

1. 序論

◆線形予測(Linear Prediction; LP)法による音声分析 [1]

- 声帯駆動信号を白色Gauss性雑音と仮定し、全極モデルによって表される声道伝達特性を観測信号から推定

$$\epsilon_i \longrightarrow \boxed{1/A(z)} \longrightarrow x_i = \sum_{p=1}^P a_p x_{i-p} + \epsilon_i$$

$$A(z) = 1 - a_1 z^{-1} - \dots - a_P z^{-P}$$

$$\epsilon_i \sim \mathcal{N}(0, \sigma^2) \quad (i = 1, \dots, I)$$

問題: 白色性の仮定は高い基本周波数の音声ほど適しなくなる (声道伝達特性の推定結果に影響を及ぼす可能性がある)

◆LP法の改良アプローチ

- 駆動音源信号がラプラス分布に従うと仮定したもの [2]
- 予め抽出した調波成分のスペクトルピークだけに対して全極スペクトルをフィッティングするもの [3]
- 声帯駆動信号を隠れマルコフモデルにより確率モデル化したもの [4]

駆動音源の周期性を陽に仮定し、声道伝達特性とともに声帯駆動信号の周期、周期性の有無についても推定できるLP法が実現できないか?

2. 従来のLP法 [1]

◆音声信号のモデル化

- P次の自己回帰過程: $x_i = \sum_{p=1}^P a_p x_{i-p} + \epsilon_i$

$$\begin{aligned} \epsilon &:= (\epsilon_1, \dots, \epsilon_I)^T \\ \mathbf{x} &:= (x_1, \dots, x_I)^T \\ \Psi &:= \begin{bmatrix} 1 & & & & 0 \\ -a_1 & & & & \\ & \ddots & & & \\ -a_P & & & & \\ 0 & & -a_P & \dots & -a_1 & 1 \end{bmatrix} \end{aligned} \quad \Psi \mathbf{x} = \epsilon \quad \dots (1)$$

- 白色声帯駆動信号の仮定: $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$

$$(1) \text{より } \mathbf{x} | \sigma^2, \mathbf{a} \sim \mathcal{N}(0, \sigma^2 (\Psi^T \Psi)^{-1})$$

$$\mathbf{a} := (a_1, \dots, a_P)^T$$

- パラメータの対数尤度: $|\Psi| = 1$ より

$$\begin{aligned} \log p(\mathbf{x}; \mathbf{a}) &= -\frac{I}{2} \log 2\pi - \frac{1}{2\sigma^2} \mathbf{x}^T \Psi^T \Psi \mathbf{x} \\ &= -\frac{I}{2} \log 2\pi - \frac{1}{2\sigma^2} (\mathbf{x} - \mathbf{X}\mathbf{a})^T (\mathbf{x} - \mathbf{X}\mathbf{a}) \end{aligned}$$

- 最尤の予測係数:

$$\mathbf{X}^T \mathbf{X} \mathbf{a} = \mathbf{X}^T \mathbf{x}$$

(正規方程式)

$$\mathbf{X} := \begin{bmatrix} 0 & \dots & 0 \\ x_1 & \dots & x_1 \\ \vdots & \ddots & \vdots \\ 0 & \dots & 0 \\ \vdots & \ddots & \vdots \\ x_{I-1} & \dots & x_{I-P} \end{bmatrix}$$

3. GPIによる声帯駆動のモデル化

- 基底関数 $\phi_1(t), \dots, \phi_H(t)$ の線形結合+ノイズ:

$$\epsilon = \underbrace{f(t) + \eta}_{\sum_{h=1}^H \phi_h(t) w_h = \phi(t)^T \mathbf{w}}$$

$$\begin{aligned} \epsilon &:= (\epsilon_1, \dots, \epsilon_I)^T \\ \Phi &:= (\phi(t_1), \dots, \phi(t_I))^T \\ \mathbf{w} &:= (w_1, \dots, w_H)^T \\ \eta &:= (\eta_1, \dots, \eta_I)^T \end{aligned} \quad \epsilon = \Phi \mathbf{w} + \eta \quad \dots (2)$$

$$\eta \sim \mathcal{N}(0, \sigma_\eta^2 I)$$

- 結合係数 \mathbf{w} の事前確率: $\mathbf{w} \sim \mathcal{N}(0, \sigma_w^2 I)$

正規分布する結合係数による基底関数の線形和で表される回帰モデルをガウシアンプロセス(GP)という

$$(2) \text{より } \epsilon \sim \mathcal{N}(0, \sigma_w^2 \Phi \Phi^T + \sigma_\eta^2 I)$$

\mathbf{K} の $\{i, j\}$ 成分 $K_{i,j}$ は、 $k(t, t') := \phi(t)^T \phi(t')$ と定義されるカーネル関数を用いて $K_{i,j} = k(t_i, t_j)$ で表される

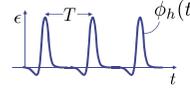
- 白色雑音に帰着する基底関数:

$$\phi_h(t) := \begin{cases} 1 & (t = c_h) \\ 0 & (t \neq c_h) \end{cases} \longrightarrow \subseteq \{c_1, \dots, c_H\} \text{ ならば}$$

$$\mathbf{K} = \mathbf{I}$$

4. 周期カーネルの導入

- 有声音の場合、声帯駆動信号は周期性をもつことが多い
→ 周期 T の基底関数を用いて ϵ をモデル化したい



- 周期カーネルの例 [5]

$$k(t, t') := \exp \left\{ -2 \sin^2 \left(\pi \frac{t-t'}{T} \right) / \ell^2 \right\}$$

- ・正定値性を満たす
- ・周期 T の基底関数を仮定したことに相当

- 本発表では以下のように周期カーネルを定義

$$\phi_h(t) := \sum_{n=1}^N \sin \left(2\pi n \frac{t - c_h}{T} \right) \longrightarrow k(t, t') := \sum_{h=1}^H \phi_h(t) \phi_h(t')$$

5. カーネル線形予測モデル

- 声帯駆動信号の周期が既知の場合のLPモデル

$$(2) \text{を}(1) \text{に代入: } \Psi \mathbf{x} = \Phi \mathbf{w} + \eta$$

周期 T の雑音 η 白色雑音

$$\longrightarrow \mathbf{x} \sim \mathcal{N}(0, \sigma_w^2 \Psi^{-1} \mathbf{K} \Psi + \sigma_\eta^2 (\Psi^T \Psi)^{-1})$$

6. マルチカーネル線形予測モデル

- 声帯駆動の周期は未知 → カーネル自体を推定したい
- マルチカーネル学習の考え方を導入:

- ・複数の周期カーネルの線形結合でカーネル関数を定義

$$k(t, t') = \sum_{m=1}^M \theta_m k_m(t, t'), \quad \sum_{m=1}^M \theta_m = 1$$

- ・周期 T_m の周期カーネル $k_m(t, t')$ の優勢度 θ_m が推定すべき未知パラメータ

- ・ $\mathbf{K} = \sum_{m=1}^M \theta_m \mathbf{K}_m$ の要素にもつ行列

- \mathbf{x} の確率密度関数:

$$\mathbf{x} \sim \mathcal{N} \left(0, \sigma_w^2 \sum_{m=1}^M \theta_m \Psi^{-1} \mathbf{K}_m \Psi^{-T} + \sigma_\eta^2 (\Psi^T \Psi)^{-1} \right)$$

- モデルパラメータ $\Theta := \{\mathbf{a}, \theta, \sigma\}$ の対数尤度:

$$\log p(\mathbf{x}; \Theta) = -\frac{1}{2} \left(I \log 2\pi + \log |\Sigma| + \mathbf{x}^T \Psi^T \Sigma^{-1} \Psi \mathbf{x} \right)$$

$$\Sigma = \sigma_w^2 \sum_{m=1}^M \theta_m \mathbf{K}_m + \sigma_\eta^2 \mathbf{I}$$

- ・ $\theta := \{\theta_m\}_{m=1}^M$ の推定を通して音声信号に含まれる優勢な基本周期を知ることができる
- ・ $\sigma := \{\sigma_w, \sigma_\eta\}$ の推定を通して周期成分・白色成分の混在比を知ることができる

7. スパース正則化パラメータ学習

- θ の事前分布: $p(\theta) = \prod_{m=1}^M \frac{\lambda}{2} \exp \{ -\lambda |\theta_m| \}$ (ラプラス分布)



- 最大事後確率パラメータの反復推定: $\mathbf{a} \leftarrow \underset{\mathbf{a}}{\operatorname{argmax}} \log p(\mathbf{x}; \Theta) + \log p(\theta) \quad \dots (3)$

$$\{\theta, \sigma\} \leftarrow \underset{\theta, \sigma}{\operatorname{argmax}} \log p(\mathbf{x}; \Theta) + \log p(\theta) \quad \dots (4)$$

- ・(3): $\mathbf{X}^T \Sigma^{-1} \mathbf{X} \mathbf{a} = \mathbf{X}^T \Sigma^{-1} \mathbf{x}$

- ・(4): 解析的に求まらないがEM (Expectation-Maximization)法で反復推定できる

$$\begin{aligned} \text{完全データ } \mathbf{y} &:= (\mathbf{y}_1^T, \dots, \mathbf{y}_M^T, \mathbf{y}_{M+1}^T) \longrightarrow \mathbf{H} := [\mathbf{I}, \dots, \mathbf{I}] \\ \mathbf{y}_m &\sim \mathcal{N}(0, \sigma_w^2 \theta_m \Psi^{-1} \mathbf{K}_m \Psi^{-T}), \quad m = 1, \dots, M \\ \mathbf{y}_{M+1} &\sim \mathcal{N}(0, \sigma_\eta^2 (\Psi^T \Psi)^{-1}) \end{aligned}$$

$$\text{完全データ対数尤度: } \log p(\mathbf{y}; \Theta) \doteq -\frac{1}{2} \left(\log |\Lambda| + \mathbf{y}^T \Lambda^{-1} \mathbf{y} \right)$$

$$Q(\Theta, \Theta') = \log p(\theta) - \frac{1}{2} \left(\log |\Lambda| + \operatorname{tr} \left(\Lambda^{-1} \mathbb{E}[\mathbf{y} \mathbf{y}^T | \mathbf{x}; \Theta'] \right) \right) \longrightarrow \text{「Q関数」}$$

$$\begin{aligned} &\doteq -\sum_{m=1}^M \lambda |\theta_m| - \frac{I}{2} \log (\sigma_w^{2M} \sigma_\eta^2 \theta_1 \dots \theta_M) \\ &\quad - \frac{1}{2\sigma_w^2} \sum_{m=1}^M \frac{1}{\theta_m} \operatorname{tr} \left(\Psi^T \mathbf{K}_m^{-1} \Psi \mathbf{R}_m \right) - \frac{1}{2\sigma_\eta^2} \operatorname{tr} \left(\Psi^T \Psi \mathbf{R}_{M+1} \right) \end{aligned}$$

$$\longrightarrow \text{Mステップ: } \theta_m \leftarrow 2\sigma_w^2 \lambda \theta_m^2 + I \sigma_w^2 \theta_m - \operatorname{tr} \left(\Psi^T \mathbf{K}_m^{-1} \Psi \mathbf{R}_m \right) = 0 \text{ の正の根}$$

$$\sigma_w^2 \leftarrow \frac{1}{IM} \sum_{m=1}^M \frac{1}{\theta_m} \operatorname{tr} \left(\Psi^T \mathbf{K}_m^{-1} \Psi \mathbf{R}_m \right) \quad \sigma_\eta^2 \leftarrow \frac{1}{I} \operatorname{tr} \left(\Psi^T \Psi \mathbf{R}_{M+1} \right)$$

8. 動作実験

- 分析条件

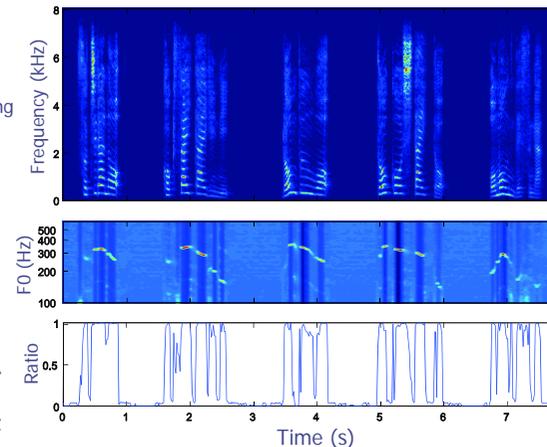
- ・音声データ: ATR音声データ ベースセットの女性話者音声
- ・標準化周波数: 16kHz
- ・フレーム分割: 64ms長のHanning窓, シフト幅32ms

- 提案法の条件

- ・反復回数: 600
- ・予測次数 P : 13
- ・周期カーネルの設定: 100Hz ~ 600Hzを6cent間隔で離散化した計518個の基本周波数値 $1/T_m$ をもとに構成

- 実験結果

- (上) 観測信号のスペクトログラム
- (中) 各時刻における θ の分布
- (下) 各時刻における周期成分と白色成分のエネルギー比



参考文献

- [1] F. Itakura and S. Saito, "Analysis synthesis telephony based upon the maximum likelihood method," In Proc. 6th Int'l Cong. Acoust. (ICA'68), C-5-5, pp. C17-20, 1968.
- [2] C.-H. Lee, "On robust linear prediction of speech," IEEE Trans. Acoust., Speech & Signal Process., 36(5), pp. 642-650, 1988.
- [3] A. El-Jaroudi and J. Makhoul, "Discrete all-pole modeling," IEEE trans. Signal Process., 39(2), pp. 411-423, 1991.
- [4] 佐宗, 田中, "HMMによる音源のモデリングと高基本周波数に頑健な声道特性抽出," 信学論D-II, J84-D-II(9), 1960-1969, 2001.
- [5] C.E. Rasmussen and C.K.I. Williams, Gaussian Processes for Machine Learning, MIT Press, Cambridge, Mass, USA, 2006.