

話者適応型 Restricted Boltzmann Machine を用いた声質変換の検討

中鹿 亘[†] 滝口 哲也^{††} 有木 康雄^{††}

[†] 神戸大学大学院システム情報学研究科

〒 657-8501 兵庫県神戸市灘区六甲台町 1-1

^{††} 神戸大学自然科学系先端融合研究環

〒 657-8501 兵庫県神戸市灘区六甲台町 1-1

E-mail: [†]nakashika@me.cs.scitec.kobe-u.ac.jp, ^{††}{takigu,ariki}@kobe-u.ac.jp

あらまし 声質変換は、入力した音声を音韻情報などを保ったまま、話者性に関する特定の情報のみを変換する技術であり、話者変換や感情変換、発話支援など様々なタスクへの応用が期待されている。従来の多くの声質変換手法は、同一発話内容の入出力音声対（パラレルデータ）を学習時に必要とするが、予め発話内容を決めておく必要がある、音声間のアライメントを取る必要があるなど、学習データを慎重に用意しなければならないという問題がある。また、変換モデルの利用は学習された話者対のみに限定されてしまう。本研究では、パラレルデータを必要としない任意話者声質変換を実現するため、確率モデルの一つである Restricted Boltzmann machine (RBM) を拡張した話者適応型 RBM (Adaptive restricted Boltzmann machine; ARBM) を新たに提案する。適応型 RBM は可視素子層と隠れ素子層からなる二層の確率モデルであり、異なる層の素子間には話者によって変化する結合重みが存在する。本稿では、適応型 RBM を用いた任意話者声質変換に関する評価実験の結果について報告する。

キーワード 声質変換, restricted Boltzmann machine, 話者適応, 非パラレル学習

Voice Conversion Using Speaker Adaptive Restricted Boltzmann Machine

Toru NAKASHIKA[†], Tetsuya TAKIGUCHI^{††}, and Yasuo ARIKI^{††}

[†] Graduate School of System Informatics, Kobe University

1-1 Rokkodaicho, Nada-ku, Kobe, Hyogo, 657-8501 Japan

^{††} Organization of Advanced Science and Technology, Kobe University

1-1 Rokkodaicho, Nada-ku, Kobe, Hyogo, 657-8501 Japan

E-mail: [†]nakashika@me.cs.scitec.kobe-u.ac.jp, ^{††}{takigu,ariki}@kobe-u.ac.jp

Abstract Voice conversion (VC) is a technique where only speaker-specific information in source speech is converted while keeping phonological information. The technique can be applied to various tasks such as speaker-identity conversion, emotion conversion and aid to speaking for people with articulation disorders. Most of the existing VC methods rely on parallel data—pairs of speech data from source and target speakers uttering the same articles. However, this approach involves several problems; firstly, the data used for the training is limited to the pre-defined articles. Secondly, the use of the trained model is limited only to the speaker pair used in the training. In this paper, we propose a novel probabilistic model called an adaptive restricted Boltzmann machine (ARBM) for VC between arbitrary speakers without use of parallel data. This model consists of a visible-unit and a hidden-unit layer with the speaker-dependent connection. In this paper, we report our experimental results of arbitrary-speaker VC using our model, an ARBM.

Key words Voice conversion, restricted Boltzmann machine, speaker adaptation, non-parallel training

1. はじめに

近年、音声信号処理の分野の中でも、声質変換技術（入力話

者音声の音韻情報を保存したまま、話者性に関する情報のみを出力話者のものへ変換させる技術）が盛んに研究されている。その背景として、雑音環境下 [1] や感情音声 [2] の音声認識精度

の向上、発話困難な障がい者のための発話補助 [3], その他様々なタスク [4], [5] への応用が可能であることが挙げられる。入力話者の音声を出力話者の声質をもつ音声へ変換するためには主に F0 特徴とスペクトル情報を変換する必要があるが, 多くの従来法と同様に, 本研究ではスペクトル情報の変換法に着目する。

これまでの声質変換法では, 統計的手法に基づくアプローチが広く研究されてきた [6], [7]. 中でも GMM (Gaussian Mixture Model) を用いた手法 [8] が最も広く用いられており, 様々な改良がなされてきた。例えば, 戸田ら [9] は, 動的特徴量とグローバルバリエーション (GV) を導入して, 変換精度を向上させる手法を提案した。Helander ら [10] は PLS (Partial Least Squares) を用いて精度を改善する手法を提案した。齋藤ら [11] は, GMM の入力として行列変量を用いる確率分布 (行列変量 GMM) を用いた声質変換法を提案した。GMM 以外のアプローチとしては, 近年 NMF (Non-negative matrix factorization) を用いた声質変換手法 [12] が提案され, 過平滑の少ない手法として注目されている。

しかしながら, これらの手法はいずれもモデルの学習時にパラレルデータ (入力話者と出力話者の, 同一発話内容による音声対) を必要とし, パラレルデータの作成には様々な制限が課せられる。第一に, 発話データは同一の発話内容でないといけないという制限があるため, 選択 (または作成) できる学習データセットの自由度は低い。第二に, フレーム単位で入出力音声の同期を取る必要があるため, 動的計画法などを用いてアライメントを取るが, 完全にフレームの同期が取れている保証がない, 伸縮の際に音声に変換が加わっているなどの問題がある。また, 学習を行っていない話者対に対して, 既存の変換モデルを利用できない。

入出力話者間のパラレルデータを必要としない, 若しくは少量のパラレルデータを用いて, 話者性を柔軟に制御するアプローチもいくつか提案されている [13]~[16]。例えば文献 [13] では, 参照話者のパラレルデータを用いて二話者間の関係性を GMM でモデル化しておき, 入力話者 (もしくは出力話者) を参照話者の特徴空間へ射影する行列を求めるため, 入力話者-出力話者間のパラレルデータは必要としない (しかしながら, 参照話者の間でパラレルデータを必要とする)。また, 文献 [15] では, 予め複数の話者によるパラレルデータを用いて固有声 (Eigenvoice) を作成し, 入力話者から固有声, 固有声から出力話者へマッピングすることで多対多声質変換を実現している (このアプローチでも, 固有声作成時に複数話者のパラレルデータを用意する必要がある)。

本研究では, 確率モデルの一つである restricted Boltzmann machine (RBM) [17] を拡張したモデル (adaptive restricted Boltzmann machine; ARBM) を用いて, 入力話者-出力話者間のパラレルデータだけではなく, 参照話者間のパラレルデータさえも必要としない任意話者声質変換法を提案する。本研究で提案する適応型 RBM は, 複数の話者が混在する音声データから, 話者に依存しない情報と話者に依存した情報に分離しながら, 潜在的な特徴を抽出する確率モデルである。このモデ

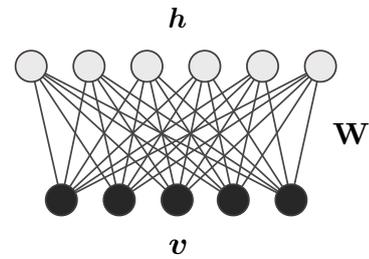


図 1 RBM のグラフ構造。

Fig. 1 Graphical representation of an RBM.

ルは可視素子層と隠れ素子層からなる無向グラフで表現され, 同層素子間の結合はなく, 異層素子間のみ話者に依存した強度 (重み) で結合が存在する。さらに, この重みは話者依存項と話者非依存項で表現され, 複数の話者が混在した音声データ (パラレルである必要はない) を用いて, それぞれが教師なし学習で同時に推定される。結果として, 話者依存重みと話者非依存重みに分離しながら潜在特徴 (隠れ素子) を得ることができる。任意話者声質変換を行う際, まず, 複数の話者 (参照話者) のデータを用いて, 上記のように話者依存重みと話者非依存重みを同時推定する。次に, 変換を行いたい話者 (入力話者) の (少量の) データを用いて, 話者非依存重みを固定しながら新たな話者依存重みを推定する。変換先の話者 (出力話者) の話者依存重みにしても同様に推定する。そして, 変換したい音声から, 入力話者の話者依存重み, 話者非依存重みを用いて潜在特徴を推定し, その後, 出力話者の話者依存重み, 話者非依存重みを用いて音響特徴ベクトルを逆推定することで変換音声を得る。

GMM や NMF など, 従来の声質変換手法の多くは線形変換をベースとしているため, 変換精度には限界がある。つまり, 我々の声道形状は非線形的であるため, 音声信号に含まれる声質の特性をより正確に捉えるためには非線形ベースのモデル化の方が線形ベースよりも適切であると考えられる。非線形ベースの声質変換手法として, Desai らによる多層 NN (Neural Networks) を用いた声質変換法 [18] や, 我々が提案してきた話者依存型 RBM (restricted Boltzmann machine) 若しくは DBN (deep belief networks [19]) を用いた多層型声質変換法 [20], 話者依存型 CRBM (conditional restricted Boltzmann machine) [21] を用いた手法 [22] や話者依存型 RTRBM (recurrent temporal restricted Boltzmann machine) [23] を用いた手法 [24], Wu らによる CRBM を用いた非線形声質変換法 [25], Joint 型 RBM を用いた手法 [26] が挙げられる。いずれの非線形変換に基づく手法においても, 線形変換ベースの手法と比べて高い精度が得られていることが報告されている [18], [20], [22], [25]。本研究で提案する声質変換法も非線形関数をベースとした変換式を用いており, 精度の高いモデル化が期待できる。

2. Restricted Boltzmann Machine

Restricted Boltzmann machine (RBM) は特殊な構造を持つ 2 層ネットワークであり, 図 1 のように, 可視層と隠れ層の確率変数分布を表現する無向グラフィカルモデルである [17]。元々

RBM はバイナリデータを入力させるモデルとして提案されていたが、後に連続値を入力させるモデル (Gaussian-Bernoulli RBM; GBRBM [27]) が考案された。しかしながらこのモデルは、分散項の影響で学習が不安定になるという問題 [19], [28] があったため、Cho らによって GBRBM の改良版 (Improved GBRBM; ImpGBRBM [29]) が提案された。この ImpGBRBM では、連続値の可視素子 $\mathbf{v} = [v_1, \dots, v_I]^T, v_i \in \mathbb{R}$ と 2 値の隠れ素子 $\mathbf{h} = [h_1, \dots, h_J]^T, h_j \in \{0, 1\}$ の同時確率 $p(\mathbf{v}, \mathbf{h})$ は、以下のように表される。

$$p(\mathbf{v}, \mathbf{h}) = \frac{1}{Z} e^{-E(\mathbf{v}, \mathbf{h})} \quad (1)$$

$$E(\mathbf{v}, \mathbf{h}) = \left\| \frac{\mathbf{v} - \mathbf{b}}{2\sigma} \right\|^2 - \mathbf{c}^T \mathbf{h} - \left(\frac{\mathbf{v}}{\sigma^2} \right)^T \mathbf{W} \mathbf{h} \quad (2)$$

$$Z = \sum_{\mathbf{v}, \mathbf{h}} e^{-E(\mathbf{v}, \mathbf{h})} \quad (3)$$

ここで、 $\|\cdot\|^2$ は L2 ノルム、括弧は要素除算を表す。 $\mathbf{W} \in \mathbb{R}^{I \times J}$, $\sigma \in \mathbb{R}^{I \times 1}$, $\mathbf{b} \in \mathbb{R}^{I \times 1}$, $\mathbf{c} \in \mathbb{R}^{J \times 1}$ はそれぞれ可視層-隠れ層間の重み行列、可視素子の偏差、可視素子のバイアス、隠れ素子のバイアスを示しており、いずれも推定すべきパラメータである。

RBM では可視素子間、または隠れ素子間の接続は存在しないため (つまり、それぞれの可視素子、隠れ素子は互いに条件付き独立であるため)、それぞれの条件付き確率 $p(\mathbf{h}|\mathbf{v})$, $p(\mathbf{v}|\mathbf{h})$ は以下の様な単純な関数で表現される。

$$p(h_j = 1|\mathbf{v}) = \mathcal{S}(c_j + \left(\frac{\mathbf{v}}{\sigma^2}\right)^T \mathbf{W}_{:,j}) \quad (4)$$

$$p(v_i = v|\mathbf{h}) = \mathcal{N}(v|b_i + \mathbf{W}_{i,:} \mathbf{h}, \sigma_i^2) \quad (5)$$

ここで、 $\mathbf{W}_{:,j}$ と $\mathbf{W}_{i,:}$ は \mathbf{W} の第 j 行ベクトル、第 i 列ベクトルを表す。また、 $\mathcal{S}(\cdot)$ は要素ごとのシグモイド関数 ($\mathcal{S}(x) = \mathbf{1} \odot (1 + e^{-x})$), $\mathcal{N}(\cdot|\mu, \sigma^2)$ は平均 μ , 分散 σ^2 の正規分布を表す。

それぞれの RBM のパラメータ $\Theta = \{\mathbf{W}, \mathbf{b}, \sigma, \mathbf{c}\}$ は、 N 個の観測データを $\{\mathbf{v}_n\}_{n=1}^N$ とするとき、この確率変数の対数尤度 $\mathcal{L} = \log \prod_n p(\mathbf{v}_n)$ を最大化するように推定される。この対数尤度をそれぞれのパラメータで偏微分すると、

$$\frac{\partial \mathcal{L}}{\partial W_{ij}} = \langle \frac{v_i h_j}{\sigma_i^2} \rangle_{data} - \langle \frac{v_i h_j}{\sigma_i^2} \rangle_{model} \quad (6)$$

$$\frac{\partial \mathcal{L}}{\partial b_i} = \langle \frac{v_i}{\sigma_i^2} \rangle_{data} - \langle \frac{v_i}{\sigma_i^2} \rangle_{model} \quad (7)$$

$$\frac{\partial \mathcal{L}}{\partial c_j} = \langle h_j \rangle_{data} - \langle h_j \rangle_{model}, \quad (8)$$

が得られる。ただし、 $\langle \cdot \rangle_{data}$ と $\langle \cdot \rangle_{model}$ はそれぞれ、観測データ、モデルデータの期待値を表す。しかし、一般に後者の期待値に関しては計算困難であるため、代わりに式 (4)(5) によって得られる再構築したデータの期待値 $\langle \cdot \rangle_{recon}$ が用いられる (CD: Contrastive Divergence 法 [19])。また、ImpGBRBM では分散を非負値に制約し、学習を安定化させるため $\sigma_i^2 = e^{z_i}$ と置き換える。これにより、 z_i に関する勾配は以下のように計算される。

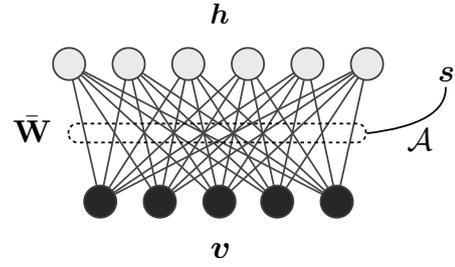


図 2 ARBM のグラフ構造。

Fig. 2 Graphical representation of an ARBM.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial z_i} = e^{-z_i} & \left\langle \frac{(v_i - b_i)^2}{2} - v_i \mathbf{W}_{i,:} \mathbf{h} \right\rangle_{data} \\ & - e^{-z_i} \left\langle \frac{(v_i - b_i)^2}{2} - v_i \mathbf{W}_{i,:} \mathbf{h} \right\rangle_{model} \end{aligned} \quad (9)$$

それぞれのパラメータは式 (6)(7)(8) から、確率的勾配法を用いて繰り返し更新される (初期値はランダムに設定される)。すなわち、

$$\Theta^{(new)} = \Theta^{(old)} - \gamma_\theta \frac{\partial \mathcal{L}}{\partial \Theta} \quad (10)$$

ここで、 γ_θ は学習率を表す。

3. 適応型 RBM と声質変換への応用

本稿では、前節で述べた RBM を拡張したモデルとして、適応型 RBM (Adaptive restricted Boltzmann machine; ARBM) を定義し、声質変換タスクへ応用する手法について述べる。

3.1 適応型 RBM の定義

適応型 RBM は、図 2 のように、通常の RBM で見られた可視素子と隠れ素子だけでなく、識別素子 $\mathbf{s} = [s_1, \dots, s_S]^T, s_k \in \{0, 1\}$ が加わったモデルとなっている (S は識別素子の数とする)。例えば声質変換において、入力 \mathbf{v} が話者 k の発話であることを示す場合、 $s_k = 1, \forall s_{k'} = 0 (k' \neq k)$ となる。このモデルでは、可視素子と隠れ素子の間には識別素子 \mathbf{s} で制御される重みの結合が存在する。この結合重みを $\mathbf{W}(\mathbf{s})$ とし、本稿ではこれを以下のように定義する。

$$\mathbf{W}(\mathbf{s}) = \mathbf{A} \otimes_3 \mathbf{s} \bar{\mathbf{W}} + \mathbf{B} \otimes_3 \mathbf{s} \quad (11)$$

ただし、 \mathbf{A} と \mathbf{B} はいずれも、不特定重み行列 $\bar{\mathbf{W}} \in \mathbb{R}^{I \times J}$ を特定化 (適応) するための 3 階のテンソルパラメータ ($\mathbf{A} \in \mathbb{R}^{I \times I \times S}, \mathbf{B} \in \mathbb{R}^{I \times J \times S}$) である。また、 $\mathcal{X} \otimes_d \mathbf{y}$ はモード d を展開した 3 階テンソル \mathcal{X} の各行列とベクトル \mathbf{y} の内積をとる演算子を表す。声質変換の場合、 $\bar{\mathbf{W}}$ が不特定話者による結合重み、 $\mathbf{A}_{:, :, k}$ と $\mathbf{B}_{:, :, k}$ が話者 k の適応行列及びバイアス行列を表す (ただし $\mathbf{A}_{:, :, k}$ は 3 階テンソル \mathbf{A} のモード 3 の第 k 行行列を表す)。

適応型 RBM では、式 (11) で定義した $\mathbf{W}(\mathbf{s})$ を用いて、可視素子 \mathbf{v} , 隠れ素子 \mathbf{h} , 識別素子 \mathbf{s} の同時確率 $p(\mathbf{v}, \mathbf{h}, \mathbf{s})$ を以下のように定義する。

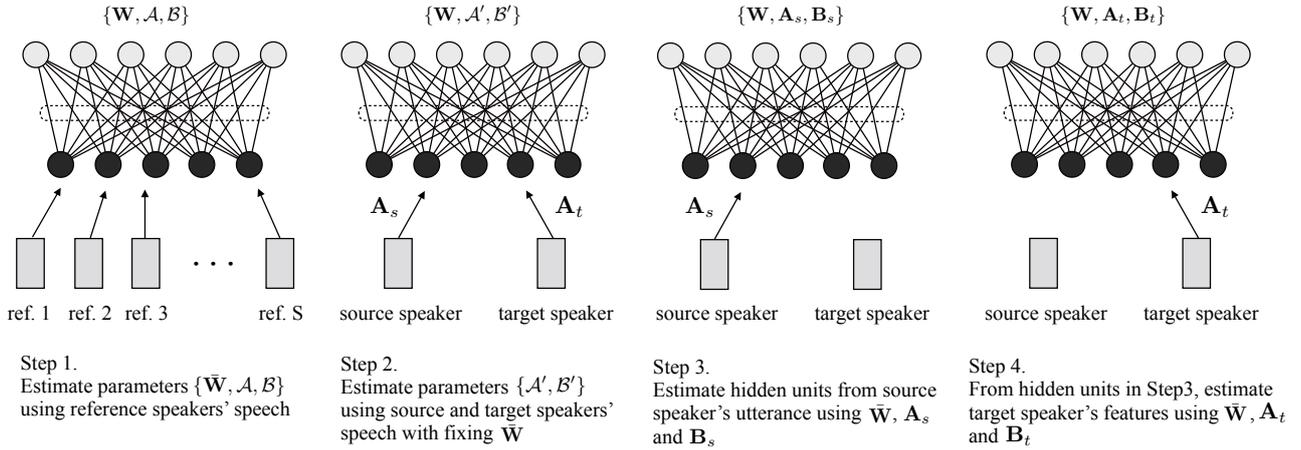


図3 適応型 RBM を用いた声質変換の流れ。

Fig. 3 Procedure of voice conversion using an ARBM.

$$p(\mathbf{v}, \mathbf{h}, \mathbf{s}) = \frac{1}{Z_A} e^{-E_A(\mathbf{v}, \mathbf{h}, \mathbf{s})} \quad (12)$$

$$E_A(\mathbf{v}, \mathbf{h}, \mathbf{s}) = \left\| \frac{\mathbf{v} - \mathbf{b}}{2\sigma} \right\|^2 - \mathbf{c}^T \mathbf{h} - \left(\frac{\mathbf{v}}{\sigma^2} \right)^T \mathbf{W}(\mathbf{s}) \mathbf{h} \quad (13)$$

$$Z_A = \sum_{\mathbf{v}, \mathbf{h}, \mathbf{s}} e^{-E_A(\mathbf{v}, \mathbf{h}, \mathbf{s})} \quad (14)$$

これらの定義により、条件付き確率 $p(\mathbf{h}|\mathbf{v}, \mathbf{s})$, $p(\mathbf{v}|\mathbf{h}, \mathbf{s})$ は以下のように計算できる。

$$p(h_j = 1|\mathbf{v}, \mathbf{s}) = \mathcal{S}(c_j + \left(\frac{\mathbf{v}}{\sigma^2} \right)^T \mathbf{W}(\mathbf{s})_{:j}) \quad (15)$$

$$p(v_i = v|\mathbf{h}, \mathbf{s}) = \mathcal{N}(v|b_i + \mathbf{W}(\mathbf{s})_{i:} \mathbf{h}, \sigma_i^2) \quad (16)$$

適応型 RBM のパラメータ $\Theta_A = \{\bar{\mathbf{W}}, \mathbf{A}, \mathbf{B}, \mathbf{b}, \sigma, \mathbf{c}\}$ は、 N 個の学習データ $\{\mathbf{v}_n, \mathbf{s}_n\}_{n=1}^N$ を用いて、対数尤度 $\mathcal{L}_A = \log \prod_n p(\mathbf{v}_n, \mathbf{s}_n) = \log \prod_n \sum_{\mathbf{h}} p(\mathbf{v}_n, \mathbf{h}_n, \mathbf{s}_n)$ を最大化するように推定される。この対数尤度を $\bar{\mathbf{W}}, \mathbf{A}, \mathbf{B}$ の要素 (\bar{W}_{ij} , $A_{ii'k}$, B_{ijjk}) で偏微分したものは、それぞれ

$$\frac{\partial \mathcal{L}_A}{\partial \bar{W}_{ij}} = \left\langle \sum_{l,k} \frac{A_{lik} v_l h_j s_k}{\sigma_l^2} \right\rangle_{data} - \left\langle \sum_{l,k} \frac{A_{lik} v_l h_j s_k}{\sigma_l^2} \right\rangle_{model} \quad (17)$$

$$\frac{\partial \mathcal{L}_A}{\partial A_{ii'k}} = \left\langle \sum_m \frac{W_{i'm} v_i h_m s_k}{\sigma_i^2} \right\rangle_{data} - \left\langle \sum_m \frac{W_{i'm} v_i h_m s_k}{\sigma_i^2} \right\rangle_{model} \quad (18)$$

$$\frac{\partial \mathcal{L}_A}{\partial B_{ijjk}} = \langle v_i h_j s_k \rangle_{data} - \langle v_i h_j s_k \rangle_{model}, \quad (19)$$

と計算できる。他のパラメータ $\mathbf{b}, \sigma, \mathbf{c}$ に関しては、それぞれ式 (7), (9), (8) と同様にして求められる。適応型 RBM においても、CD 法を適用することができるため、各偏微分値の第二項 $\langle \cdot \rangle_{model}$ を観測データの再構築値 $\langle \cdot \rangle_{recon}$ として計算することで効率よくパラメータを推定することができる。

3.2 適応型 RBM を用いた声質変換

適応型 RBM を声質変換へ応用する場合、図3のようにまず複数 (S 人) の参照話者によるデータを用いて適応型 RBM の各パラメータを同時推定する (Step 1)。次に、 $\bar{\mathbf{W}}$ など話者に依存し

ないパラメータを固定して、適応データを用いて入力話者と出力話者の適応パラメータ $\mathcal{A}' \ni \{\mathbf{A}_s = \mathbf{A}_{::(S+1)}, \mathbf{A}_t = \mathbf{A}_{::(S+2)}\}$, $\mathcal{B}' \ni \{\mathbf{B}_s = \mathbf{B}_{::(S+1)}, \mathbf{B}_t = \mathbf{B}_{::(S+2)}\}$ を式 (18)(19) より推定する (Step 2)。そして、入力話者の変換したい音声のフレーム音響特徴量 \mathbf{v}_s から、次式のように潜在特徴量 (隠れ素子) を推定する (Step 3)。

$$\hat{\mathbf{h}} = \underset{\mathbf{h}}{\operatorname{argmax}} p(\mathbf{h}|\mathbf{v}_s, \mathbf{s}_s) = \mathcal{S}\left(\mathbf{c} + \left(\frac{\mathbf{v}_s}{\sigma^2}\right)^T \mathbf{W}(\mathbf{s}_s)\right) \quad (20)$$

ただし、 \mathbf{s}_s は第 $S+1$ 要素のみ 1, 他を 0 とするベクトルとする。また、同時に変数 \mathbf{s} の長さを $S+2$ へ拡張し、 \mathbf{A}, \mathbf{B} をモード 3 に沿ってそれぞれ \mathbf{A}', \mathbf{B}' を追加するものとする。式 (20) を書き直すと、

$$\hat{\mathbf{h}} = \mathcal{S}\left(\mathbf{c} + \left(\frac{\mathbf{v}_s}{\sigma^2}\right)^T (\mathbf{A}_s \bar{\mathbf{W}} + \mathbf{B}_s)\right) \quad (21)$$

が得られ、話者に依存しない項 $\bar{\mathbf{W}}$ を入力話者に適応させた結合重みを用いて潜在特徴量を推定していることになる。また式 (21) は、一度適応型 RBM の学習が終われば $\hat{\mathbf{h}}$ は変数 \mathbf{v} の関数となるので、 $\hat{\mathbf{h}}$ は話者に依存しない潜在特徴量であることを示唆している。すなわち、話者性は \mathbf{s} のみで制御され、 $\hat{\mathbf{h}}$ は話者に依存しない音韻に近い情報を表すと考えられる。したがって、出力話者の話者性を持つ音声を得たい場合、音韻情報 $\hat{\mathbf{h}}$ から、 \mathbf{s}_t (第 $S+2$ 要素のみ 1, 他が 0 となるベクトル) を用いて音響特徴量を復元すればよい。すなわち、出力話者の変換先のフレーム特徴量 \mathbf{v}_t を以下のように計算する (Step 4)。

$$\begin{aligned} \mathbf{v}_t &= \underset{\mathbf{v}}{\operatorname{argmax}} p(\mathbf{v}|\hat{\mathbf{h}}, \mathbf{s}_t) \\ &= \mathbf{b} + \mathbf{W}(\mathbf{s}_t)^T \hat{\mathbf{h}} \\ &= \mathbf{b} + (\mathbf{A}_t \bar{\mathbf{W}} + \mathbf{B}_t)^T \hat{\mathbf{h}} \end{aligned} \quad (22)$$

これは、入力話者音声から得られた音韻情報を基に、話者非依存項を出力話者に適応した基底を用いて、出力話者の音響特徴量を生成していることを表している。また、式 (21)(22) にもあるように、入力話者の音響特徴量 \mathbf{v}_s を出力話者の音響特徴量 \mathbf{v}_t へ変換する際、 $\hat{\mathbf{h}}$ の推定に非線形関数を用いているため、提

表 1 提案手法による声質変換の客観評価値 (SDIR [dB]).

Table 1 Performance of our VC method (SDIR [dB]).

| # of hidden units | 128 | 192 | 256 | 512 |
|-------------------|------|------|-------------|------|
| female-to-female | 7.18 | 7.26 | 7.30 | 7.14 |
| female-to-male | 7.64 | 7.81 | 7.81 | 7.82 |
| male-to-female | 7.50 | 7.54 | 7.61 | 7.48 |
| male-to-male | 7.86 | 8.00 | 8.03 | 8.06 |
| avg. | 7.54 | 7.65 | 7.69 | 7.63 |

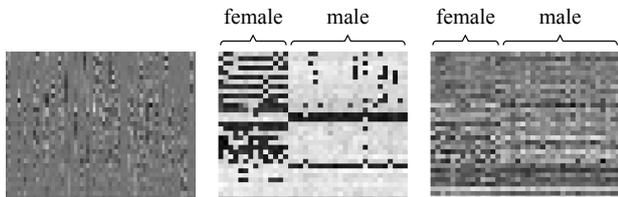


図 4 左から順に, 推定された $\hat{\mathbf{W}}$, \mathbf{A} , \mathbf{B} の一部を表す.

Fig. 4 From left to right: estimated weights $\hat{\mathbf{W}}$, \mathbf{A} , and \mathbf{B} .

案法は非線形変換ベースの声質変換だと言える.

なお, 現実の音声データを使って適応型 RBM を学習する場合, 話者は豊富に存在するが, それぞれの発話データは少ないといったケースがある. この場合, $\hat{\mathbf{W}}$ の推定に用いられるデータは十分存在するが, 適応パラメータ \mathbf{A} , \mathbf{B} を推定するためのデータが少量となるため, 誤推定もしくは過学習の要因となる. そこで本稿で述べる実験では, $\mathbf{A}_{::k}$ を対角行列, $\mathbf{B}_{::k}$ を各列が等しい行列で近似することでパラメータ数を抑える.

4. 評価実験

4.1 実験条件

本実験では, 英語圏の複数の話者による音声が含まれたコーパスである TIMIT [30] を用いて, 提案手法である適応型 RBM を用いた声質変換の精度を調べた. このコーパスから, 話者非依存パラメータの推定のために, 参照話者として 38 名 (内女性 14 名男性 24 名) を選んだ. 各話者からは, 5 文の発話データを学習に用いている (学習に用いた総フレーム数はおおよそ 27 万). 提案手法を評価するために, 女性 4 名, 男性 4 名の音声を用いて入力話者・出力話者のペア (計 28 ペア) を作成し, 異性間及び同性間の声質変換の性能比較を行った. このとき, 入力・出力話者のパラレルデータ (同一発話内容による, 学習データには含まれない 2 文のデータから動的計画法によって作成) を用いて SDIR (spectral distortion improvement ratio) による評価をおこなっている. 音響特徴量として, STRAIGHT スペクトル [31] から計算された 32 次元の MFCC を用いた. 適応型 RBM における学習率, バッチサイズ, 繰り返し回数はそれぞれ 0.005, 50, 500 とした. 隠れ素子数を 128, 192, 256, 512 と変えて比較を行った.

4.2 実験結果と考察

提案手法による声質変換の結果を表 1 に示す. 例えば female-to-female では評価用の女性 4 名の音声を, それぞれ他の女性 3 名へ変換し, 全フレームの SDIR の平均をとったものを表す.

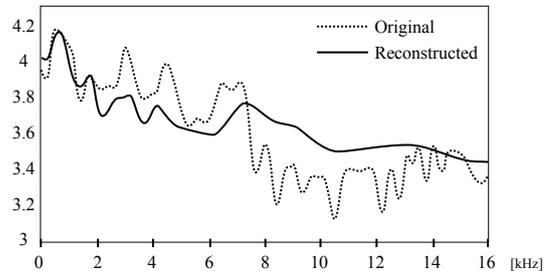
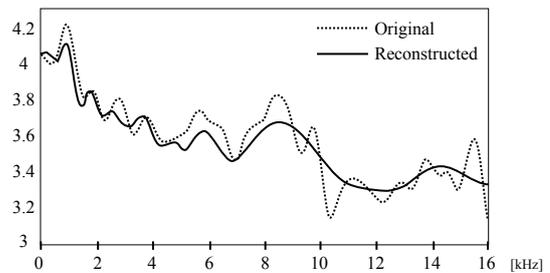


図 5 入力話者の音声スペクトルと再構築されたスペクトル (上図), 出力話者の音声スペクトルと入力話者から出力話者へ変換されたスペクトル (下図).

Fig. 5 Log-spectrum from a source speaker and the reconstructed spectrum (above), and log-spectrum from a target speaker and the converted spectrum from the source speech to the target speech (below).

“avg.” は全組み合わせの平均値である. 表 1 から, 一部を除いて隠れ素子数が増加すれば変換精度が向上していることが分かる. 隠れ素子数が 512 と 256 の結果を比較すると, 512 の場合は男性への変換 (female-to-male, male-to-male) で優っているが, 女性への変換 (female-to-female, male-to-female) で精度が下がってしまい, 結果として全平均の SDIR 値が低くなってしまっている. この理由として, パラメータ数の増加に伴い, モデルが過学習しているためだと考えられる (男性と女性の話者数は 14 対 28 であり, 隠れ素子数 512 のモデルでは変換音声は男性側へ強く反応していることから過学習が窺える).

提案法によって, 実際に推定されたパラメータ $\hat{\mathbf{W}}$, \mathbf{A} および \mathbf{B} の一部を図 4 に示す. \mathbf{A} に関しては, 対角行列として近似した $\mathbf{A}_{::k}$ の対角成分を列ベクトルとして話者ごとに並べた行列を示しており (\mathbf{B} も同様に話者ごとに並べた列ベクトルを示している), 左 14 列ベクトルは女性話者, 右 24 列ベクトルは男性話者に相当する. この図から分かるように, \mathbf{A} (または \mathbf{B}) の各々の列ベクトルは同性間で類似性が高く, 異性間で類似性が低いベクトルとなっている. これは, 音声を聴いて話者の違いを認識する際, 個人の差異よりも性別の違いをより大きく感じ取っているという直感と一致する.

最後に, 提案手法によって女性話者音声 (コーパスでは FCJF0) を男性話者音声 (MWAR0) へ変換した例を図 5 に示す. この例では, FCJF0 のある時刻における対数スペクトル (図上段点線) から MFCC を計算し, FCJF0 の適応型 RBM によって $\hat{\mathbf{h}}$ を推定した後, MWAR0 の適応パラメータを用いて変換された音響特徴量を対数スペクトルへ復元した (図下段実線). 参考として, $\hat{\mathbf{h}}$ の推定後 FCJF0 の適応パラメータによ

て復元したスペクトル（図上段実線），目標となる MWAR0 のスペクトル（図下段点線）を載せている．この図より，FCJF0 の音声から FCJF0 の音声へ再構築したスペクトルのみならず，別の話者である MWAR0 へ変換した音声スペクトルにおいても，低域におけるスペクトルピークの周波数（フォルマント）がおおよそ目標と一致するなど，その話者の特徴を捉えていることが分かる．高周波数域に関してはいずれも目標と大きく異なっているが，MFCC からスペクトルを復元しているため，高域における情報が損失してしまうことに起因する．パラレルデータを学習時に一切使用せず，かつ FCJF0 から MWAR0 への変換モデルを学習していないにも関わらず FCJF0 から MWAR0 へ変換できていることは提案手法の大きな利点であると言える．

5. おわりに

本研究では，潜在的な特徴量を抽出する RBM を拡張して，話者に依存する項と依存しない項に分離してモデル化することで学習時にパラレルデータを必要としない任意話者声質変換手法を提案した．本研究で提案する RBM の拡張モデル（適応型 RBM）は声質変換のみならず，音声の感情付与や物体認識など，様々なタスクへの応用が考えられる．また，このモデルにおいて識別素子 s を推定することで，例えば話者認識へ応用することも可能であると考えられる．音韻情報と話者情報が混在した音声からそれぞれを分離し，話者性を制御できる点が適応型 RBM の強みであり，今後は適応型 RBM を用いた話者認識と音声認識の同時推定法について検討していきたい．

文 献

- [1] A. Kain and M. W. Macon: "Spectral voice conversion for text-to-speech synthesis", ICASSP, pp. 285–288 (1998).
- [2] C. Veaux and X. Robet: "Intonation conversion from neutral to expressive speech", Interspeech, pp. 2765–2768 (2011).
- [3] K. Nakamura, T. Toda, H. Saruwatari and K. Shikano: "Speaking-aid systems using gmm-based voice conversion for electrolaryngeal speech", Speech Communication, **54**, 1, pp. 134–146 (2012).
- [4] L. Deng, A. Acero, L. Jiang, J. Droppo and X. Huang: "High-performance robust speech recognition using stereo training data", ICASSP, pp. 301–304 (2001).
- [5] A. Kunikoshi, Y. Qiao, N. Minematsu and K. Hirose: "Speech generation from hand gestures based on space mapping", Interspeech, pp. 308–311 (2009).
- [6] R. Gray: "Vector quantization", IEEE ASSP Magazine, **1**, 2, pp. 4–29 (1984).
- [7] H. Valbret, E. Moulines and J.-P. Tubach: "Voice transformation using PSOLA technique", Speech Communication, **11**, 2, pp. 175–187 (1992).
- [8] Y. Stylianou, O. Cappé and E. Moulines: "Continuous probabilistic transform for voice conversion", IEEE Trans. Speech and Audio Process., **6**, 2, pp. 131–142 (1998).
- [9] T. Toda, A. W. Black and K. Tokuda: "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory", IEEE Trans. Audio, Speech, and Lang. Process., **15**, 8, pp. 2222–2235 (2007).
- [10] E. Helander, T. Virtanen, J. Nurminen and M. Gabbouj: "Voice conversion using partial least squares regression", IEEE Trans. Audio, Speech, and Lang. Process., **18**, 5, pp. 912–921 (2010).
- [11] N. M. Daisuke Saito, Hidenobu Doi and K. Hirose: "Application of matrix variate gaussian mixture model to statistical voice conversion", Interspeech, pp. 2504–2508 (2014).
- [12] R. Takashima, T. Takiguchi and Y. Ariki: "Exemplar-based voice conversion in noisy environment", SLT, pp. 313–317 (2012).
- [13] A. Mouchtaris, J. Van der Spiegel and P. Mueller: "Non-parallel training for voice conversion based on a parameter adaptation approach", IEEE Trans. Audio, Speech, and Lang. Process., **14**, 3, pp. 952–963 (2006).
- [14] C.-H. Lee and C.-H. Wu: "Map-based adaptation for speech conversion using adaptation data selection and non-parallel training", Interspeech, pp. 2254–2257 (2006).
- [15] T. Toda, Y. Ohtani and K. Shikano: "Eigenvoice conversion based on gaussian mixture model", Interspeech, pp. 2446–2449 (2006).
- [16] D. Saito, K. Yamamoto, N. Minematsu and K. Hirose: "One-to-many voice conversion based on tensor representation of speaker space", Interspeech, pp. 653–656 (2011).
- [17] Y. Freund and D. Haussler: "Unsupervised learning of distributions of binary vectors using two layer networks", Computer Research Laboratory (1994).
- [18] S. Desai, E. V. Raghavendra, B. Yegnanarayana, A. W. Black and K. Prahallad: "Voice conversion using artificial neural networks", ICASSP, pp. 3893–3896 (2009).
- [19] G. E. Hinton, S. Osindero and Y. W. Teh: "A fast learning algorithm for deep belief nets", Neural computation, **18**, 7, pp. 1527–1554 (2006).
- [20] T. Nakashika, R. Takashima, T. Takiguchi and Y. Ariki: "Voice conversion in high-order eigen space using deep belief nets", Interspeech, pp. 369–372 (2013).
- [21] G. W. Taylor, G. E. Hinton and S. T. Roweis: "Modeling human motion using binary latent variables", Advances Neural Info. Process. Systems, pp. 1345–1352 (2006).
- [22] T. Nakashika, T. Takiguchi and Y. Ariki: "Voice conversion in time-invariant speaker-independent space", ICASSP, pp. 7939–7943 (2014).
- [23] I. Sutskever, G. E. Hinton and G. W. Taylor: "The recurrent temporal restricted boltzmann machine", Advances Neural Inform. Process. Systems, pp. 1601–1608 (2009).
- [24] T. Nakashika, T. Takiguchi and Y. Ariki: "High-order sequence modeling using speaker-dependent recurrent temporal restricted boltzmann machines for voice conversion", Interspeech, pp. 2278–2282 (2014).
- [25] Z. Wu, E. S. Chng and H. Li: "Conditional restricted boltzmann machine for voice conversion", ChinaSIP (2013).
- [26] L. H. Chen, Z. H. Ling, Y. Song and L. R. Dai: "Joint spectral distribution modeling using restricted boltzmann machines for voice conversion", Interspeech, pp. 3052–3056 (2013).
- [27] G. E. Hinton and R. R. Salakhutdinov: "Reducing the dimensionality of data with neural networks", Science, **313**, 5786, pp. 504–507 (2006).
- [28] A. Krizhevsky and G. Hinton: "Learning multiple layers of features from tiny images", Computer Science Department, University of Toronto, Tech. Rep (2009).
- [29] K. Cho, A. Ilin and T. Raiko: "Improved learning of gaussian-bernoulli restricted boltzmann machines", ICANN, Springer, pp. 10–17 (2011).
- [30] J. S. Garofolo, L. D. Consortium, et al.: "TIMIT: acoustic-phonetic continuous speech corpus", Linguistic Data Consortium (1993).
- [31] H. Kawahara, M. Morise, T. Takahashi, R. Nisimura, T. Irino and H. Banno: "TANDEM-STRAIGHT: A temporally stable power spectral representation for periodic signals and applications to interference-free spectrum, F0, and aperiodicity estimation", ICASSP, pp. 3933–3936 (2008).