適応型 Restricted Boltzmann Machine を用いた パラレルデータフリーな任意話者声質変換*

中鹿亘,滝口哲也,有木康雄(神戸大)

1 はじめに

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,様々なタスク[1,2]への応用が可能であることから近年盛んに研究されている.これまでの声質変換法として,GMM(Gaussian Mixture Model)を用いた手法[3]が最も広く用いられており,様々な改良がなされてきた[4,5].

しかしながら,ほとんどの従来手法ではモデルの 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[6,7,8,9]. 例えば文献[6]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 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[8]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) [10] 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

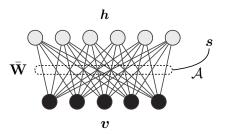


Fig. 1 Graphical representation of an ARBM.

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,入力話 者と出力話者の少量データ(適応データ)を用いて, 話者非依存重みを固定しながらそれぞれの話者依存 重みを推定する.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.

2 適応型 RBM

本稿で定義する適応型 RBM は,Fig. 1 のように,従来の RBM [10] で見られた可視素子と隠れ素子だけでなく,識別素子 $s=[s_1,\cdots,s_S]^T,s_k\in\{0,1\}$ が加わったモデルである(S は識別素子の数とする).例えば声質変換において,入力 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

^{*}Parallel-dictionary-free voice conversion using adaptive restricted Boltzmann machine. by Toru NAKASHIKA, Tetsuya TAKIGUCHI, Yasuo ARIKI (Kobe University)

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sで制御される重みの結合が存在する.この結合重み $\mathbf{W}(s)$ を以下のように定義する.

$$\mathbf{W}(s) = \mathcal{A} \otimes_3 s \bar{\mathbf{W}} + \mathcal{B} \otimes_3 s \tag{1}$$

ただし,A と 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($A \in \mathbb{R}^{I \times I \times S}, \mathcal{B} \in \mathbb{R}^{I \times J \times S}$)である.また, $\mathcal{X} \otimes_d y$ 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 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, $\mathbf{A}_{:,i,k}$ と $\mathbf{B}_{:,i,k}$ が話者 k の適応行列及びバイアス行列を表す(ただし $\mathbf{A}_{:,i,k}$ は 3 階テンソル A のモード 3 の第k 行列を表す).

適応型 RBM では , 式 (1) で定義した $\mathbf{W}(s)$ を用いて , 可視素子 v , 隠れ素子 h , 識別素子 s の同時確率 p(v,h,s) を以下のように定義する .

$$p(\boldsymbol{v}, \boldsymbol{h}, \boldsymbol{s}) = \frac{1}{Z_{A}} e^{-E_{A}(\boldsymbol{v}, \boldsymbol{h}, \boldsymbol{s})}$$
(2)
$$E_{A}(\boldsymbol{v}, \boldsymbol{h}, \boldsymbol{s}) = \left\| \frac{\boldsymbol{v} - \boldsymbol{b}}{2\boldsymbol{\sigma}} \right\|^{2} - \boldsymbol{c}^{T} \boldsymbol{h} - \left(\frac{\boldsymbol{v}}{\boldsymbol{\sigma}^{2}} \right)^{T} \mathbf{W}(\boldsymbol{s}) \boldsymbol{h}$$
(3)

$$Z_{\mathbf{A}} = \sum_{\boldsymbol{v}, \boldsymbol{h}, \boldsymbol{s}} e^{-E_{\mathbf{A}}(\boldsymbol{v}, \boldsymbol{h}, \boldsymbol{s})} \tag{4}$$

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

$$p(h_j = 1 | \boldsymbol{v}, \boldsymbol{s}) = \mathcal{S}(c_j + (\frac{\boldsymbol{v}}{\sigma^2})^T \mathbf{W}(\boldsymbol{s})_{:j})$$
 (5)

$$p(v_i = v | \boldsymbol{h}, \boldsymbol{s}) = \mathcal{N}(v | b_i + \mathbf{W}(\boldsymbol{s})_{i:} \boldsymbol{h}, \sigma_i^2)$$
 (6)

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

$$\frac{\partial \mathcal{L}_{A}}{\partial \bar{W}_{ij}} = \langle \sum_{l,k} \frac{A_{lik} v_{l} h_{j} s_{k}}{\sigma_{l}^{2}} \rangle_{data} - \langle \sum_{l,k} \frac{A_{lik} v_{l} h_{j} s_{k}}{\sigma_{l}^{2}} \rangle_{model}$$

$$\frac{\partial \mathcal{L}_{\mathbf{A}}}{\partial A_{ii'k}} = \langle \sum_{m} \frac{W_{i'm} v_{i} h_{m} s_{k}}{\sigma_{i}^{2}} \rangle_{data} - \langle \sum_{m} \frac{W_{i'm} v_{i} h_{m} s_{k}}{\sigma_{i}^{2}} \rangle_{model}$$
(8)

$$\frac{\partial \mathcal{L}_{A}}{\partial B_{ijk}} = \langle v_i h_j s_k \rangle_{data} - \langle v_i h_j s_k \rangle_{model}, \tag{9}$$

と計算できる.ただし, $\langle \cdot \rangle_{
m data}$ と $\langle \cdot \rangle_{
m model}$ はそれぞれ,観測データ,モデルデータの期待値を表す.他のパラメータ b , σ , c に関しては,従来の m RBM の更

新式 [10] と同様である . 適応型 RBM においても CD (Contrastive Divergence) 法を適用することができるため , 各偏微分値の第二項 $\langle \cdot \rangle_{\mathrm{model}}$ を観測データの 再構築値 $\langle \cdot \rangle_{\mathrm{recon}}$ 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 声質変換への応用

適応型 RBM を声質変換へ応用する場合,Fig. 2 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 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 (8)(9) 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 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).

$$\hat{\boldsymbol{h}} = \operatorname*{argmax}_{\boldsymbol{h}} p(\boldsymbol{h}|\boldsymbol{v}_s, \boldsymbol{s}_s) = \mathcal{S}(\boldsymbol{c} + (\frac{\boldsymbol{v}_s}{\boldsymbol{\sigma}^2})^{\mathrm{T}} \mathbf{W}(\boldsymbol{s}_s))$$
(10)

ただし, 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 s の長さを S+2 へ拡張し,A, B をモード 3 に沿ってそれぞれ A', B' を追加するものとする.式 (10) を書き直すと,

$$\hat{\boldsymbol{h}} = \mathcal{S}(\boldsymbol{c} + (\frac{\boldsymbol{v}_s}{\boldsymbol{\sigma}^2})^{\mathrm{T}}(\mathbf{A}_s \bar{\mathbf{W}} + \mathbf{B}_s))$$
 (11)

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(11)は,一度適応型 \mathbf{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 $\mathbf{S}+2$ 要素のみ $\mathbf{1}$,他が $\mathbf{0}$ 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 \mathbf{S} (Step 4).

$$v_t = \underset{\boldsymbol{v}}{\operatorname{argmax}} p(\boldsymbol{v}|\hat{\boldsymbol{h}}, \boldsymbol{s}_t)$$
$$= \boldsymbol{b} + (\mathbf{A}_t \bar{\mathbf{W}} + \mathbf{B}_t)^{\mathrm{T}} \hat{\boldsymbol{h}}$$
(12)

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 (11)(12) にもあるように,入力話者の音響特徴量 v_t へ変換す

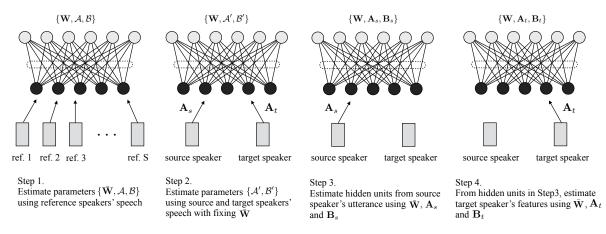


Fig. 2 Procedure of voice conversion using an ARBM.

Table 1 Performance of our 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 |

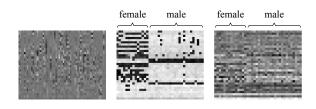


Fig. 3 Left to right: estimated $\bar{\mathbf{W}}$, \mathcal{A} , and \mathcal{B} .

る際 , h の推定に非線形関数を用いているため , 提案 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, $\bar{\mathbf{W}}$ 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 A, \mathcal{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 $TIMIT^1$ を用いて,提案手法

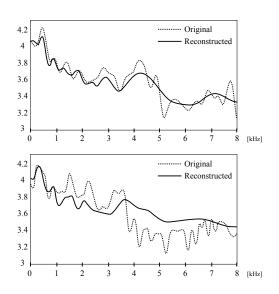


Fig. 4 Log-spectrum from a source speaker and the reconstructed spectrum (above), and log-spectrum from a target speaker and the converted spectrum from the source speech to the target speech (below).

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

¹https://catalog.ldc.upenn.edu/LDC93S1

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 Table 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. "avg." は全組 み合わせの平均値である. Table 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-tomale, 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が男性側へ強く反応 していることからも過学習が窺える).

提案法によって,実際に推定されたパラメータ $\bar{\mathbf{W}}$,A および \mathcal{B} の一部を $\mathrm{Fig.}$ 3 に示す.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(\mathcal{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,A (または \mathcal{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 Fig. 4 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 , $\mathrm{FCJF0}$ の適応型 RBM によって $\hat{m{h}}$ 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, $\hat{m h}$ の推定後 ${
m FCJF0}$ の適応パラ メータによっ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.

参考文献

- [1] C. Veaux and X. Robet, "Intonation conversion from neutral to expressive speech," Interspeech, pp. 2765– 2768, 2011.
- [2] K. Nakamura et al., "Speaking-aid systems using GMM-based voice conversion for electrolaryngeal speech," Speech Commun., vol. 54, no. 1, pp. 134– 146, 2012.
- [3] Y. Stylianou et al., "Continuous probabilistic transform for voice conversion," IEEE Trans. Speech and Audio Process., vol. 6, no. 2, pp. 131–142, 1998.
- [4] T. Toda et al., "Voice conversion based on maximum-likelihood estimation of spectral parameter trajectory," IEEE Trans. Audio, Speech, and Lang. Process., vol. 15, no. 8, pp. 2222–2235, 2007.
- [5] D. Saito et al., "Application of matrix variate Gaussian mixture model to statistical voice conversion," Interspeech, pp. 2504–2508, 2014.
- [6] A. Mouchtaris et al., "Nonparallel training for voice conversion based on a parameter adaptation approach," IEEE Trans. Audio, Speech, and Lang. Processs, vol. 14, no. 3, pp. 952–963, 2006.
- [7] 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.
- [8] T. Toda et al., "Eigenvoice conversion based on Gaussian mixture model," Interspeech, pp. 2446– 2449, 2006.
- [9] D. Saito et al., "One-to-many voice conversion based on tensor representation of speaker space," Interspeech, pp. 653–656, 2011.
- [10] K. Cho et al., "Improved learning of Gaussian-Bernoulli restricted Boltzmann machines," ICANN, pp. 10–17, 2011.