Factored 3-Way Restricted Boltzmann Machine を用いたマルチモーダル音 声認識の検討\*

高島悠樹 (神戸大), 中鹿亘 (電通大), 滝口哲也 (神戸大/JST さきがけ), 有木康雄 (神戸大)

# 1 はじめに

現在,我が国の障害者手帳を持つ18歳以上の人口 は350万人を超えており,聴覚・言語障害者の数は 36万人とされている[1].文献[2]では,構音障害者 音声を対象とした音響モデル適応の検証を行ってい るが,言語障害者などの障害者を対象としている研 究は非常に少ない.本研究は,コミュニケーション手 段として口話を用いる重度難聴者を対象として,音 声と唇形状によるマルチモーダル音声認識を実現し, ユビキタス社会における彼らの生活の支援をするこ とを目的としている.

人間は発話内容を理解する際,種々の情報を統合的 に利用している.音声が聞き取り難い場合,発話者の 顔,特に唇の動きに注目して発話内容を理解しよう とし,逆に,唇の動きと音声が不一致の場合,唇の動 きに影響されて発話内容を誤って理解してしまうこ ともある.これは,McGurk effect(マガーク効果)と 呼ばれ,音韻知覚が音声の聴覚情報のみで決まるの ではなく,唇の動きといった視覚情報からも影響を受 けることが報告されている [3].このように人間によ る発話内容の理解には,唇の画像と音声の情報の統 合的利用が極めて重要である.

唇の動きからの発話内容の読み取りは、リップリー ディング(読唇)と呼ばれ、聴覚障害者にとって重要 なコミュニケーション手段の一つである、リップリー ディングは、背景雑音に影響されることがないため、 計算機上での実現が期待されている、例えば、監視カ メラに収録された会話映像のように音声が聞き取り にくい場合であっても、リップリーディングであれば 発話内容の分析が可能であり、犯罪の防止や抑止に繋 がると考えられる、そのため、音声の雑音に対して頑 健な発話認識を行う手法の一つとして、音声情報に唇 動画像情報を併用して認識を行うマルチモーダル音 声認識が注目され、研究が進められている [4, 5, 6].

重度難聴者は耳で音を聞くことができないため,正 確な発音をすることが難しく,発話スタイルが健常 者と異なる.彼らのコミュニケーション手段の一つ として口話があり,訓練により意図した発話の唇の 形状を作ることが可能である.そこで,彼らの音声 を認識するために,唇画像を併用した音声認識シス テムの構築が望まれる.重度難聴者を対象としたマ ルチモーダル音声認識として, CNN (convolutional neural network)を用いた手法 [7] が提案されている. しかし,この手法はニューラルネットワークを用いて いるためパラメータ数が多く,比較的大量のデータを 必要とする.また,音声と唇画像のモダリティギャッ プが考慮されていない.

本研究では,音声と唇画像を統合した特徴量抽 出のために,factored 3-way restricted Boltzmann machine (F3WRBM [8])を用いる.このモデルは RBM (restricted Boltzmann machine [9])を拡張し たものであり,エネルギー関数に基づく確率モデル である.近年,RBM を用いた特徴抽出(特徴学習) 法 [8,10,11]が多く提案されており,物体認識や音 素認識に応用されている.F3WRBM は2つの観測 変数と1つの潜在変数からなるモデルであり,音響 特徴量と画像特徴量を観測変数とすることで,音声 と唇画像を相補的に考慮できる特徴が潜在変数とし て現れると期待できる.

# 2 Restricted Boltzmann Machine

RBM は, Fig. 1 左図に示すように,可視素子 $v \in \mathbb{R}^{D}$ と隠れ素子 $h \in \mathbb{B}^{H}$ ( $\mathbb{B}$ は0または1のみを取り得る空間)からなる無向グラフィカルモデルである.入力として連続値を定義した IGB (Improved Gaussian-Bernoulli)-RBM [12](以下,この IGB-RBM を単にRBM とする)では,その同時確率とエネルギー関数は以下の式で表される.

$$p(\boldsymbol{v}, \boldsymbol{h}) = \frac{1}{Z} e^{-E_{RBM}(\boldsymbol{v}, \boldsymbol{h})}$$
$$E_{RBM}(\boldsymbol{v}, \boldsymbol{h}) = \sum_{i} \frac{(v_i - b_i)^2}{2\sigma_i^2} - \sum_{i,j} \frac{v_i}{\sigma_i^2} W_{ij} h_j - \sum_{j} c_j h_j$$
(1)

ここで,  $[b_1, \dots, b_D] \in \mathbb{R}^D$ ,  $[c_1, \dots, c_H] \in \mathbb{R}^H$ ,  $\{W_{ij}\}_{D \times H} \in \mathbb{R}^{D \times H}$ ,  $[\sigma_1^2, \dots, \sigma_D^2] \in \mathbb{R}^D$  はそれぞれ,可視素子バイアス,隠れ素子バイアス,可視層-隠れ層間の結合重み,可視素子の分散を表し,いずれも推定すべきパラメータである.i,j はそれぞれ,可視素子及び隠れ素子のインデクスである.ここで,

<sup>\*</sup>Multi-modal speech recognition using factored 3-way restricted Boltzmann machine, by Yuki Takashima (Kobe University), Toru Nakashika (UEC), Tetsuya Takiguchi (Kobe University/JST PRESTO), Yasuo Ariki (Kobe University)



Fig. 1 Graphical representation of an RBM (left) and factored 3-way RBM (right)

 $Z = \int^D \Sigma_h e^{-E(\boldsymbol{v}, \boldsymbol{h})} d^D \boldsymbol{v}$ は全域での確率を1にするための正規化項である.

RBM では,可視素子間,及び隠れ素子間の接続は 存在せず,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は以下のよう な単純な式で表現される.

$$p(v_i = v | \boldsymbol{h}) = \mathcal{N}(v | b_i + \sum_j h_j W_{ij}, \sigma_i^2) \qquad (2)$$

$$p(h_j = 1 | \boldsymbol{v}) = \mathcal{S}(c + \sum_i W_{ij} \frac{v_i}{\sigma_i^2})$$
(3)

ここで,  $\mathcal{N}(\cdot|\mu, \sigma^2)$  は平均  $\mu$ , 分散共分散  $\sigma^2$  の正規 分布,  $\mathcal{S}(\cdot)$  は要素ごとのシグモイド関数を表す.

RBM の各パラメータは, N 個の観測データを  $\{v_n\}_{n=1}^N$  とするとき,この確率変数の対数尤度  $\mathcal{L} = \log \prod_n p(v_n)$ を最大化するように推定される.この 対数尤度をパラメータ  $\theta$  で偏微分すると,

$$\frac{\partial \mathcal{L}}{\partial \theta} = \langle \frac{\partial E(\boldsymbol{v}, \boldsymbol{h})}{\partial \theta} \rangle_{\text{data}} - \langle \frac{\partial E(\boldsymbol{v}, \boldsymbol{h})}{\partial \theta} \rangle_{\text{model}}, \quad (4)$$

が得られる.ここで、 $\langle \cdot \rangle_{data} \geq \langle \cdot \rangle_{model}$ 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、後 者は一般に計算困難なため Contrastive Divergence 法 [13] を用いて求められる.各パラメータは式(4) から、確率的勾配法 (SGD)を用いて繰り返し更新さ れる.

## 3 提案手法

F3WRBM はバイナリな可視素子を仮定し定義されている [8].本節では,可視素子が正規分布に従う と仮定し,モデルを再定義する.さらに,F3WRBM を用いた特徴量抽出と,音声認識への応用について 述べる.

### 3.1 Factored 3-way RBM

RBM は 2 つの確率変数からなるモデルであ るが,これは一般的により高次なモデルへと拡張 が可能である [14].本研究では,音響特徴量を表 す $s = [s_1, \dots, s_D] \in \mathbb{R}^D$ ,画像特徴量を表す  $v = [v_1, \dots, v_L] \in \mathbb{R}^L$ ,潜在特徴量を表すh =  $[h_1, \dots, h_H] \in \mathbb{B}^H$ の3変数間の関係性を3-way RBM を用いて表現する.このとき,3つの確率変数の結合 を表現するエネルギー関数は以下のように定義される.

$$E_{3WRBM}(\boldsymbol{s}, \boldsymbol{v}, \boldsymbol{h}) = -\sum_{i,j,k} s_i v_j h_k W_{ijk} \qquad (5)$$

ここで,*i*,*j*,*k*はそれぞれ*s*,*v*,*h*のインデクスである. $\{W_{ijk}\}_{D \times L \times H} \in \mathbb{R}^{D \times L \times H}$ は3階のテンソルで表現されており,非常に膨大なパラメータ数となる.そこで,3つの確率変数の間に"ファクター"と呼ばれる概念を設け,結合行列を $W_{ijk} = \sum_f C_{if}K_{jf}P_{kf}$ と近似する.ここで, $\{C_{if}\}_{D \times F} \in \mathbb{R}^{D \times F}$ , $\{K_{jf}\}_{L \times F} \in \mathbb{R}^{L \times F}$ ,  $\{P_{kf}\}_{H \times F} \in \mathbb{R}^{H \times F}$ はそれぞれ,音響特徴量-ファクター,画像特徴量-ファクター,ファクター-潜在特徴量の間の結合行列であり,*F*はファクターの数である.これにより,パラメータの数を削減することができ,式(5)のエネルギー関数は以下のように表現される.

$$-E_{F3WRBM}(\boldsymbol{s}, \boldsymbol{v}, \boldsymbol{h})$$
  
=  $\sum_{f} (\sum_{i} s_{i}C_{if}) (\sum_{j} v_{j}K_{jf}) (\sum_{k} h_{k}C_{kf})$  (6)

音響特徴量及び画像特徴量をファクター空間へ射影 してから結合するため,このファクター空間が2つ のモダリティギャップを埋める空間の役割を果たすこ とが期待される.(Fig.1右図)さらに,本稿では音 響特徴量と画像特徴量の確率変数が正規分布に従う と仮定し,エネルギー関数を以下のように定義する.

$$-E(\mathbf{s}, \mathbf{v}, \mathbf{h}) =$$

$$-\sum_{i} \frac{(s_i - b_i)^2}{2\sigma_i^2} - \sum_{j} \frac{(v_j - d_j)^2}{2\delta_j^2} + \sum_{k} c_k h_k +$$

$$\sum_{f} \left(\sum_{i} \frac{s_i}{\sigma_i^2} C_{if}\right) \left(\sum_{j} \frac{v_j}{\delta_j^2} K_{jf}\right) \left(\sum_{k} h_k C_{kf}\right)$$

ここで,  $[b_1, \dots, b_D] \in \mathbb{R}^D$ ,  $[d_1, \dots, d_D] \in \mathbb{R}^L$  はそ れぞれ*s*, *v* に関するバイアス,  $[\sigma_1^2, \dots, \sigma_D^2] \in \mathbb{R}^D$ ,  $[\delta_1^2, \dots, \delta_L^2] \in \mathbb{R}^L$  はそれぞれ*s*, *v* に関する分散を表 す.RBM と同様に,それぞれの変数の要素同士には 結合が存在しないと仮定している.そのため,*s*,*v*, *h* の条件付き確率はそれぞれ以下のように簡単に計 算することができる.

$$p(s_{i} = s | \boldsymbol{v}, \boldsymbol{h})$$

$$= \mathcal{N}(s | b_{i} + \sum_{f} C_{if}(\sum_{j} \frac{v_{j}}{\delta_{j}^{2}} K_{jf})(\sum_{k} h_{k} P_{kf})s, \sigma_{i}^{2})$$

$$p(v_{j} = v | \boldsymbol{s}, \boldsymbol{h})$$

$$= \mathcal{N}(v | d_{i} + \sum_{f} K_{jf}(\sum_{i} \frac{s_{i}}{\sigma_{i}^{2}} C_{if})(\sum_{k} h_{k} P_{kf}), \delta_{j}^{2})$$
(8)

$$p(h_k = 1 | \boldsymbol{s}, \boldsymbol{v})$$
  
=  $\mathcal{S} \left( c_k + \sum_f P_{kf} \left( \sum_i \frac{s_i}{\sigma_i^2} C_{if} \right) \left( \sum_j \frac{v_j}{\delta_j^2} K_{jf} \right) \right)$  (10)

3.2 パラメータ推定

提案モデルのパラメータは,同期の取れたNフレームの音声データと画像データに対する対数尤度

$$\mathcal{L}' = \log \prod_{n} p(\boldsymbol{s}_{n}, \boldsymbol{h}_{n}) = \sum_{n} \log \sum_{\boldsymbol{h}} p(\boldsymbol{s}_{n}, \boldsymbol{h}_{n}, \boldsymbol{v}_{n})$$
(11)

を最大化するように同時に推定することが可能である、本研究では確率的勾配法を用いて各パラメータを更新する.式(11)を各パラメータθに関して偏微分を計算すると以下の式が得られる.

$$\frac{-\partial E(\boldsymbol{s}, \boldsymbol{v}, \boldsymbol{h})}{\partial C_{i'f'}} = \frac{s_{i'}}{\sigma_{i'}^2} \left(\sum_j \frac{v_j}{\delta_j^2} K_{jf'}\right) \left(\sum_k h_k P_{kf'}\right)$$
$$\frac{-\partial E(\boldsymbol{s}, \boldsymbol{v}, \boldsymbol{h})}{\partial P_{k'f'}} = h_{k'} \left(\sum_i \frac{s_i}{\sigma_i^2} C_{if'}\right) \left(\sum_j \frac{v_j}{\delta_j^2} K_{jf'}\right)$$
$$\frac{-\partial E(\boldsymbol{s}, \boldsymbol{v}, \boldsymbol{h})}{\partial b_{i'}} = \frac{s_{i'} - b_{i'}}{\sigma_{i'}^2}$$
$$\frac{-\partial E(\boldsymbol{s}, \boldsymbol{v}, \boldsymbol{h})}{\partial c_{k'}} = h_{k'}$$

 $K_{j'f'}$ ,  $d_{j'}$ に関しても同様である.分散パラメータは 学習の安定化と非負制約のため,  $\sigma_i^2 = e^{z_i}$ ,  $\delta_j^2 = e^{p_j}$ とおき,それぞれ $z_i$ ,  $p_j$ で偏微分しパラメータの更 新を行う.

$$\frac{-\partial E(\boldsymbol{s}, \boldsymbol{v}, \boldsymbol{h})}{\partial z_{i'}} = e^{-z_{i'}} \left( \frac{1}{2} (s_{i'} - b_{i'})^2 - s_{i'} \left( \sum_{f} C_{i'f} (\sum_{j} \frac{v_j}{\delta_j^2} K_{jf}) (\sum_{k} h_k P_{kf}) \right) \right)$$

 $p_j$ に関しても同様である.本稿では,h, v, sを順に サンプリングする.例えばhのサンプリングでは,式 (10)を用いて, $\tilde{h} \sim p(h|s, v)$ とすることでサンプル  $\tilde{h}$ を得る.このようにすることで,既知の特徴量s, vから $\tilde{h} \sim p(h|s, v), \tilde{v} \sim p(v|s, \tilde{h}), \tilde{s} \sim p(s|\tilde{v}, \tilde{h}), \tilde{h} \sim$  $p(h|\tilde{s}, \tilde{v}) \cdots$ とGibbsチェインを繋げていくことがで きる.

### 3.3 F3WRBM による特徴量抽出

本研究では,F3WRBMを用いて音声情報と画像情 報を統合した特徴量の抽出を行う.具体的には,音響 特徴量 s と画像特徴量 v から式 (10)を用いて潜在特 徴量 h を推定し,これを音響モデル HMM (hidden Marcov model)の入力特徴量とする.一般的なマル チモーダル音声認識手法として,結果統合やマルチ ストリーム HMM が用いられているが,統合時の重 みの調節や計算コストが大きいという問題点がある. 本手法では,特徴量抽出器において2つのモダリティ の統合を行うため,音響モデルは単純な構成で済み 計算コストを抑えることができる.また,パラメータ を教師なし学習により学習するため,音声と画像の 相補的な関係が自動的に学習できると期待される.

### 4 評価実験

#### 4.1 実験条件

データセットとして,重度難聴者の男性1名の音声 及び唇画像を収録し使用した.発話内容は,ATR 音素 バランス単語Aセット [15]から選択した.F3WRBM 及び音響モデルの学習データとして2620単語,評価 データとして216単語を用いた.重度難聴者の発話 スタイルは,健聴者の発話スタイルと大きく異なる ため,特定話者モデルにより認識を行う.音声の標本 化周波数は16kHz,語長16bitであり,音響分析には Hamming窓を用いた.STFTにおけるフレーム幅, シフト幅はそれぞれ25ms,5msである.本稿で用い る音響モデルは,54音素のmonophone-HMMで,各 HMMの状態数は5,状態あたりの混合分布数は6で ある.

F3WRBM への入力音響特徴量として 39 次元のメ ル周波数スペクトル,入力画像特徴量として唇画像に 対して DCT を行い上位 30 次元を用いた.また,ファ クター数を 96,潜在特徴量の数を 196 とした.学習 率 0.001,モーメント係数 0.9,バッチサイズ 128,繰 り返し回数 500 の確率的勾配法を用いてモデルを学 習した.

#### 4.2 実験結果と考察

Fig. 2 にマルチモーダル音声認識結果を示す. ベースラインとして,従来の音響特徴量 MFCC(12 次元)+AMFCC と画像特徴量 DCT(30 次元) を 連結して HMM の入力とする初期統合を行った ("MFCC+△+DCT").また, RBM を用いた単純な 統合方法として,メル周波数スペクトル(39次元) と DCT(30 次元)を連結して RBM を学習し,潜在 特徴量を新たな特徴量として認識する手法を行った ("RBM"). この RBM は潜在特徴量の数を 128 とし, PCA (principal component analysis) を行い60次元 に圧縮して認識に用いた."F3WRBM"が提案手法に より抽出した特徴量を用いたものであり , これも潜在 特徴量に対し PCA を行い, 60 次元に圧縮して用い ている. Fig. 2 より,提案手法は他の手法より精度が 大きく劣化していることが分かる.これは,提案モ デルが音声と唇画像を対等に扱っているためだと考 えられる.一般に,クリーン環境下においては,音響 特徴量と唇画像特徴量を用いて認識を行うと,音響



Fig. 2 Word recognition accuracies of each feature.

特徴量のみを用いて認識を行う場合と比べて精度は 劣化し, 唇画像は逆にノイズとして影響してしまう. しかし, 唇画像も認識に好影響な特徴を持っている はずであり, 提案手法はその相補的な関係を学習し, 特徴量に反映できると期待したが,本稿のモデルでは その効果は見られなかった.原因として,学習された 3つの結合重みの学習の程度にばらつきがあり,全体 として十分な学習が行えていないことが考えられる. そこで,各結合重みに対して,何らかの正規化を制約 として加えることで,安定した学習が行えるのでは ないかと考えられる.

次に, MFCC と提案手法で得られた特徴 量("F3WRBM") を連結して認識を行った ("MFCC+F3WRBM").動的特徴量を含むべー スラインよりも良い精度が得られていることから, 提案手法により得られた特徴量は従来の特徴量には ない特徴が表現されていると考えられる.

## 5 おわりに

本稿では,factored 3-way RBM を用いて,音声特 徴量と画像特徴量を統合した特徴量抽出法を提案し た.重度難聴者の男性1名の音声と唇画像を用いた マルチモーダル音声認識実験により,提案手法の評価 を行った.教師なし学習により,音声と画像の相補的 な関係のモデル化を期待したが,評価実験において, 提案手法で得られた特徴量のみを用いた場合は良い 精度が得られず,効果を確認できなかった.本研究で は,パラメータの学習時に制約を入れていないため, 3 つある結合行列のうち,ファクターと潜在変数間の 結合行列に学習が偏る傾向が見られた.今後は,それ ぞれの結合行列がバランス良く学習できる正規化を 導入し,精度の改善を図りたい.

## 参考文献

- [1] 内閣省, "平成 25 年版障害者白書,".
- [2] 中村圭吾 et al., "発話障害者音声を対象にした 健常者音響モデルの適応と検証,"日本音響学会

講演論文集, pp. 109–110, 2015.

- [3] H. McGurk and J. MacDonald, "Hearing lips and seeing voices," Nature, vol. 264, pp. 746– 748, 1976.
- [4] G. Potamianos *et al.*, "Audio-visual automatic speech recognition: An overview," Issues in visual and audio-visual speech processing, vol. 22, pp. 23, 2004.
- [5] G. Potamianos *et al.*, "Recent advances in the automatic recognition of audiovisual speech," Proceedings of the IEEE, vol. 91, no. 9, pp. 1306–1326, 2003.
- [6] Y. Mroueh *et al.*, "Deep multimodal learning for audio-visual speech recognition.," in *ICASSP*, 2015, pp. 2130–2134.
- [7] Y. Takashima *et al.*, "Audio-visual speech recognition using convolutive bottleneck networks for a person with severe hearing loss," IPSJ Transactions on Computer Vision and Applications, vol. 7, pp. 64–68, 2015.
- [8] M. Ranzato *et al.*, "Factored 3-way restricted boltzmann machines for modeling natural images.," in *AISTATS*, Y. W. Teh and D. M. Titterington, Eds. 2010, vol. 9 of *JMLR Proceedings*, pp. 621–628, JMLR.org.
- [9] Y. Freund and D. Haussler, "Unsupervised learning of distributions on binary vectors using two layer networks," Tech. Rep., 1994.
- [10] G. E. Dahl *et al.*, "Phone recognition with the mean-covariance restricted boltzmann machine.," in *NIPS*, J. D. Lafferty *et al.*, Eds. 2010, pp. 469–477, Curran Associates, Inc.
- [11] M. Ranzato and G. E. Hinton, "Modeling pixel means and covariances using factorized thirdorder boltzmann machines.," in *CVPR*. 2010, pp. 2551–2558, IEEE Computer Society.
- [12] A. L. K. Cho and T. Raiko, "Improved learning of Gaussian-Bernoulli restricted Boltzmann machines," in Artificial Neural Networks and Machine Learning, 2011, pp. 10–17.
- [13] G. E. Hinton *et al.*, "A fast learning algorithm for deep belief nets," Neural Comput., vol. 18, no. 7, pp. 1527–1554, July 2006.
- [14] T. J. Sejnowski, "Higher-order boltzmann machines," in AIP Conference Proceedings, 1986, vol. 151, pp. 398–403.
- [15] A. Kurematsu *et al.*, "ATR Japanese speech database as a tool of speech recognition and synthesis," Speech Communication, vol. 9, no. 4, pp. 357–363, 1990.