

背景

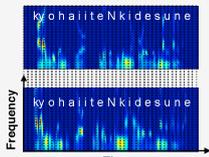
声質変換とは

話している内容はそのままに、音声を実在の人の声に変換する技術。福祉やエンターテインメント用途の他に、喉頭摘出者のコミュニケーション補助等に用いられる。



パラレル型と非パラレル型

モデルの学習にパラレルデータ（同一文言による音声データ対）を用いるアプローチをパラレル声質変換、任意の音声データで学習を行うアプローチを非パラレル声質変換と呼ぶ。パラレル声質変換は比較的高品質な変換が可能であるが、使用する学習データが限られてしまい、利便性が損なわれる。一方、非パラレル声質変換は任意発話で学習できるため利便性・実用性は高いが、品質面で課題がある。



パラレルデータの例

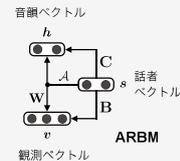
従来手法：ARBMによる声質変換

[T. Nakashika et al., 2016]

ARBM: Adaptive restricted Boltzmann machine

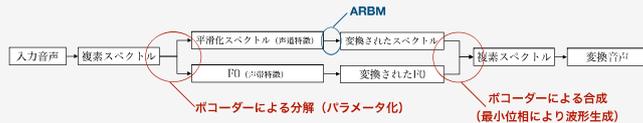
アプローチ

適応型制限ボルツマンマシン（ARBM）を用いた非パラレル声質変換手法。観測ベクトル、潜在的な音韻、話者を表すワンホットベクトルに存在する接続重みを自動学習させる。結果、観測ベクトルが与えられた時、音韻と話者を推定でき、話者を切替ることで声質変換を実現。



考えられる問題点

ボコーダーを使用する必要があり、最小位相の仮定やF0推定エラーで品質が低下してしまう。

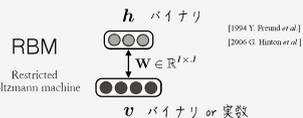


複素RBMの検証

[T. Nakashika et al., 2017, 2018]

提案手法を検証する前に、まず、一般化したモデル（RBMを複素拡張したモデル：複素RBM）について評価し、実験的に有効性を確認した。

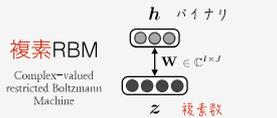
複素RBMと、通常のRBMの違い



$$p(v, h) = \frac{1}{N} e^{-E(v, h)}$$
$$E(v, h) = -v^T W h$$
$$E(v, h) = -\frac{1}{2} v^T \Sigma^{-1} v - v^T \Sigma^{-1} W h$$

$$p(v|h) = \mathcal{N}(v; W h, \Sigma)$$

正規分布 (観測期待値は実数となる)

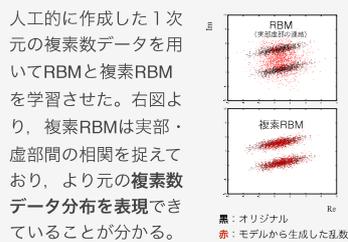


$$p(z, h) = \frac{1}{N} e^{-E(z, h)}$$
$$E(z, h) = \frac{1}{2} \begin{bmatrix} z \\ \bar{z} \end{bmatrix}^H \Phi^{-1} \begin{bmatrix} z \\ \bar{z} \end{bmatrix}$$
$$- \begin{bmatrix} z \\ \bar{z} \end{bmatrix}^H \Phi^{-1} \begin{bmatrix} W \\ \bar{W} \end{bmatrix} h$$

$$p(z|h) = \mathcal{N}_c(z; W h, \Delta(\gamma), \Delta(\delta))$$

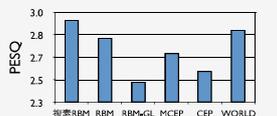
複素正規分布 (観測期待値は複素数となる)

人工データを用いた実験



音声データを用いた実験

200文の音声データを用いて複素RBMを学習させたところ、従来のRBMだけでなくメルケプストラム、ケプストラム、WORLDなど代表的な音声特徴表現手法を上回る結果が得られた。

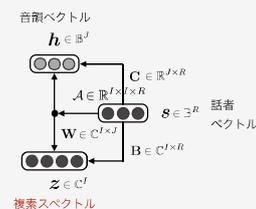


提案手法：複素ARBMによる声質変換

従来ではボコーダーを使用する必要があり、品質低下を招いていた。そこでボコーダーを使用せず、複素拡張モデルによって複素スペクトルを直接変換する手法を提案する。



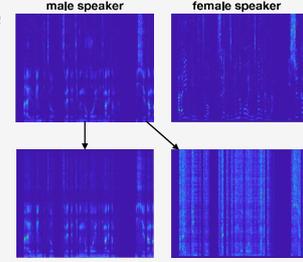
$$p(z, h|s) = \frac{1}{Z} e^{-E(z, h|s)}$$
$$E(z, h|s) = \frac{1}{2} \begin{bmatrix} z \\ \bar{z} \end{bmatrix}^H \Phi^{-1} \begin{bmatrix} z \\ \bar{z} \end{bmatrix} - \begin{bmatrix} b + B s \\ \bar{b} + \bar{B} \bar{s} \end{bmatrix}^H \Phi^{-1} \begin{bmatrix} z \\ \bar{z} \end{bmatrix} - 2(c + C s)^T h$$
$$- \begin{bmatrix} z \\ \bar{z} \end{bmatrix}^H \Phi^{-1} \begin{bmatrix} (A \circ s) W \\ (A \circ s) \bar{W} \end{bmatrix} h$$



- ①入力話者から音韻の推定 $p(h|z, s) = \mathcal{B}(h; \sigma(2\Re\{W^H(A \circ s)^T z\} + 2c + 2Cs))$
- ②目標話者のスペクトル推定 $p(z|h, s) = \mathcal{N}_c(z; (A \circ s) W h + b + B s, \Delta(\gamma), \Delta(\delta))$

評価実験と現状の課題

男女1名それぞれ50文の音声を用いて複素ARBMを学習させ、実際に声質変換を実行したところ、右図のようなスペクトルが得られた。同一話者のリコンストラクションは正しく行われた一方、話者の特徴を表すフォルマントやF0など、変換音声では正しく表現されていないことが分かる。この主な原因として、



課題1. 複素スペクトルをテンプレートの線形変換で表現していることが不十分

課題2. 同じテキストでも潜在変数（音韻）が異なるなどが挙げられる。

$$p(z|h, s) = \mathcal{N}_c(z; (A \circ s) W h + b + B s, \Delta(\gamma), \Delta(\delta))$$

複素スペクトルテンプレートWを行列Aで射影

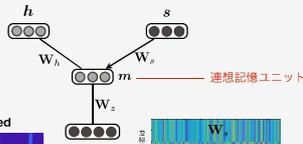
課題1. の解決策 連想記憶モデル

複素ARBMではスペクトルがテンプレートの線形変換で表現される仮定をしているため、F0やフォルマントを正しく表現できなかった。そこで複素ARBMに連想記憶を追加し、複素スペクトルを非線形的に表現できるモデル（連想記憶モデル）について検討した。

複素スペクトルの条件付き確率

$$p(z|m) = \mathcal{N}_c(z; W_z m(z, h, s) + b_z, \Delta(\gamma), \Delta(\delta))$$

どのテンプレートWを使用するかをmで選択



連想記憶による非線形変換により、多少品質が改善された。推定されたWsは一方が0他方が1となる重みが多く、mがそれぞれの話者に関する複素スペクトルテンプレートを選択することを示唆している。一方Whでは、各音韻素子に大きな差異が見られない。

推定されたパラメータ

課題2. の解決策 複素CRBMによるパラレル声質変換

複素ARBMでは潜在的な音韻情報の共通化・差別化が困難であるため、パラレルデータを用いて音韻情報を強制的に一致させた場合の複素スペクトルの声質変換について検討した。CRBMの複素拡張（複素CRBM）を用いてモデル化し、パラレル声質変換実験を行った。

CRBM: Conditional restricted Boltzmann machine

複素CRBM



x: 入力話者の複素スペクトル y: 目標話者の複素スペクトル

モデル学習後は、yの初期値を適当に設定し、

- ① $p(h|x, y) = \mathcal{B}(h; \sigma(2\Re\{V^H x + W^H y\} + 2c))$
 - ② $p(y|x, h) = \mathcal{N}_c(y; U x + W h + b, \Delta(\gamma), \Delta(\delta))$
- を繰り返して目標話者の複素スペクトルyを推定

パラレルデータを使用することで品質が著しく改善し、真の目標話者に近い音声を得られた。また、線型変換と、潜在変数yを用いた非線形変換の組み合わせによる変換も品質改善に貢献。複素ARBMでも何らかの方法で潜在音韻を一致させることができれば十分な声質変換性能を期待できる。



正解に近い変換音声を得られた