpha (\partial \log L_{\text{ms}} / \partial \xi_A) |_{\xi_A = \xi_A^{(i)}}$  にて反復的に推定される.  $i$  は反復インデックス、 $\alpha$  は学習係数である. 各結合パラメータに関する一次微分は次式で与えられる.

$$\frac{\partial \log L_{\text{ms}}}{\partial \xi_A} = \mathbf{S}_{\hat{q}}^T \text{diag}^{-1}[\mathbf{D}_{\hat{q}}^{-1} \mathbf{W} (\mathbf{y} - \hat{\mathbf{y}}_{\hat{q}}) \mathbf{X}^T + \frac{\omega_s^T}{D'_s} \mathbf{D}_{\hat{q}}^{-1} \mathbf{W} \mathbf{R}_{\hat{q}}^{-1} \mathbf{s}_{\hat{q}} \mathbf{X}^T] \quad (8)$$

$$\frac{\partial \log L_{\text{ms}}}{\partial \xi_b} = \frac{\partial \log L_{\text{ms}}}{\partial \mu} = \mathbf{S}_{\hat{q}}^T (\mathbf{D}_{\hat{q}}^{-1} \mathbf{W} (\mathbf{y} - \hat{\mathbf{y}}_{\hat{q}}) + \frac{\omega_s^T}{D'_s} \mathbf{D}_{\hat{q}}^{-1} \mathbf{W} \mathbf{R}_{\hat{q}}^{-1} \mathbf{s}_{\hat{q}}) \quad (9)$$

$$\frac{\partial \log L_{\text{ms}}}{\partial \Sigma^{-1}} = \frac{1}{2} \mathbf{S}_{\hat{q}}^T \text{diag}^{-1}[\mathbf{W} (\mathbf{R}_{\hat{q}}^{-1} + \hat{\mathbf{y}}_{\hat{q}} \hat{\mathbf{y}}_{\hat{q}}^T - \mathbf{y} \mathbf{y}^T) - \mathbf{E}_{\hat{q}} (\hat{\mathbf{y}}_{\hat{q}} - \mathbf{y})^T \mathbf{W}^T - \mathbf{W} (\hat{\mathbf{y}}_{\hat{q}} - \mathbf{y}) \mathbf{E}_{\hat{q}}^T + \frac{2\omega_s^T}{D'_s} \mathbf{W} \mathbf{R}_{\hat{q}}^{-1} \mathbf{s}_{\hat{q}} (\mathbf{E}_{\hat{q}} - \mathbf{W} \hat{\mathbf{y}}_{\hat{q}})] \quad (10)$$

\* Modulation Spectrum-Constrained Trajectory Training Algorithm in Statistical Parametric Speech Synthesis, by TAKAMICHI, Shinnosuke (NAIST/CMU), TODA, Tomoki (NAIST), Black Alan (CMU), NAKAMURA, Satoshi (NAIST)

Table 1 実験条件

	HMM 音声合成	GMM 声質変換
学習データ	593 文	50 文
評価データ	100 文	100 文
話者	男性	男性及び女性
音響モデル	5 状態 HSM	64 混合 GMM
提案法	スペクトル, $F_0$	スペクトル

ただし,  $\mathbf{S}_{\hat{q}}$  は  $\hat{q}$  を表現する行列である [5].  $\mathbf{s}_{\hat{q}}$  は次式で与えられる.

$$\mathbf{s}_{\hat{q}} = [\mathbf{s}'_1, \dots, \mathbf{s}'_T]^\top, \quad (11)$$

$$\mathbf{s}'_t = [s_t(1), \dots, s_t(D)]^\top, \quad (12)$$

$$s_t(d) = 2\mathbf{f}_t(d) \mathbf{p}_s^{(d)\top} (\mathbf{s}(\mathbf{y}) - \mathbf{s}(\hat{\mathbf{y}})), \quad (13)$$

$$\mathbf{f}_t(d) = [f_{t,d}(0), \dots, f_{t,d}(D'_s - 1)]^\top, \quad (14)$$

$$f_{t,d}(f) = \hat{R}_{d,f} \cos kt + \hat{I}_{d,f} \sin kt, \quad (15)$$

ただし,  $\hat{R}_{d,f}$  と  $\hat{I}_{d,f}$  は  $\hat{\mathbf{y}}_{\hat{q}}$  から計算される. 本稿では, 通常のトラジェクトリ学習による結果を初期値として, 最急降下法による反復更新を行う. HMM 音声合成における  $F_0$  モデリングには連続  $F_0$  モデル [10] を使用し, [1] と同様に系列平均を 0 とした MS を計算する. GMM 声質変換においては, 結合 GMM パラメータの内, 分布重みと, 入力特徴量の平均ベクトル及び共分散行列については, 更新しない.

### 3.4 考察

学習部における MS 制約により生成パラメータ  $\hat{\mathbf{y}}_{\hat{q}}$  の変調スペクトルは補償されるため, 生成部における補償は不要である. 故に, 式 (2) の生成処理を使用しつつ, 音質改善が可能である. また, MS は GV を内包するため, 提案法は暗黙的に GV も補償する.

## 4 実験的評価

### 4.1 実験条件

CMU ARCTIC 音声データベース [11] から学習データと評価データを選択する. 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 [12] 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 [13] を使用する.  $\omega_s$  は 1.0 に設定する.  $D_s$  は 2048 とする.  $D'_s$  は, スペクトルにおいて  $D_s/2$  (= 50Hz) [2],  $F_0$  において  $D_s/10$  (= 10Hz) [1] に設定する. 詳細な実験条件は, 表 1 参照とする.

以下の学習法を比較する. パラメータ生成法には式 (2) を使用する.

**BASIC**: 通常の学習 [3, 4]

**TRJ**: トラジェクトリ学習 [5]<sup>1</sup>

**GV**: GV 制約付き (トラジェクトリ) 学習 [14, 15]<sup>2</sup>

**MS**: MS 制約付きトラジェクトリ学習 (提案法)

評価として, HMM 音声合成と GMM 声質変換における音質に関するプリファレンス AB テスト, 及び, GMM 声質変換における話者性に関する XAB テストを実施する. XAB テストのリファレンス音声は, ターゲット話者の分析再合成音である. 被験者は, 各評価に対して 6 人である. 非周期成分における提案法の知覚的影響は小さいため, 非周期成分の学習には “BASIC” を使用する.

<sup>1</sup>トラジェクトリ HMM は [5] による定式化, トラジェクトリ GMM は本稿による定式化である.

<sup>2</sup>[15] はトラジェクトリ GMM を使用していない.

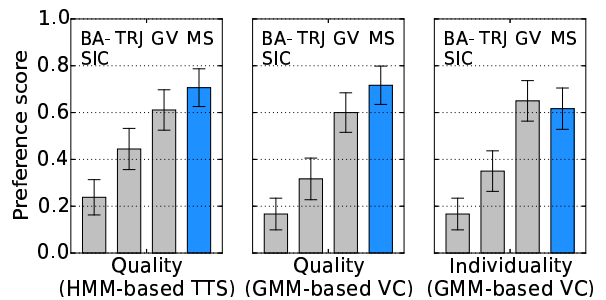


Fig. 1 主観評価結果 (エラーバーは 95% 信頼区間)

### 4.2 実験結果

各主観評価結果を Fig. 1 に示す. HMM 音声合成と GMM 声質変換の音質評価において提案法 (“MS”) のスコアが最も高い. 一方で, “MS” の話者性のスコアは, “TRJ” より高いものの “GV” と同程度のスコアとなっている. これらの結果より, 提案する変調スペクトル制約付き学習により, 解析的なパラメータ生成処理を保持しつつ, 変調スペクトル補償による音質改善効果が得られることが明らかになった.

## 5 まとめ

本稿では, 統計的パラメトリック音声合成において, 単遅延合成処理能力を維持しつつ音質を改善させる手法として, 変調スペクトル制約付きトラジェクトリ学習法を提案し, 実験的評価により有効性を確認した. 今後は, 任意話者音声合成への応用を検討する.

謝辞 本研究の一部は, JSPS 特別研究員奨励費 26・10354, JSPS 科研費 26280060, 及び, 頭脳循環を加速する若手研究者戦略的海外派遣プログラムの助成を受け実施した.

## 参考文献

- [1] S. Takamichi et al., *Proc. ICASSP*, pp. 290–294, 2014.
- [2] 高道 他, 音講論 (春), 2-2-2, 2015.
- [3] K. Tokuda et al., *Proc. IEEE*, Vol. 111, No. 5, pp. 1234–1252, 2013.
- [4] T. Toda et al., *IEEE Trans. on Audio, Speech and Language Processing*, Vol. 15, No. 8, pp. 2222–2235, 2007.
- [5] H. Zen et al., *Computer Speech and Language*, Vol. 21, No. 1, pp. 153–173, 2007.
- [6] K. Tokuda et al., *Proc. ICASSP*, pp. 1315–1318, 2000.
- [7] T. Muramatsu et al., *Proc. INTERSPEECH*, pp. 1076–1079, 2008.
- [8] H. Zen et al., *IEEE Trans.*, Vol. 19, pp. 417–430, 2011.
- [9] T. Toda et al., *IEICE on Trans. Inf. and Syst.*, Vol. E90-D, No. 5, pp. 816–824, 2007.
- [10] K. Yu et al., *IEEE Trans. on Audio, Speech and Language*, Vol. 19, No. 5, pp. 1071–1079, 2011.
- [11] J. Kominek et al., *Tech. Rep. CMU-LTI-03-177*, 2003.
- [12] H. Kawahara et al., *Speech Comm.*, Vol. 27, No. 3–4, pp. 187–207, 1999.
- [13] 大谷 他, 信学論, Vol. J91-D, No. 4, pp. 1082–1091, 2008.
- [14] T. Toda et al., *Proc. ICASSP*, pp. 4025–4028, 2009.
- [15] H. Hwang et al., *Proc. APSIPA ASC*, pp.1–6, 2014.