

複数のパラメータ共有構造を用いた音響モデリング

塩田 さやか モデリング研究系 特任助教

【研究背景：統計的機械学習に基づく音声認識】

- ◆ 音声媒体としたインタフェースの普及（携帯電話やカーナビゲーションなど）
- ◆ 音声の要素技術への統計的アプローチ

統計的アプローチ

音声合成

音声認識

話者認識

話者適応

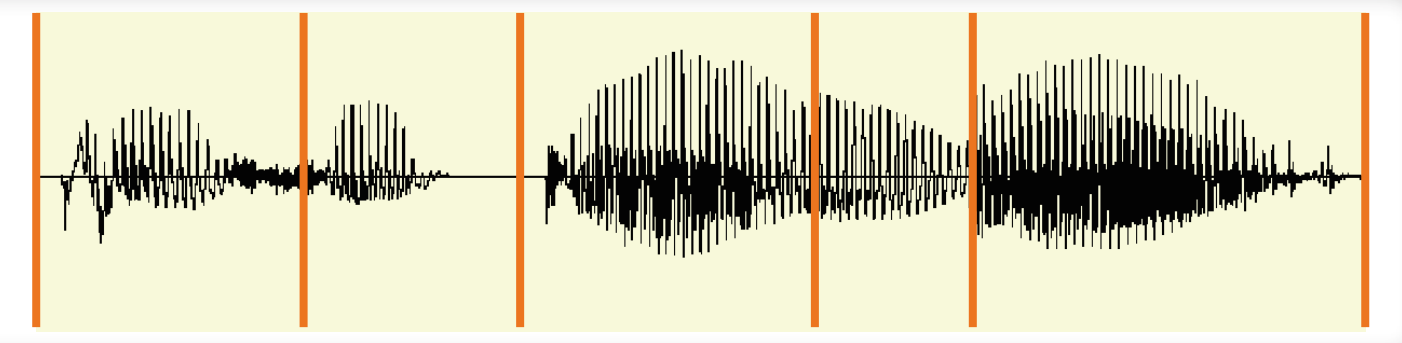
音声翻訳

- ◆ ウェブ上等の膨大な音声データの存在
⇒ データの膨大化、複雑化

- ◆ 音響モデル

音声信号から音の特徴をモデル化したもの

音声信号



様々な要素技術の核となるモデル

⇒ 音声に関する要素技術全てに影響

- ◆ 音響モデルの性能はまだ不十分であるため
更なる性能向上を目指した研究が重要となります

【研究内容：ベイズ基準に基づくモデル構造選択】

- ◆ 概要

- 音響モデルとして広く使われている隠れマルコフモデル(HMM)
 - ・ 学習アルゴリズム: EMアルゴリズム
 - ・ 学習基準: 尤度最大化(ML)基準
⇒ 初期値の影響が大きい改善の必要あり
- 音声認識に必要なモデル構造
 - ・ モデルの複雑化や複数のモデル構造の組み合わせた研究の活発化(e.g., ROVER, random forest)

本研究では、複数のモデル構造を用い**ベイズ基準と確定的アニーリングEMアルゴリズム**を適用することで初期値に依存しにくい高性能な音響モデルの推定を行う。

- ◆ 目的関数の導出

- 複数のモデル構造を隠れ変数として含む周辺尤度関数

$$P(O) = \sum_m \sum_Z \int P(O, Z, m, \Lambda_m) d\Lambda_m$$

$$P(O, Z, m, \Lambda_m) = P(O, Z | m, \Lambda_m) P(\Lambda_m | m) P(m)$$

Z : 状態系列
 O : 観測
 $m \in \{1, \dots, M\}$: モデル構造
 $\Lambda \in \{\Lambda_1, \dots, \Lambda_M\}$: モデルパラメータセット

- ◆ 確定的アニーリングEMアルゴリズムの適用

- 自由エネルギー関数の再定義

$$\bar{F}_\beta = -\frac{1}{\beta} \log \sum_m \sum_Z \int P^\beta(O, Z, m, \Lambda_m) d\Lambda_m \quad \beta: \text{Temperature parameter}$$

- 変分近似を用いることでモデルパラメータ・状態系列・モデル構造の変分事後確率を導出

$$\tilde{Q}(\Lambda_m | m) = C_{\Lambda_m} P^\beta(\Lambda_m | m) \exp \left\langle \log P^\beta(O, Z | m, \Lambda_m) \right\rangle_{\tilde{Q}(Z)}$$

$$\tilde{Q}(Z) = C_Z \exp \left\langle \log P^\beta(O, Z | m, \Lambda_m) \right\rangle_{\tilde{Q}(\Lambda_m | m)} \tilde{Q}(m)$$

$$\tilde{Q}(m) = C_m P^\beta(m) \exp \left\langle \log P^\beta(O, Z | m, \Lambda_m) \right\rangle_{\tilde{Q}(Z)} + \log \frac{P^\beta(\Lambda_m | m)}{\tilde{Q}(\Lambda_m | m)} \right\rangle_{\tilde{Q}(\Lambda_m | m)}$$

- 温度パラメータを0(高温)から1(低温)に変えていきながら各温度で $\tilde{Q}(Z)$, $\tilde{Q}(m)$, $\tilde{Q}(\Lambda_m | m)$ を繰り返し推定していくことで**初期値に依存しにくい高精度なモデルを推定することが可能**

アニーリングの枠組みを用いたモデル学習



$$\text{温度スケジュール} \quad \beta(i) = \left(\frac{i}{I}\right)^\alpha, \quad i = 0, \dots, I, \quad \alpha = 1/8, 1/4, \dots, 4, 8$$

- ◆ 実験結果

