

話者空間のテンソル表現に基づく任意話者声質変換

齋藤 大輔[†] 山本 敬介[†] 峯松 信明[†] 広瀬 啓吉[†]

[†] 東京大学大学院情報理工学系研究科 〒113-0033 東京都文京区本郷 7-3-1

E-mail: dsaito@hil.t.u-tokyo.ac.jp, yama@nii.ac.jp, {mine,hirose}@gavo.t.u-tokyo.ac.jp

あらまし 本稿では、話者空間をテンソル形式によって表現することにより、柔軟に話者性を制御することが可能となる新しい手法を提案する。声質変換の研究において、任意話者の音声を入力または出力として、変換を実現する手法はアプリケーション応用の観点からも非常に重要な技術であるといえる。任意話者声質変換を目的とする技術として、固有声混合正規分布モデル (EV-GMM) に基づく固有声変換法 (EVC) が提案されている。EVC においては、話者認識でよく用いられるアプローチと同様に、各話者 GMM の正規分布の平均ベクトルを連結して得られる GMM スーパーベクトルをもとに話者空間が構築される。構築された話者空間上において、個々の話者は固有スーパーベクトルに対する少数の重みパラメータによって表現することが可能となる。本稿では、話者空間を構築するための事前学習話者データに対して、テンソル解析を導入することによって話者空間を構築することを検討する。本研究における提案手法では、個々の話者はスーパーベクトルではなく行列によって表現される。この話者を表す行列の行及び列は、それぞれ音響特徴量の平均ベクトルの次元及びガウス分布の要素に対応する。ここで、これらの行列のセットに対してテンソル解析を導入することで話者空間が構築される。提案法は、話者情報のスーパーベクトル表現に内在する問題点に対する解法となっており、任意話者声質変換の性能向上が期待できる。本稿では、一対多声質変換において、提案する話者空間表現を導入することで、その有効性を示す。

キーワード 声質変換, 混合正規分布モデル, 固有声, テンソル解析, Tucker 分解

Voice conversion from/to arbitrary speakers based on tensor representation of speaker space

Daisuke SAITO[†], Keisuke YAMAMOTO[†], Nobuaki MINEMATSU[†], and Keikichi HIROSE[†]

[†] Graduate School of Information Science and Technology, The University of Tokyo
7-3-1 Hongo, Bunkyo-ku, Tokyo 113-0033, Japan

E-mail: dsaito@hil.t.u-tokyo.ac.jp, yama@nii.ac.jp, {mine,hirose}@gavo.t.u-tokyo.ac.jp

Abstract This paper describes a novel approach to flexible control of speaker characteristics using tensor representation of speaker space. In voice conversion studies, realization of conversion from/to an arbitrary speaker's voice is one of the important objectives. For this purpose, eigenvoice conversion (EVC) based on an eigenvoice Gaussian mixture model (EV-GMM) was proposed. In the EVC, similarly to speaker recognition approaches, a speaker space is constructed based on GMM supervectors which are high-dimensional vectors derived by concatenating the mean vectors of each of the speaker GMMs. In the speaker space, each speaker is represented by a small number of weight parameters of eigen-supervectors. In this paper, we revisit construction of the speaker space by introducing the tensor analysis of training data set. In our approach, each speaker is represented as a matrix of which the row and the column respectively correspond to the Gaussian component and the dimension of the mean vector, and the speaker space is derived by the tensor analysis of the set of the matrices. Our approach can solve an inherent problem of supervector representation, and it improves the performance of voice conversion. Experimental results of one-to-many voice conversion demonstrate the effectiveness of the proposed approach.

Key words Voice conversion, Gaussian mixture model, eigenvoice, tensor analysis, Tucker decomposition

1. はじめに

声質変換は、入出力の対応関係を記述する変換モデルに基づいて、任意の文に対して入力音声の声質を所望の声質へ変換する技術である。声質変換は、広義には異なる二つの特徴量空間のマッピング技術と考えられ、テキスト音声合成における話者性の制御をはじめとして [1], 雑音環境下音声の音声強調や身体運動から音声への変換など多岐にわたる応用が検討されている [2], [3]。入出力の対応関係を記述する変換モデルの構築に関しては、多くの手法が提案されている。そのなかでも統計的変換手法は盛んに研究されており、コードブックマッピング法をはじめ [4], ニューラルネットワークを用いた手法 [5] や混合正規分布モデル (GMM) に基づく変換法 [1], [6] など数多く提案されている。その中でも GMM に基づく変換法はその柔軟性から近年主流となっている。

しかし、変換モデルの構築の際には、基本的に同一発話内容の入出力音声対からなるパラレルデータを用いる必要がある。また、変換モデルの利用は学習時の入出力話者対に限定される。すなわち話者性を柔軟に制御することは声質変換における重要な課題といえる。そのためには他の話者の音声データを事前知識として用いることが有効と考えられ、上記のようなパラレルデータを必要としない手法もいくつか提案されている [7], [8]。これらの手法は、結合確率密度分布に含まれるパラメータを非パラレルデータを用いて適応するアプローチとなっている。一方、多数の話者の音声データをより効果的に用いる手法として、固有声変換法 (Eigenvoice conversion; EVC) が提案されている [9]。EVC では、参照話者と多数の事前収録話者との間の複数のパラレルデータを用いて、固有声に基づく混合正規分布モデル (EV-GMM) を構築する。複数のパラレルデータより得られた結合確率密度分布から、事前収録話者の話者 GMM をそれぞれ抽出し、ガウス分布の平均ベクトルを連結した GMM スーパーベクトルを用いて話者空間を構築する。話者認識の場合と同様に、任意の話者はこの話者空間の一点で表され、基底に対する少数の重みパラメータを推定することで、話者性を柔軟に制御することができる。

しかし、GMM スーパーベクトルによる話者空間表現は、複数要因からの音響的な変動を一つの特徴量空間に含んでいる。すなわち、GMM のガウス分布の要素と平均ベクトルの次元が混在した高次の特徴量空間となっている。本研究では、話者性を柔軟に制御するという観点から、話者空間をテンソルで表現した一対多声質変換を提案する。提案法では、任意の話者はスーパーベクトルではなく、行および列がそれぞれ GMM の要素と平均ベクトルに対応するような行列の形で表現される。このような話者表現を用いることで事前収録話者のデータセットが 3 階のテンソルで表現でき、テンソル解析を導入することで話者空間を構築することができる。テンソル解析は複数要因からの変動を適切に扱うことが可能であり [10], 提案する話者空間表現による性能の向上が期待できる。また、本稿では一対多声質変換にテンソルを用いた話者空間表現を導入するが、提案法は多対一声質変換や話者認識にも応用することができる。ま

た、提案法は主に話者空間表現に着眼しているため、話者正規化学習を用いた EVC [11] や非パラレルデータを用いた学習 [12] との統合も可能である。本稿では提案法の定式化について述べ、従来の EVC と比較して提案法の有効性を示す。

以下、2. において、従来の固有声変換法について概説する。3. で、提案するテンソル表現に基づく話者空間表現について述べる。4. において提案法の実験的評価について述べ、最後に 5. で、本稿の結論と今後の課題について述べる。

2. 固有声変換法 (EVC)

2.1 固有声に基づく混合正規分布モデル

本章では、一対多 EVC について概説する。今、参照話者の音響特徴量を $\mathbf{X}_t = [\mathbf{x}_t^\top, \Delta \mathbf{x}_t^\top]^\top$, s 番目の事前収録話者の音響特徴量を $\mathbf{Y}_t^{(s)} = [\mathbf{y}_t^{(s)\top}, \Delta \mathbf{y}_t^{(s)\top}]^\top$ と表す。ただし \top は転置を表す。ここで、音響特徴量は D 次元の静的および動的特徴量を結合した $2D$ 次元の音響特徴量となる。参照話者と事前収録話者の結合確率密度は、EV-GMM として以下のようにモデル化される。

$$P(\mathbf{X}_t, \mathbf{Y}_t^{(s)} | \boldsymbol{\lambda}^{(EV)}, \mathbf{w}^{(s)}) = \sum_{m=1}^M \alpha_m \mathcal{N}([\mathbf{X}_t^\top, \mathbf{Y}_t^{(s)\top}]^\top; \boldsymbol{\mu}_m^{(Z)}(\mathbf{w}^{(s)}), \boldsymbol{\Sigma}_m^{(Z)}) \quad (1)$$

$$\boldsymbol{\mu}_m^{(Z)}(\mathbf{w}^{(s)}) = \begin{bmatrix} \boldsymbol{\mu}_m^{(X)} \\ \mathbf{B}_m \mathbf{w}^{(s)} + \mathbf{b}_m^{(0)} \end{bmatrix}, \boldsymbol{\Sigma}_m^{(Z)} = \begin{bmatrix} \boldsymbol{\Sigma}_m^{(XX)} & \boldsymbol{\Sigma}_m^{(XY)} \\ \boldsymbol{\Sigma}_m^{(YX)} & \boldsymbol{\Sigma}_m^{(YY)} \end{bmatrix} \quad (2)$$

ここで $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ は、平均ベクトルを $\boldsymbol{\mu}$, 分散共分散行列を $\boldsymbol{\Sigma}$ とする正規分布を表す。 m 番目の要素の重みは α_m で表し、混合数を M とする。EV-GMM では、 S 人の事前収録話者を利用して、出力話者の平均ベクトル $\boldsymbol{\mu}_m^{(Y)}$ をバイアスベクトル $\mathbf{b}_m^{(0)}$ と $K (< S)$ 個の表現ベクトルの線形結合で表す。このとき、出力話者の話者性は K 次元の重みベクトル $\mathbf{w}^{(s)}$ で制御できる。すなわち話者空間が K 個の基底スーパーベクトル $\mathbf{B} = [\mathbf{B}_1^\top, \mathbf{B}_2^\top, \dots, \mathbf{B}_m^\top]^\top \in \mathcal{R}^{2DM \times K}$ とバイアススーパーベクトル $\mathbf{b} = [\mathbf{b}_1^{(0)\top}, \mathbf{b}_2^{(0)\top}, \dots, \mathbf{b}_m^{(0)\top}]^\top \in \mathcal{R}^{2DM \times 1}$ によって構築される。

2.2 EVC における話者空間構築

主成分分析 (PCA) に基づいて、EV-GMM を構築する場合、最初に出力話者非依存の GMM (TI-GMM) を、全ての参照話者と事前収録話者とのパラレルデータを用いて学習する。次に、対応するパラレルデータを用いて TI-GMM の出力話者の平均ベクトルを更新することで、話者依存のモデルを得る。話者空間の特徴量ベクトルとして、事前収録話者の GMM の各要素の平均ベクトルを連結し、スーパーベクトルを生成する。得られたスーパーベクトルを用いて PCA を行うことで、バイアスベクトル \mathbf{b} と表現ベクトル \mathbf{B} を得ることができる。

2.3 EV-GMM の教師なし適応

任意の話者に対する EV-GMM は、出力話者の音声データを用いて、最尤基準に基づいて重みベクトル \mathbf{w} を推定することで適応できる [9]。今、出力話者の音響特徴量系列を $\mathbf{Y}^{(tar)}$ とすると、 \mathbf{w} は以下のように推定できる。

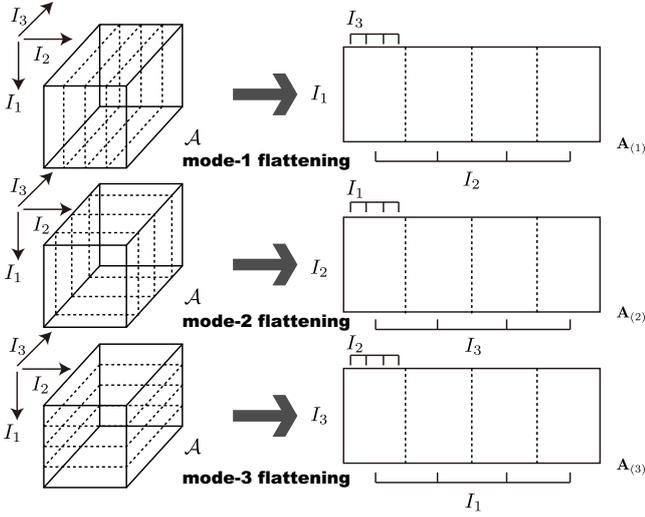


図1 $I_1 \times I_2 \times I_3$ テンソル \mathcal{A} の行列 $\mathbf{A}_{(1)}$, $\mathbf{A}_{(2)}$, $\mathbf{A}_{(3)}$ への平坦化
Fig.1 Flattening of the $(I_1 \times I_2 \times I_3)$ -tensor \mathcal{A} to the flattened matrices $\mathbf{A}_{(1)}$, $\mathbf{A}_{(2)}$ and $\mathbf{A}_{(3)}$.

$$\begin{aligned} \hat{\mathbf{w}} &= \operatorname{argmax}_{\mathbf{w}} \int P(\mathbf{X}, \mathbf{Y}^{(tar)} | \boldsymbol{\lambda}^{(EV)}, \mathbf{w}) d\mathbf{X} \\ &= \operatorname{argmax}_{\mathbf{w}} P(\mathbf{Y}^{(tar)} | \boldsymbol{\lambda}^{(EV)}, \mathbf{w}) \end{aligned} \quad (3)$$

ここで、出力の確率密度分布は GMM で表される。よって以下の補助関数を導入し、EM アルゴリズムを用いることで、重みを繰り返し最適化していく。

$$Q(\mathbf{w}, \hat{\mathbf{w}}) = \sum_m P(m | \mathbf{Y}^{(tar)}, \boldsymbol{\lambda}^{(EV)}, \mathbf{w}) \log P(\mathbf{Y}^{(tar)} | m | \boldsymbol{\lambda}^{(EV)}, \hat{\mathbf{w}}) \quad (4)$$

式 (4) から、 $\hat{\mathbf{w}}$ に関する以下の更新式を得る。

$$\hat{\mathbf{w}} = \left\{ \sum_{m=1}^M \bar{\gamma}_m^{(tar)} \mathbf{B}_m^\top \boldsymbol{\Sigma}_m^{(YY)^{-1}} \mathbf{B}_m \right\}^{-1} \sum_{m=1}^M \mathbf{B}_m^\top \boldsymbol{\Sigma}_m^{(YY)^{-1}} \bar{\mathbf{Y}}_m^{(tar)} \quad (5)$$

$$\bar{\gamma}_m^{(tar)} = \sum_{t=1}^T \gamma_{m,t}, \quad \bar{\mathbf{Y}}_m^{(tar)} = \sum_{t=1}^T \gamma_{m,t} (\mathbf{Y}_t^{(tar)} - \mathbf{b}_m^{(0)}) \quad (6)$$

$$\gamma_{m,t} = P(m | \mathbf{Y}_t^{(tar)}, \boldsymbol{\lambda}^{(EV)}, \mathbf{w}) \quad (7)$$

式 (5) はおよそ話者空間の基底ベクトルへの射影重みを推定していることに相当する。なお、式 (7) の初期化に際しては、TI-GMM を用いる。適応後のパラメータ生成については [13] と同様である。

EVC における重みパラメータの推定は、出力話者の発話内容を知る必要がないため、完全な教師なし適応となる。推定パラメータ数が少ないため、極少量の適応データを用いて変換モデルを構築することが可能である。しかし、表現ベクトルが構築する空間は高次元の GMM スーパーベクトルに基づいているため、複数の音響的変動要因が内在している。以下では、各話者の表現にスーパーベクトルではなく、行列形式の表現を用いる。これにより、テンソル解析を用いて複数の変動要因を適切に扱う手法を検討する。

3. 話者空間のテンソル表現

3.1 多重線形解析

本章では、提案するテンソル解析に基づく話者空間の構築について述べる。まず、提案法に関連する多重線形解析について説明する [14]。テンソルは、行列表現を一般化した多次元配列表現である。テンソルにおける個々のインデックスはモードと呼ばれる。今、 $\mathcal{A} \in \mathcal{R}^{I_1 \times I_2 \times I_3}$ を 3 階のテンソルとする。一般に高階のテンソルは Fig. 1 に示すようなモード n の平坦化操作によって行列の形で表現できる。モード n の平坦化では、モード n のインデックスに沿うようにテンソルをスライスし、行列の形に連結する。この平坦化操作を用いることで、テンソルと行列の間の積を定義できる。モード n 積は $\mathcal{A} = \mathbf{G} \times_n \mathbf{B}$ と表現し、モード n の平坦化行列を用いて $\mathbf{A}_{(n)} = \mathbf{B} \cdot \mathbf{G}_{(n)}$ のように演算できる。

行列代数において最も重要な演算として特異値分解 (SVD) がある。ここで、行列は 2 階のテンソルと見なすことができ、 \mathcal{A} の特異値分解は、モード n 積を用いて以下のように表される。

$$\mathcal{A} = \mathbf{U} \mathbf{S} \mathbf{V}^\top = \mathbf{S} \times_1 \mathbf{U} \times_2 \mathbf{V} \quad (8)$$

この SVD の表現を高階テンソルに拡張することで、テンソルに対する下記の分解を得る。

$$\mathcal{A} = \mathbf{S} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{U}_3 \quad (9)$$

この分解においてテンソル \mathcal{S} をコアテンソルとよび、行列の SVD における行列 \mathbf{S} に相当する。行列の SVD と異なり、ここで $\mathbf{U}_1, \mathbf{U}_2$ および \mathbf{U}_3 が直交行列の場合、コアテンソル \mathcal{S} は密なテンソルとなる^(注1)。このとき式 (9) を高階特異値分解または Tucker 分解とよぶ [14], [15]。PCA はデータ行列に対する SVD と見なせるため、データテンソルを導入した場合、特徴空間の構築は Tucker 分解によって拡張できる。

3.2 Tucker 分解による話者空間の構築

Tucker 分解に基づいて話者空間を構築するため、各事前収録話者を $M \times D'$ の行列で表現する [16]。ここで M は混合数であり、 $D' = 2D$ とする。まずはじめに全データ行列の平均を求め、バイアス行列 $\mathbf{b}' = [\mathbf{b}_1^{(0)}, \mathbf{b}_2^{(0)}, \dots, \mathbf{b}_m^{(0)}]^\top$ とする。これを各事前収録話者の行列からあらかじめ減算しておく。ここで事前収録話者の話者数を S とすると、話者空間を構築するデータセットは 3 階のテンソル $\mathcal{M} \in \mathcal{R}^{M \times D' \times S}$ で表される。このデータテンソルを Tucker 分解すると以下のように表される。

$$\mathcal{M} = \mathbf{G}^{M \times D' \times S} \times_1 \mathbf{U}^{(M)} \times_2 \mathbf{U}^{(D')} \times_3 \mathbf{U}^{(S)} \quad (10)$$

ここで、 $\mathbf{U}^{(M)} \in \mathcal{R}^{M \times M}$, $\mathbf{U}^{(D')} \in \mathcal{R}^{D' \times D'}$, $\mathbf{U}^{(S)} \in \mathcal{R}^{S \times S}$ である。データテンソル \mathcal{M} からのこれらの基底行列の導出については付録 1. にて詳細を述べる。これらの行列は、それぞれ

(注1)：行列の SVD では分解される行列 \mathbf{U}, \mathbf{V} は直交行列であり、行列 \mathbf{S} は対角要素のみ成分を持つスパースな行列となる。しかし一般にテンソルにおける式 (9) のような分解においては行列 \mathbf{U}_n の直交性とコアテンソル \mathcal{S} の対角性は同時には満たされない。

GMMの混合要素, 平均ベクトルの次元, 話者インデックスの効果を捉えており, コアテンソル \mathcal{G} がこれらを統合している. ここで, 第3モードのインデックスを固定することで, 話者 n を表す行列を表現することができる.

$$\boldsymbol{\mu}^{(n)} = \mathcal{G} \times_1 \mathbf{U}^{(M)} \times_2 \mathbf{U}^{(D')} \times_3 \mathbf{U}^{(S)}(n, :) \quad (11)$$

式(11)のうち, 話者空間の基底およびその重みパラメータの組み合わせとして複数の候補が考えられる. $\mathbf{U}^{(S)}(n, :) \in \mathcal{R}^{1 \times S}$ を重みパラメータとし, その他の項のモード積を基底であるとした場合, 式(11)はEVCにおける話者表現と等価になる. すなわち Tucker 分解を用いたデータテンソルの分解は, 固有声に基づく話者空間表現を一般化したものと考えることができる. 一方, 本研究では, 特にGMMの混合要素を用いることで効率的な話者表現が可能であると考え, [16]と同様に以下のようにグルーピングする.

$$\boldsymbol{\mu}^{(n)} = \mathbf{U}^{(M)} \left\{ \mathcal{G} \times_2 \mathbf{U}^{(D')} \times_3 \mathbf{U}^{(S)}(n, :) \right\}^{\top} = \mathbf{U}^{(M)} \mathbf{W}_n^{\top} \quad (12)$$

ここで $\mathbf{U}^{(M)}$ を基底とし, $\mathbf{W}_n \in \mathcal{R}^{D' \times M}$ を重み行列とする. 次元圧縮の観点から, 基底を縮約することで, 任意の話者は以下のような行列で表される.

$$\boldsymbol{\mu}^{(new)} = \mathbf{U}^{(M)} \mathbf{W}_{(new)}^{\top} + \mathbf{b}' \quad (13)$$

$\mathbf{U}^{(M)} \in \mathcal{R}^{M \times K}$ ($K \leq S$), $\mathbf{W}_{(new)} \in \mathcal{R}^{D' \times K}$ はそれぞれ, 表現行列および重み行列となる. ゆえに提案法では, K 次元の重みベクトルを推定するEVCと異なり, $D' \times K$ の重み行列を推定することになる.

[16]では, 最小平均二乗誤差に基づく適応アルゴリズムが提案されているが, 本研究ではEVCにおける適応と同様に最尤基準を導入する. 最尤基準では, 以下のように \mathbf{W} を推定する.

$$\begin{aligned} \hat{\mathbf{W}} &= \underset{\mathbf{W}}{\operatorname{argmax}} \int P(\mathbf{X}, \mathbf{Y}^{(tar)} | \boldsymbol{\lambda}^{(EV)}, \mathbf{W}) d\mathbf{X} \\ &= \underset{\mathbf{W}}{\operatorname{argmax}} P(\mathbf{Y}^{(tar)} | \boldsymbol{\lambda}^{(EV)}, \mathbf{W}) \end{aligned} \quad (14)$$

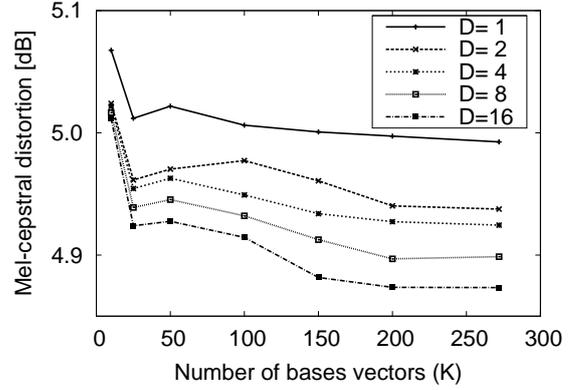
上記の確率密度分布も式(3)と同様GMMとなるため, 式(14)と同形の \mathbf{W} に関する補助関数を導入することで, 最終的に以下の更新式により重み行列を推定する.

$$\operatorname{vec}(\mathbf{W}) = \left[\sum_{m=1}^M \bar{\gamma}_m^{(tar)} \mathbf{U}_m^{\top} \mathbf{U}_m \otimes \boldsymbol{\Sigma}_m^{(YY)^{-1}} \right]^{-1} \operatorname{vec}(\mathbf{C}) \quad (15)$$

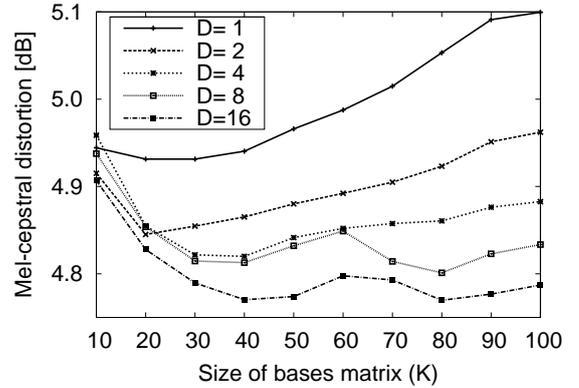
$$\mathbf{C} = \sum_{t=1}^T \sum_{m=1}^M \gamma_{m,t} \boldsymbol{\Sigma}_m^{(YY)^{-1}} (\mathbf{Y}_t^{(tar)} - \mathbf{b}_m^{(0)}) \mathbf{U}_m \quad (16)$$

$$\mathbf{U}_m = \mathbf{U}^{(M)}(m, :) \in \mathcal{R}^{1 \times K} \quad (17)$$

ここで, $\operatorname{vec}()$ は行列を列ベクトルに展開する演算子である. 式(5)と比較すると, 基底の組み合わせは異なるものの近い形をしている. $\mathbf{U}^{(M)}$ はGMMの異なる要素分布の関係性を記述していると考えられ, \mathbf{W} を推定することで, M 混合のGMMを効率的に推定していると解釈できる. 推定するパラメータ数について, EVCでは K 個のパラメータを推定するが, 提案法では $D' \times K$ 個のパラメータを推定することになり, 提案法はより柔軟にデータに対して適応しようとえられる.



(a): Number of bases vectors vs. mel-cepstral distortion (EVC).



(b): Size of bases matrix vs. mel-cepstral distortion (proposed).

図2 基底のサイズに対する変換精度の変化; D は適応文数を表す.
Fig. 2 Mel-cepstral distortion as a function of the number of bases. D denotes the number of adaptation sentences.

4. 実験

4.1 実験条件

提案手法の有効性を確かめるため, 一対多声質変換の実験を行った. 参照話者として ATR 日本語音声データベース [17] から男性 1 名のデータを用いた. また事前収録話者として JNAS から男性話者 137 名, 女性話者 136 名の計 273 名の発声を用いた [18]. 各事前収録話者は 50 文を読み上げている. 評価対象話者として男女 3 名ずつを選んだ. 適応文数を 1 文から 16 文で変化させ, 各話者 21 文を評価に用いた.

スペクトル特徴量として, STRAIGHT 分析に基づくスペクトルから得られた 24 次のメルケプストラムを用いた [19]. STRAIGHT による合成に用いる非周期性指標については全周波数において -30 dB とした. パワーおよび基本周波数については平均と標準偏差を考慮した単純な線形変換によって変換した. また GMM の混合数 (M) は 128 とした.

変換性能について, 提案法に基づく一対多声質変換および一対多 EVC を比較した. また参考として従来のパラレル学習に基づく声質変換の性能についても検証した [13]. なお本研究では話者適応学習は行っていないが, 提案法は話者空間表現に着眼しているため話者適応学習によるモデルの精緻化も導入可能である.

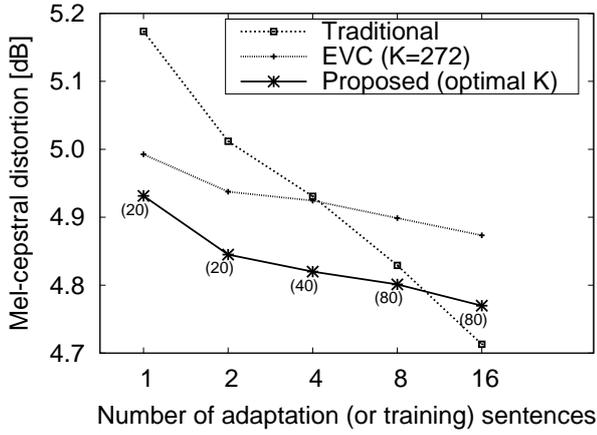


図3 客観評価実験結果; 提案法の () は最適な K を表す。

Fig. 3 Results of objective evaluation.

4.2 客観評価実験

EVC の表現ベクトルの数及び、提案法における基底行列のサイズを変化させた時のメルケプストラム歪みの値を Fig. 2 に示す。Fig. 2(a) より、適応データの量 (D) に関わらず、EVC においては事前学習話者のデータから得られた表現ベクトルを全て用いる場合が最も性能が高いことがわかる。一方、提案法においては適応データの量に応じて最適な基底サイズが変化している。適応データが多い場合には基底サイズが大きくなり、より表現力の高いモデルになっていることがわかる。その場合でも、元の混合数 ($M = 128$) に対しては、よりコンパクトなモデルとなっている。

次に、適応データ数に対するメルケプストラム歪みに基づく客観評価の結果を Fig. 3 に示す。従来のパラレル学習の結果は、それぞれのデータ数で最適な混合数を選択している。従来のパラレル学習と比べると、適応データ数が少ない場合は提案法、EVC とともに高い変換精度となっている。すなわち事前収録話者のデータが効果的に作用しているといえる。学習データ数が多くなると、従来のパラレル学習の性能が提案法、EVC を上回ってくる。これは相互共分散行列の学習によって、より精緻な変換行列が学習されるためと考えられ、提案法においても、話者適応学習を導入することによってデータ数が多い場合でもパラレル学習に迫る性能が期待する。一方提案法は、いずれの適応文数でも EVC の性能を上回っている。これは提案する話者空間表現がスーパーベクトルによる表現に比べて、より声質変換において有効であるといえる。すなわち提案法は、GMM の要素の関係性を適切に捉え、話者空間を効果的に表現しているといえる。

4.3 主観評価実験

聴取実験により変換音声の自然性と話者性について評価した。被験者は 8 名で、自然性の評価には一対比較法、話者性の評価には RAB 法を用いた。一対比較法では被験者は 2 種類の変換音声を聴取し、どちらの音声により自然かを評価した。RAB 法では、2 種類の変換音声を目標話者の参照音声の後に提示し、どちらが参照音声に近いかを評価した。各被験者はランダムに

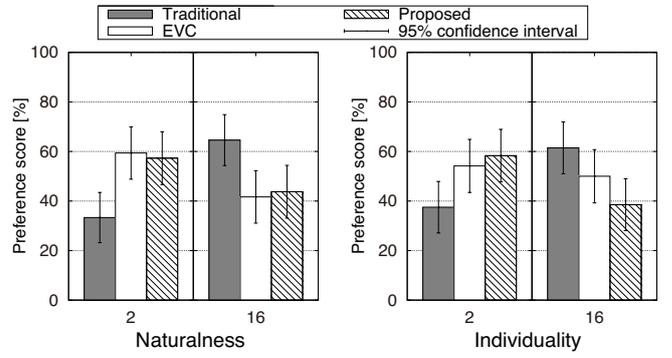


図4 主観評価実験結果

Fig. 4 Results of subjective evaluation.

選ばれた 36 文を評価した。

Fig. 4 に主観評価実験の結果を示す。適応文数が 2 文の場合は、提案法、EVC とともに従来のパラレル学習より自然性が向上していることがわかる。EVC と提案法を比較すると適応文数が 16 文の話者性評価を除いて、提案法が EVC と同等か若干よい性能が得られることがわかる。

5. おわりに

本稿では、話者空間のテンソル表現を用いた声質変換法を提案し、その適応手法について述べた。提案法では各話者を行列の形で表現し、事前収録話者のデータテンソルに対してテンソル解析を適用することでより柔軟に話者を表現することが可能となる。今回は GMM の混合要素を捉えた行列を基底として用い、提案法が EVC に比べて有効に機能することを示した。今後の課題として、話者適応学習や非パラレルデータの利用等、他の有用な手法と提案法を統合することが上げられる。合成音声の品質向上の観点から非周期性指標に対して提案法を適用することも有効と考えられる。またテンソルの基底と重みの選択についても検討の余地がある。今回、適応データ数によって差が見られた表現行列のサイズの最適化についても検討していく予定である。

謝 辞

本研究は科研費・研究活動スタート支援 (11025460) の助成を受けたものである。

付 録

1. Tucker 分解における基底行列の導出

式 (10) における基底行列 $U^{(M)}$ 、 $U^{(D')}$ 及び $U^{(S)}$ の導出について考える。今、事前収録話者から構築される 3 階のデータテンソル $\mathcal{M} \in \mathbb{R}^{M \times D' \times S}$ について、これを各モード毎に平坦化した行列をそれぞれ $M_{(M)}$ 、 $M_{(D')}$ 、 $M_{(S)}$ とする。これらの平坦化行列を以下のように特異値分解する。

$$M_{(M)} = U^{(M)} \mathbf{S}^{(M)} \mathbf{V}^{(M)\top} \quad (\text{A-1})$$

$$M_{(D')} = U^{(D')} \mathbf{S}^{(D')} \mathbf{V}^{(D')\top} \quad (\text{A-2})$$

$$M_{(S)} = U^{(S)} \mathbf{S}^{(S)} \mathbf{V}^{(S)\top} \quad (\text{A-3})$$

上記の特異値分解によって得られた左特異値行列が、式 (10) における基底行列に対応する。このように得られた基底行列を用いて、以下のようにコアテンソル \mathcal{G} を導出する。

$$\mathcal{G} = \mathcal{M} \times_1 \mathbf{U}^{(M)\top} \times_2 \mathbf{U}^{(D')\top} \times_3 \mathbf{U}^{(S)\top} \quad (\text{A.4})$$

ただし \top は行列の転置を表す。

文 献

- [1] A. Kain and M. W. Macon, “Spectral voice conversion for text-to-speech synthesis,” Proc. ICASSP, vol. 1, pp. 285–288, 1998.
- [2] L. Deng, A. Acero, L. Jiang, J. Droppo, and X. Huang, “High-performance robust speech recognition using stereo training data,” Proc. ICASSP, pp. 301–304, 2001.
- [3] A. Kunikoshi, Y. Qiao, N. Minematsu, and K. Hirose, “Speech generation from hand gestures based on space mapping,” Proc. INTERSPEECH, pp. 308–311, 2009.
- [4] M. Abe, S. Nakamura, K. Shikano, and H. Kuwabara, “Voice conversion through vector quantization,” Proc. ICASSP, pp. 655–658, 1988.
- [5] S. Desai, E. V. Raghavendra, B. Yegnanarayana, A. W. Black, and K. Prahallad, “Voice conversion using artificial neural networks,” Proc. ICASSP, pp. 3893–3896, 2009.
- [6] Y. Stylianou, O. Cappé, and E. Moulines, “Continuous probabilistic transform for voice conversion,” IEEE Trans. on Speech and Audio Processing, vol. 6, no. 2, pp. 131–142, 1998.
- [7] A. Mouchtaris, J. V. der Spiegel, and P. Mueller, “Non-parallel training for voice conversion based on a parameter adaptation approach,” IEEE Trans. on Audio, Speech, and Language Processing, vol. 14, no. 3, pp. 952–963, 2006.
- [8] C. H. Lee and C. H. Wu, “Map-based adaptation for speech conversion using adaptation data selection and non-parallel training,” Proc. INTERSPEECH, pp. 2254–2257, 2006.
- [9] T. Toda, Y. Ohtani, and K. Shikano, “Eigenvoice conversion based on Gaussian mixture model,” Proc. INTERSPEECH, pp. 2446–2449, 2006.
- [10] M. A. O. Vasilescu and D. Terzopoulos, “Multilinear analysis of image ensembles: TensorFaces,” Proc. ECCV, pp. 447–460, 2002.
- [11] Y. Ohtani, T. Toda, H. Saruwatari, and K. Shikano, “Speaker adaptive training for one-to-many eigenvoice conversion based on Gaussian mixture model,” Proc. INTERSPEECH, pp. 1981–1984, 2007.
- [12] Y. Ohtani, T. Toda, H. Saruwatari, and K. Shikano, “Non-parallel training for many-to-many eigenvoice conversion,” Proc. ICASSP, pp. 4822–4825, 2010.
- [13] T. Toda, A. W. Black, and K. Tokuda, “Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory,” IEEE Trans. on Audio, Speech, and Language Processing, vol. 15, no. 8, pp. 2222–2235, 2007.
- [14] L. De Lathauwer, B. De Moor and J. Vandewalle, “A multilinear singular value decomposition,” SIAM Journal on Matrix Analysis and Applications, vol. 21, No. 4, pp. 1253–1278, 2000.
- [15] L. R. Tucker, “Some mathematical notes on three-mode factor analysis,” Psychometrika, vol. 31, no. 3, pp. 279–311, 1966.
- [16] Y. Jeong, “Speaker adaptation based on the multilinear decomposition of training speaker models,” Proc. ICASSP, pp. 4870–4873, 2010.
- [17] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara, and K. Shikano, “ATR Japanese speech database as a tool of speech recognition and synthesis,” Speech Communication, vol.9, pp.357–363, 1990.
- [18] “Jnas: Japanese newspaper article sentences,”

- <http://www.milab.is.tsukuba.ac.jp/jnas/instruct.html>
- [19] H. Kawahara, I. Masuda-Katsuse, and A. de Cheveigné, “Restructuring speech representations using a pitch-adaptive time-frequency smoothing and an instantaneous-frequency-based F0 extraction: Possible role of a repetitive structure in sounds,” Speech Communication, vol.27, pp.187–207, 1999.