

ヒトの聴覚末梢系モデルを用いた音声情報処理

山口大学大学院理工学研究科 システム情報科学研究室

背景と目的

音声処理システム

強力な統計モデルに基づいて音声処理を行っているものがほとんどである

大きな騒音などがある場合には対応できない場合がある

ヒトの聴覚系

環境や距離、媒体が違っても必要な情報のみに注目して会話ができる

様々な環境に対応できる

ヒトの聴覚系のような柔軟性をもつ音声処理システムの開発が望まれている

背景
ヒトの聴覚系のような柔軟性をもつ音声処理システムの開発が望まれている

本研究の目的

- ヒトの聴覚系の一部である聴覚末梢系の数理モデルを構築する
- 聴覚末梢系モデルの音声処理システム(話者識別)への応用を考える

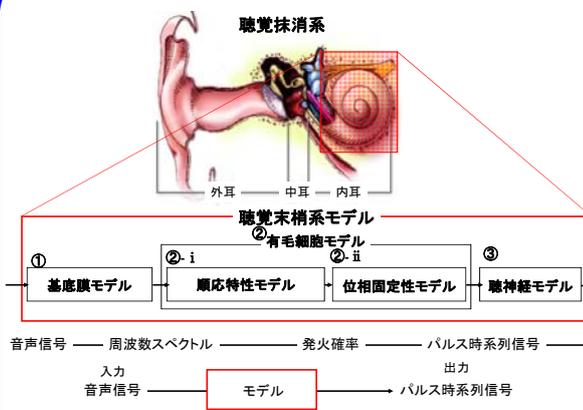
本研究で構築するモデルとその応用

入力が音声信号、出力がパルス時系列信号である聴覚末梢系モデルを構築する

モデルの出力であるパルス時系列信号に基づく特徴量を用いて話者識別を行う

話者識別精度の向上を図るための特徴量変換を行う

聴覚末梢系モデル



① 基底膜モデル

基底膜...周波数ごとに共振部位が異なり、周波数分析機能を有する周波数分析...フィルタバンクで行うことができる

Gammatone フィルタバンクを使用して基底膜のモデル化を行う

音声信号 → g_1 (中心周波数 f_1 のスペクトル) → g_2 (中心周波数 f_2 のスペクトル) → g_N (中心周波数 f_N のスペクトル)

基底膜のモデル化

$$g_n(t) = a_n^2 \exp(-2\beta_n t) \cos(2\pi f_n t + \phi) \quad (n=1,2,\dots,N)$$

FFT

② 有毛細胞モデル

有毛細胞の振動変換機能をMeddis inner hair cell modelで表現する

③ 聴神経モデル

聴神経の膜電位は神経伝達物質により上昇した後、減衰する単純な関数 $t \cdot \exp(-t)$ を用いてモデル化する

時刻 t'_{nj} における膜電位を振幅値 a_{nj} に比例させる

$$a_{nj}(t - t'_{nj}) \cdot \exp\{-t - t'_{nj}\}$$

膜電位のモデル化

$$V(t) = \sum_{j=1}^{J_n} \sum_{n=1}^{N_j} a_{nj} [t - t'_{nj} - h(\tau_{nj})] \exp\{-[t - t'_{nj} - h(\tau_{nj})] T(\tau_{nj})\}$$

膜電位の出力

②-i Meddis inner hair cell model

factory, reprocessing store, cleft, reservoir, free transmitter pool, reprocessing store, cleft, h+c(t)

基底膜振動, シナプス小胞, 神経伝達物質, 順応反応

$$\frac{dq}{dt} = \gamma(1-q(t)) + xw(t) - k(t)q(t)$$

$$\frac{dx}{dt} = k(t)q(t) - l(x) - r(x)$$

$$\frac{dr}{dt} = r(x) - xv(t)$$

②-ii 位相固定性モデル

Meddis inner hair cell modelでは表現できない位相固定性の消失を解決するため、低域通過フィルタを用いず、故らのモデルを導入する

極大点の振幅値(振幅情報)

t'_{nj} : 時間(位相情報)

極大値に隣接する極小点間の時間間隔(周波数情報)

揺らぎを加える

位相固定性の消失をモデル化

$$t'_{nj} = t_{nj} + N(0, f(\tau_{nj}) \cdot s(a_{nj}))$$

$$f(\tau_{nj}) = \tau_{nj} (\omega_1 \tau_{nj}^{\omega_2} + \omega_3)$$

$$s(a_{nj}) = 1 + r_1 e^{-r_2 a_{nj}}$$

有毛細胞モデルの出力

$$\{a_{nj}, t'_{nj}, \tau_{nj}\} \quad (j=1,2,\dots,J_n)$$

パルス生成のモデル化

$$S(t) = \begin{cases} 1 & V(t) \geq U(\alpha, \beta) \\ \text{and } S'(t) = 0 & \text{for } t' \in [t - \tau, t], \quad t' \sim N(\mu, \sigma^2) \\ 0 & \text{otherwise} \end{cases}$$

膜電位モデルの出力信号

検証実験

構築したモデルからヒトの聴覚系に有する聴神経の特性を検証する主な2つを紹介する

マスキング

ある音の最小可聴域が他の音の存在によって上昇する現象

構築したモデルで聴覚系で見られるマスキングを再現できていることがわかる

周波数選択性

ある聴神経はある特定周波数の音に対して発火しやすい現象

構築したモデルで聴覚系で見られる周波数選択性を再現できていることがわかる

話者識別実験

パルス時系列信号からの特徴量抽出

Post-Stimulus Time Histogram (PSTH)

パルス列のラスタ表示において、時間軸上を一定幅の小区間で分割

各区間に入っているパルス数を試行全体にわたって平均化

この平均値を区間の代表値としてヒストグラムを作成

時刻 t における P 次元の特徴量 C_t を用いる

特徴量変換

話者識別率を向上させるため、特徴量変換を行う

- 標準化変換 $C_t \rightarrow \hat{C}_t$
特徴量を、平均値および共分散行列を用いて標準化する。但し、共分散行列は対角成分のみをもつものとして扱う。
- 正規化変換 $\hat{C}_t \rightarrow \hat{\hat{C}}_t$
特徴量をそのノルムで正規化する。
- 標準・正規化変換 $C_t \rightarrow \hat{\hat{C}}_t$
標準化後に正規化変換を行う。

話者識別実験

男性9名、女性3名の母音発話各5回をモデルに入力し、特徴量 C_t を求める

特徴量変換により、それぞれの特徴量 $\hat{C}_t, \hat{\hat{C}}_t, \hat{\hat{\hat{C}}}_t$ を求める

最近隣法を用いて、5 hold leave 2 out cross-validation により話者識別率の評価を行う

比較のためメルLPCスペクトルを特徴量とする話者識別実験も行った

話者識別結果

話者識別率の平均値[%]

	無変換	標準化	正規化	標準・正規化
LPC	80.0 (3.42)	78.1 (3.96)	84.8 (2.28)	80.6 (2.52)
Proposed method	78.1 (3.47)	86.1 (2.12)	79.8 (2.28)	86.6 (3.79)

()は標準偏差を表す

今後の課題

多次元パルス信号から得られる様々な特徴量を用いて、さらなる話者識別率の向上を目指す