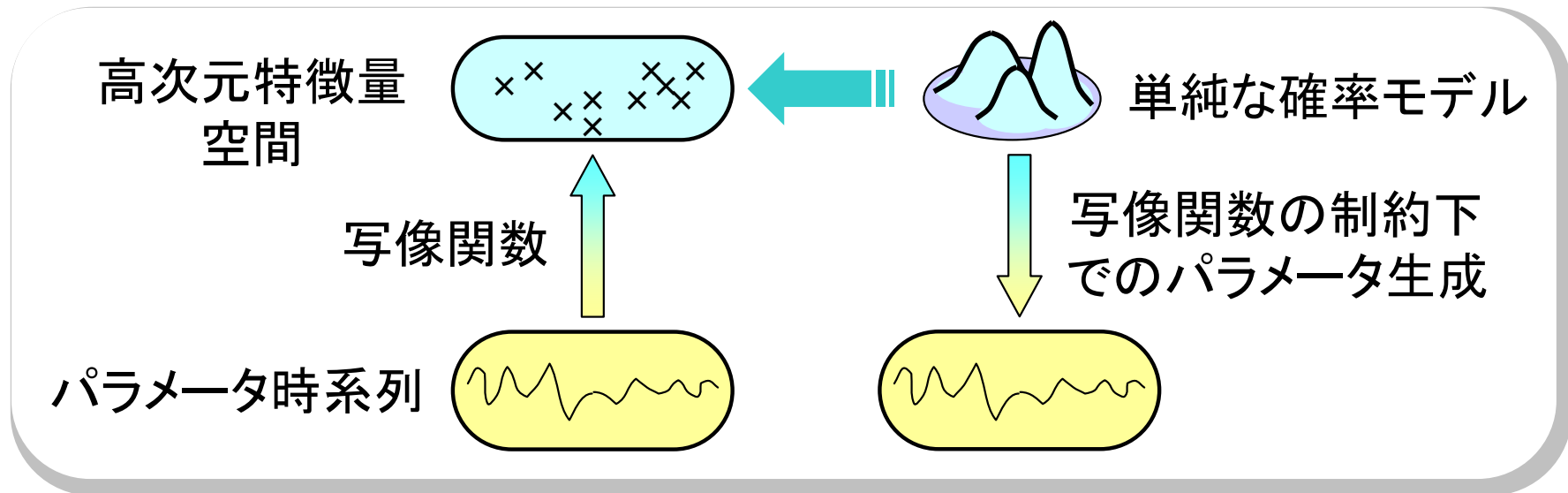


# 音声音響信号処理

## ～統計的手法による音声変換～



戸田 智基

奈良先端科学技術大学院大学

情報科学研究科

2012年1月23日(月)

# 内容

---

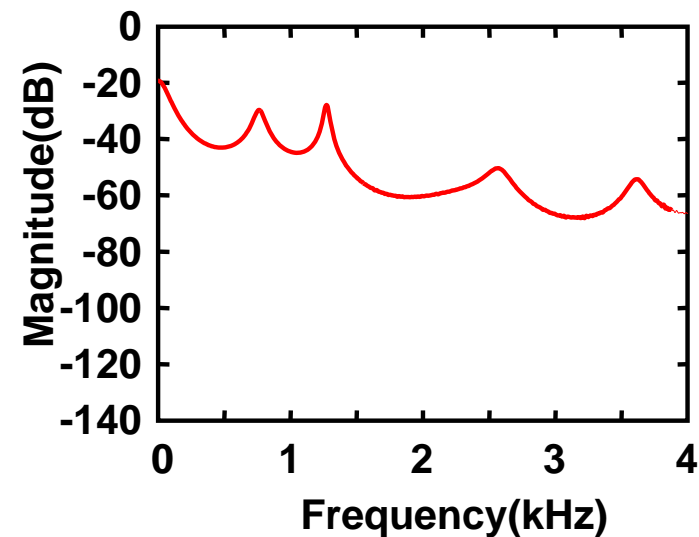
1. 音声変換のしくみ
2. 統計的手法による声質変換
  - 2.1. 基本的な枠組み
  - 2.2. フレームベース変換法
  - 2.3. 系列ベース変換法
3. 応用例

# 音声の特徴量

## 声道特徴量

音韻性や声質を調節する.

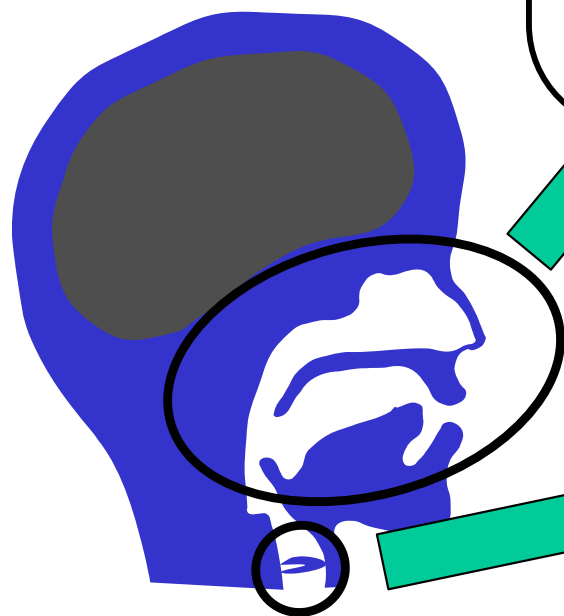
➤ **スペクトル包絡**



## 音源特徴量

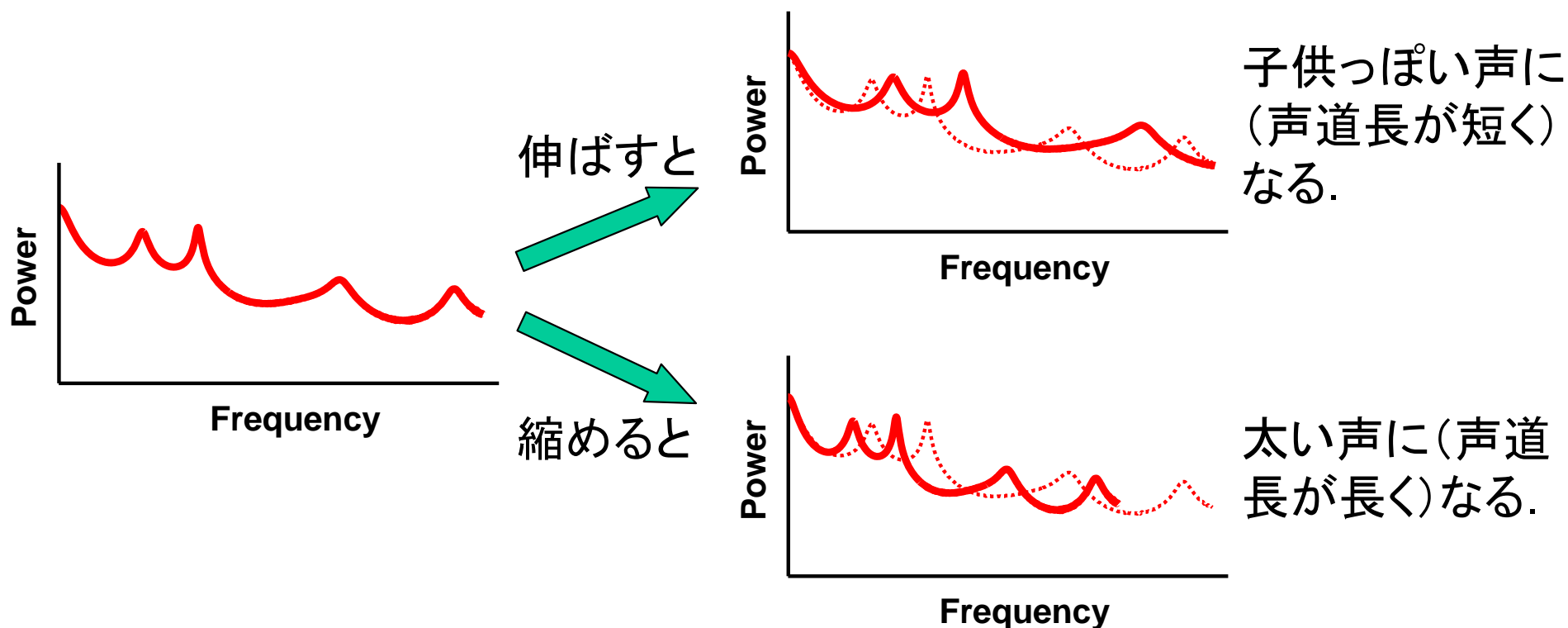
韻律 (イントネーション) を調節する.

➤ **基本周波数 ( $F_0$ )**



# 声道特徴量の変換

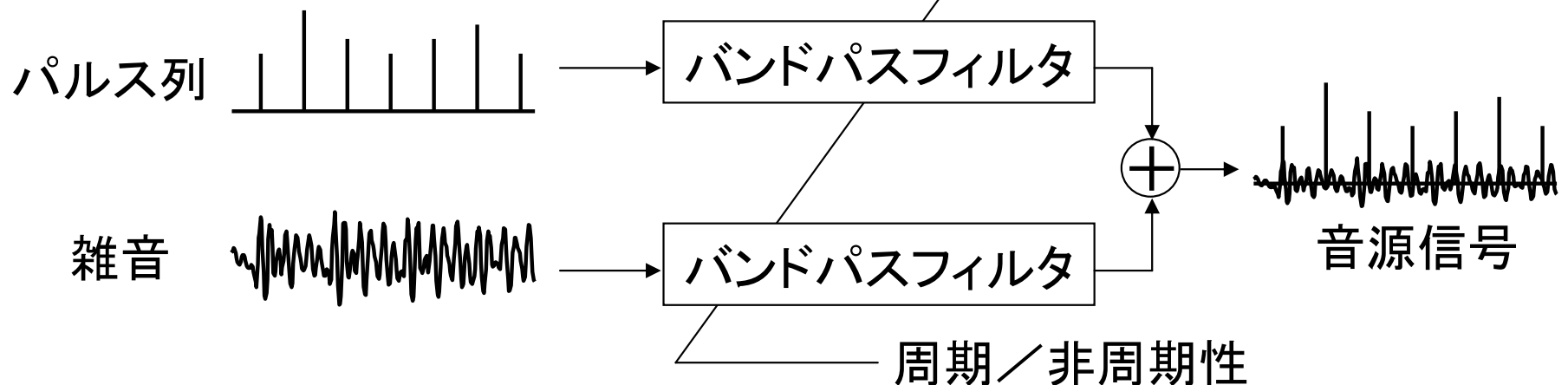
- 特徴量パラメータとして**スペクトル包絡**が用いられる。
- **音韻性**や**声質**などを表す。
- 周波数軸方向に一律に伸縮させることで、音韻性を保ったまま声質を変換することができる。



# 音源特徴量の変換

- 特徴量パラメータとして**基本周波数 ( $F_0$ )**や**周期／非周期性** (混合励振源使用時) が用いられる.
- **声の高さ**や**声のかすれ**などを表す.
- 基本周波数を高くすれば高い声に, 低くすれば低い声に変換することができる.
- 混合励振源使用時には, 非周期性を大きくすればかすれた声に変換できる.

混合励振源



# リアルタイム音声変換デモ

## リアルタイム声質変換ソフト

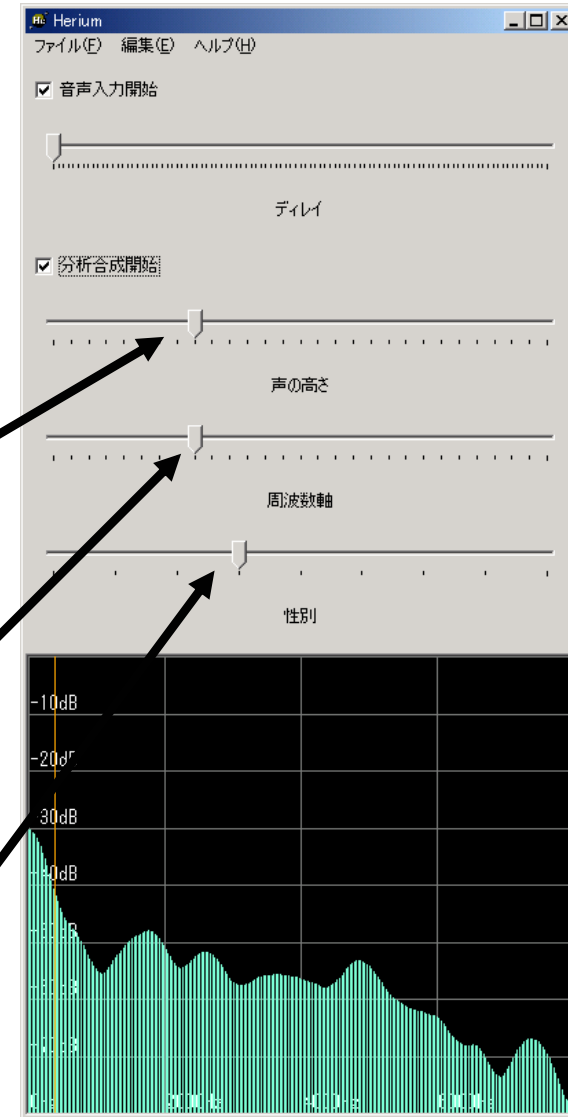
Herium: High Entertaining Real-time  
Input Utterance Modifier

作成者: 名城大学 坂野秀樹先生

声の高さの調節  
➤  $F_0$  の変換

声質の調節  
➤ スペクトルの変換  
(周波数軸の伸縮)

性別の調節  
➤  $F_0$  とスペクトルの変換



<http://www.sp.m.is.nagoya-u.ac.jp/people/banno/spLibs/herium/index-j.html>

# 内容

---

1. 音声変換のしくみ
2. 統計的手法による声質変換
  - 2.1. 基本的な枠組み
  - 2.2. フレームベース変換法
  - 2.3. 系列ベース変換法
3. 応用例

# 統計的手法に基づく声質モデリング

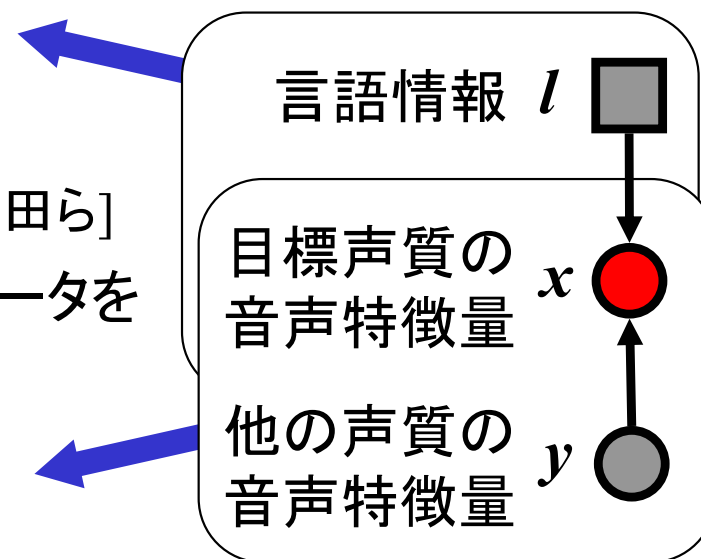
- 確率モデルによる音声特徴量のモデル化

- テキスト音声合成 (TTS)

- 確率密度関数 (*p.d.f.*)  $P(x | l)$  のモデル化
    - 隠れマルコフモデル (HMM) による手法 [徳田ら]
    - 言語情報が付与された目標声質の音声データを用いて学習

- 声質変換

- $P(x | y)$  のモデル化
    - 混合正規分布モデル (GMM) による手法 [Stylianou *et al.*]
    - 言語情報は同一で、所望の声質成分のみが異なる音声データ (パラレルデータ) を用いて学習



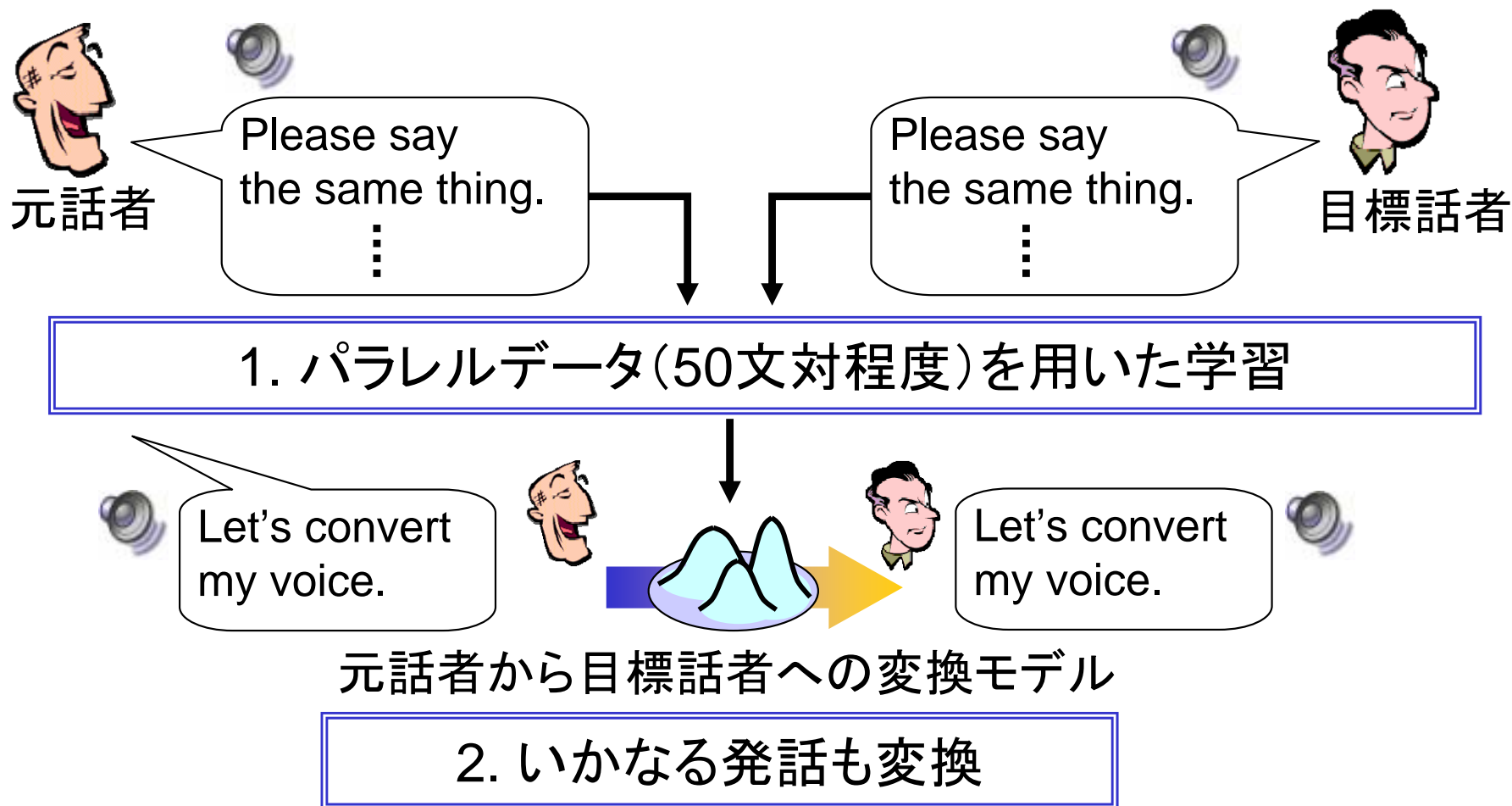
- モデル化される声質成分は**学習データ**により決定



# 統計的手法に基づく声質変換

[Abe et al.]

- 異なる話者や発話様式間における変換を実現できる。
- 入・出力話者による同一内容発声の音声データ(パラレルデータ)を用いて入・出力音声特徴量間の対応関係を統計的にモデル化する。



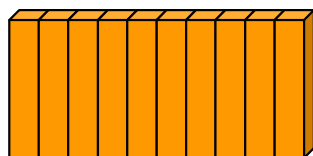
# 学習手順

音声波形 → 音声特徴量 → 対応付けられた音声特徴量 → 変換規則

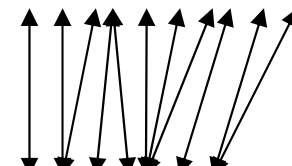
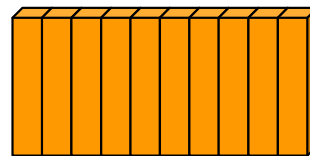
入力話者音声



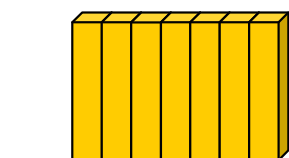
入力音声特徴量



対応づけられた  
入力音声特徴量

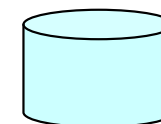


出力話者音声



出力音声特徴量

対応づけられた  
出力音声特徴量

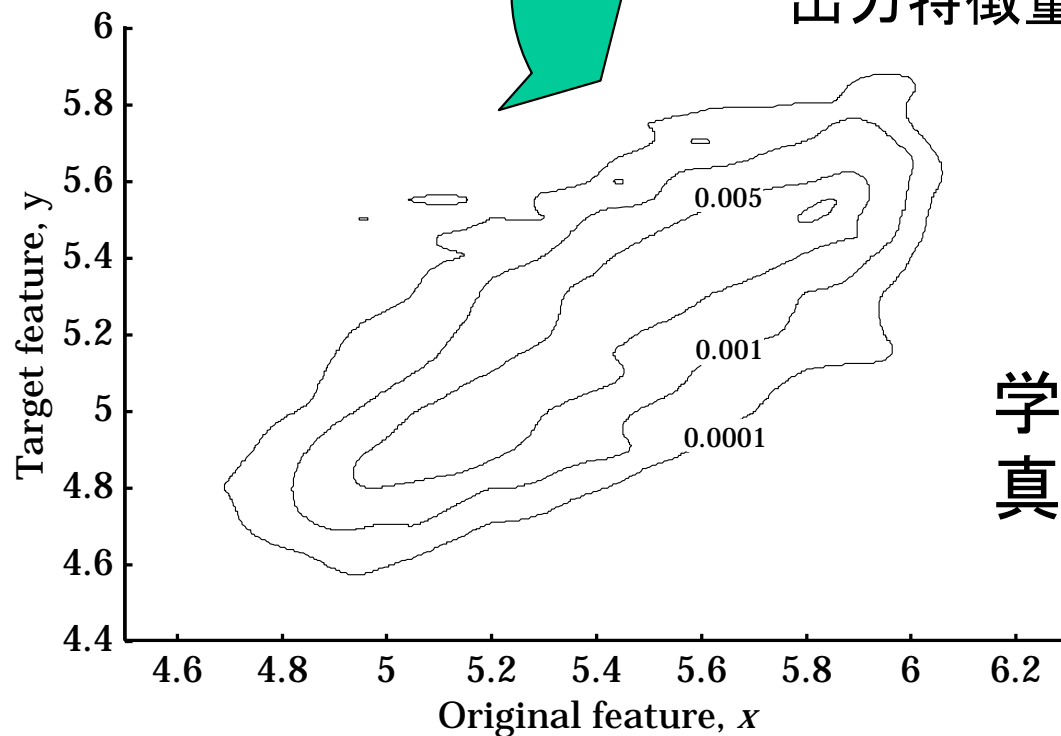


変換規則の  
統計的モデル化

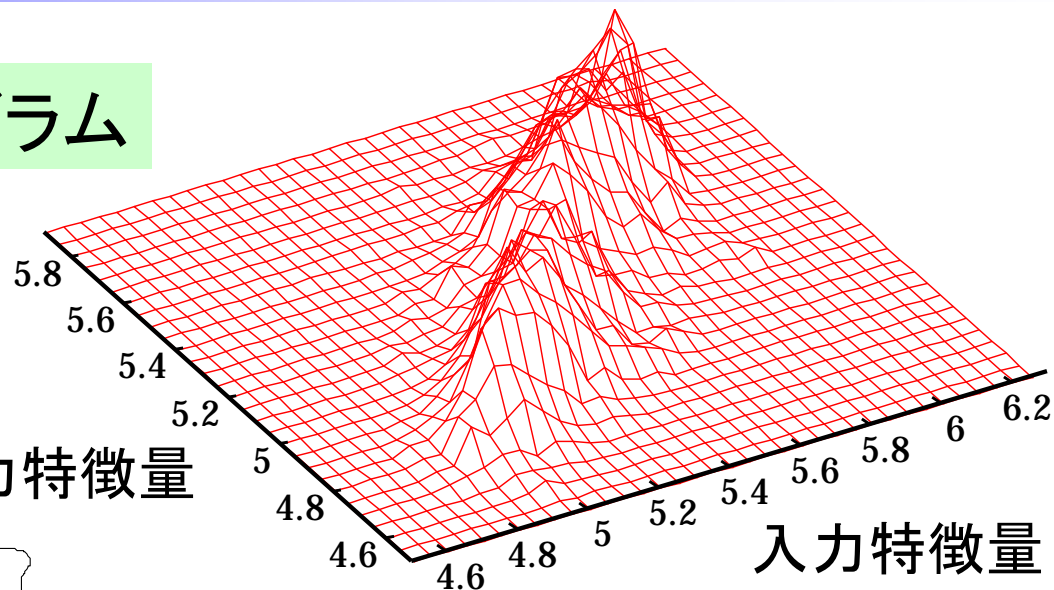
# 学習データの一例

## 学習データのヒストグラム

上から見ると



出力特徴量



学習データ量が限られているため  
真の分布は分からない。

# 内容

---

1. 音声変換のしくみ
2. 統計的手法による声質変換
  - 2.1. 基本的な枠組み
  - 2.2. フレームベース変換法
  - 2.3. 系列ベース変換法
3. 応用例

# フレームベースの変換関数

- 各フレーム(各時間)において独立に変換処理を行う.
  - 入力特徴量ベクトル:  $\mathbf{x}_t$
  - 出力特徴量ベクトル:  $\mathbf{y}_t$
  - 統計モデル:  $\lambda$
  - 変換特徴量ベクトル:  $\hat{\mathbf{y}}_t$

$$\hat{\mathbf{y}}_t = F_\lambda(\mathbf{x}_t)$$

例えば

$$= E_\lambda[\mathbf{y}_t | \mathbf{x}_t] = \int \mathbf{y}_t P(\mathbf{y}_t | \mathbf{x}_t, \lambda) d\mathbf{y}_t$$

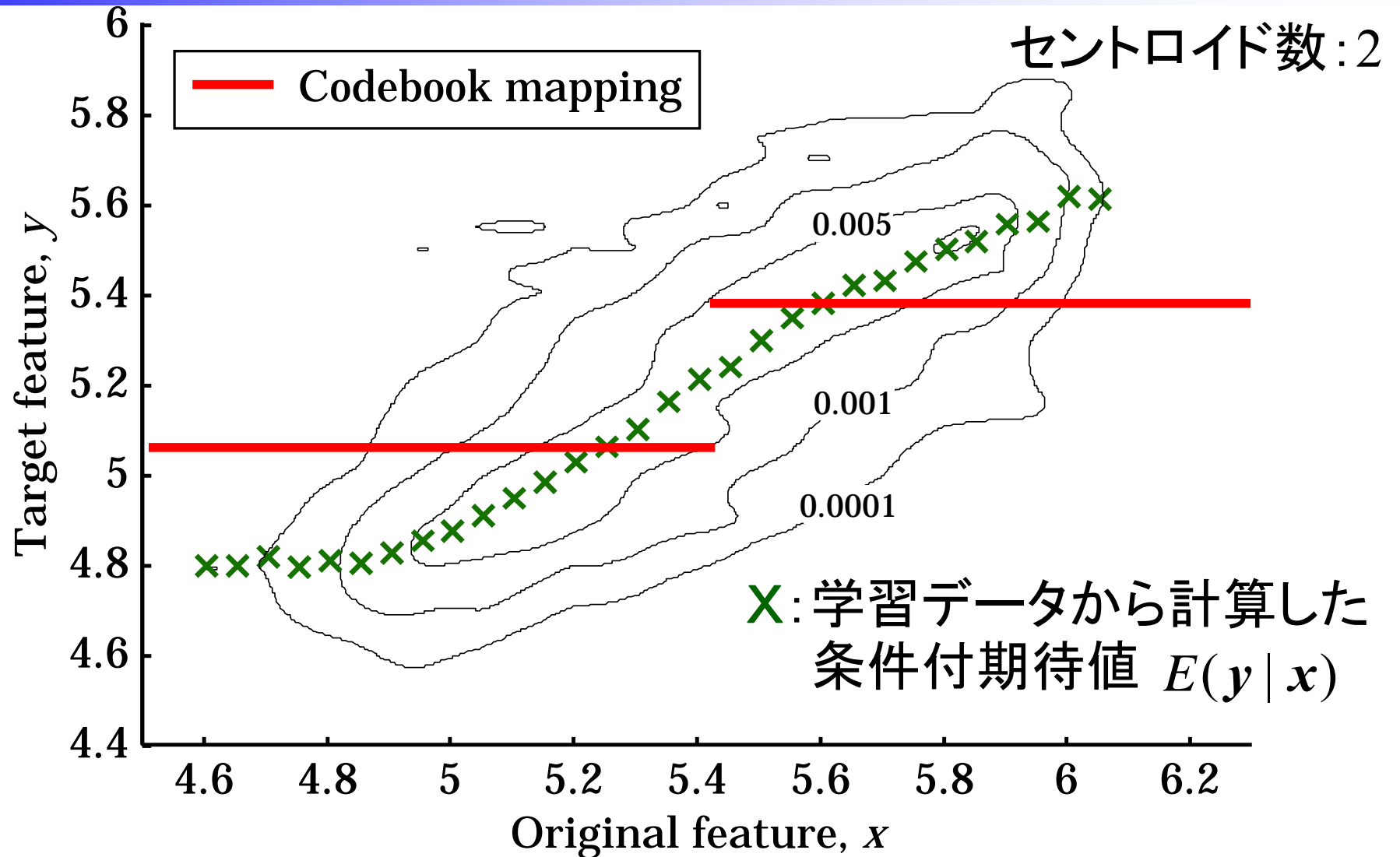
とか

$$= \arg \max_{\mathbf{y}_t} P(\mathbf{y}_t | \mathbf{x}_t, \lambda)$$

# 代表的な変換法

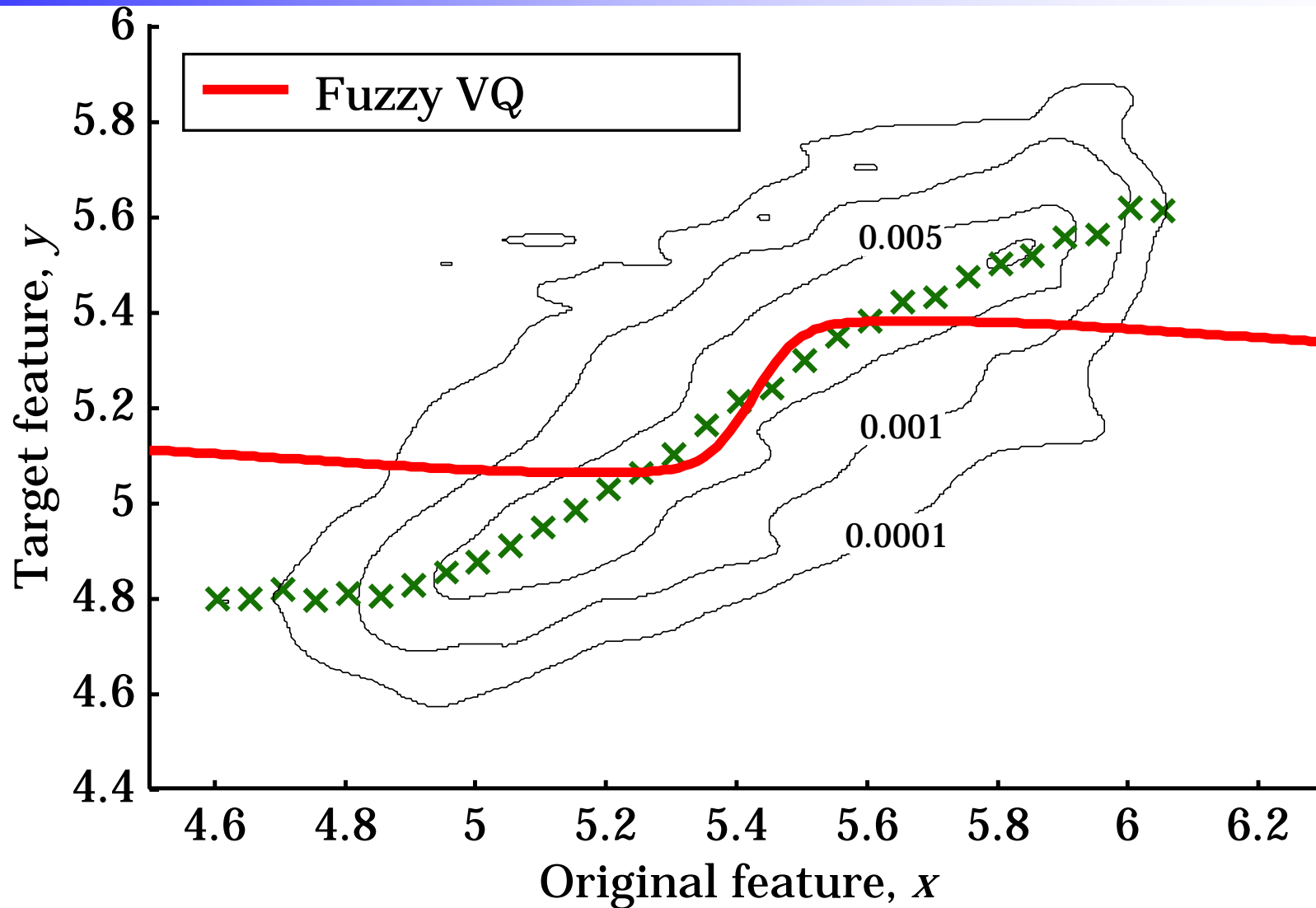
1. コードブックマッピング法 [Abe <i>et al.</i> ]	$\hat{\mathbf{y}}_t = \boldsymbol{\mu}_m^{(y)}$ ハードクラスタリングと離散的なマッピング
2. ファジーベクトル量子化に基づくコードブックマッピング法 [中村 他]	$\hat{\mathbf{y}}_t = \sum_{m=1}^M \gamma_{m,t}^{(x)} \boldsymbol{\mu}_m^{(y)}$ ソフトクラスタリングと離散的なマッピング
3. ファジーベクトル量子化と差分ベクトルに基づくコードブックマッピング法 [Matsumoto <i>et al.</i> ]	$\hat{\mathbf{y}}_t = \mathbf{x}_t + \sum_{m=1}^M \gamma_{m,t}^{(x)} (\boldsymbol{\mu}_m^{(y)} - \boldsymbol{\mu}_m^{(x)})$ ソフトクラスタリングと連続的なマッピング
4. 線形回帰法 [Valbret <i>et al.</i> ]	$\hat{\mathbf{y}}_t = \mathbf{A}_m \mathbf{x}_t + \mathbf{b}_m$ ハードクラスタリングと連続的かつ高精度な変換
5. 混合正規分布モデルに基づく変換法 [Stylianou <i>et al.</i> ]	$\hat{\mathbf{y}}_t = \sum_{m=1}^M \gamma_{m,t}^{(x)} (\mathbf{A}_m \mathbf{x}_t + \mathbf{b}_m)$ ソフトクラスタリングと連続的かつ高精度な変換

# コードブックマッピング法の変換関数



$$\hat{y}_t = \mu_m^{(y)} \quad \text{ハードクラスタリングと離散的なマッピング}$$

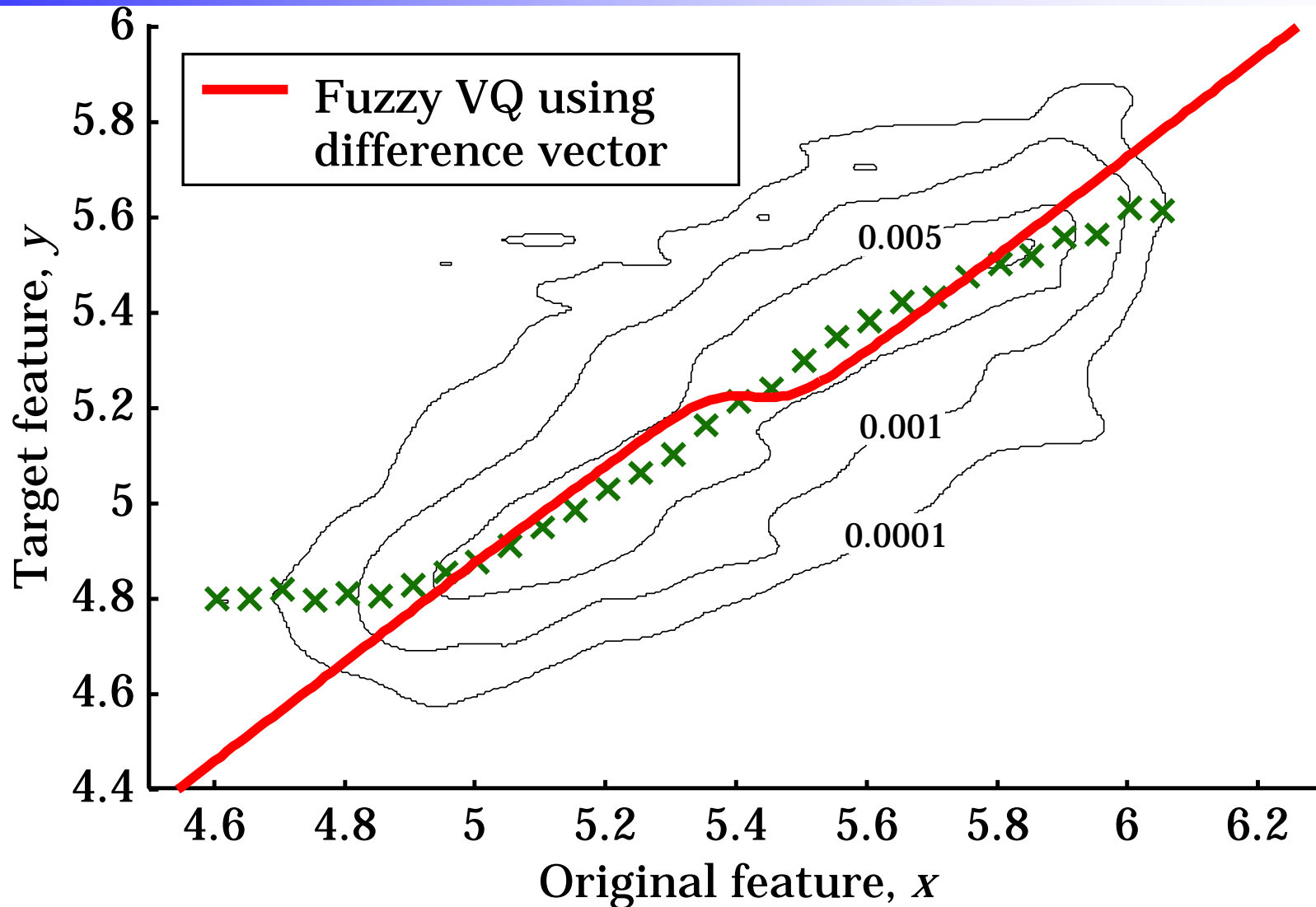
# ファジーVQ使用時の変換関数



$$\hat{y}_t = \sum_{m=1}^M \gamma_{m,t}^{(x)} \mu_m^{(y)} \quad \text{ソフトクラスタリングと離散的なマッピング}$$

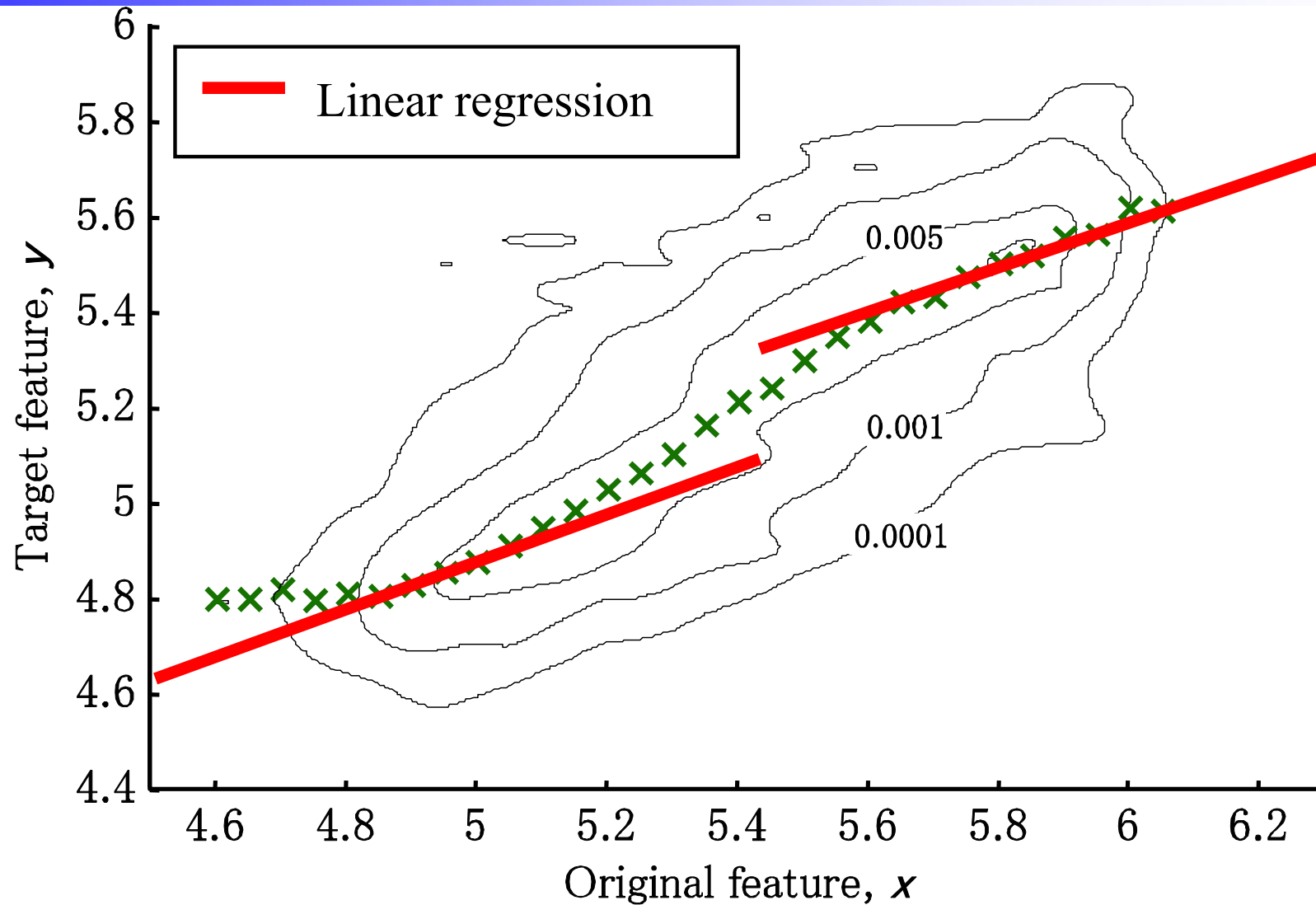


# 差分ベクトル使用時の変換関数



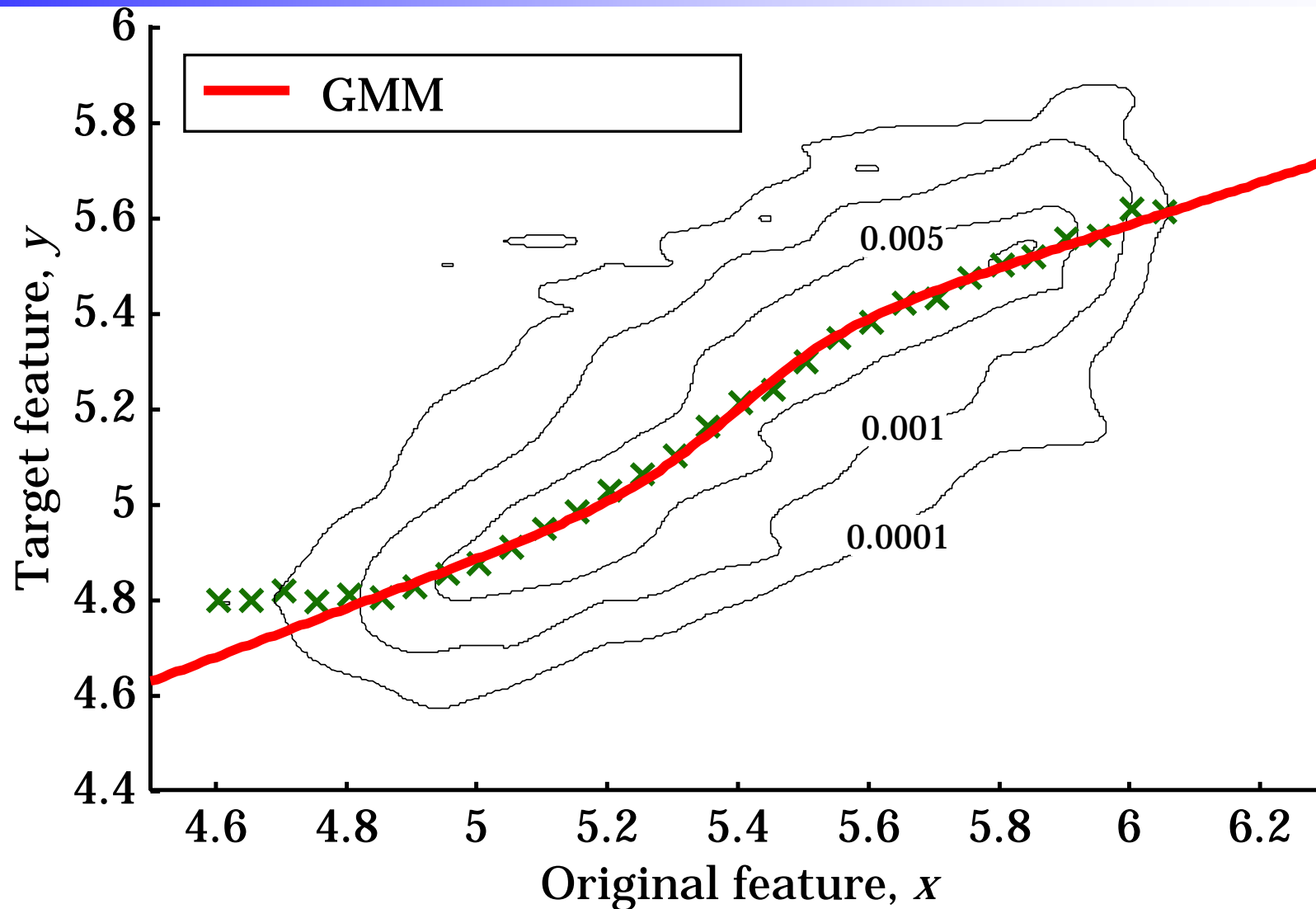
$$\hat{y}_t = x_t + \sum_{m=1}^M \gamma_{m,t}^{(x)} \left( \mu_m^{(y)} - \mu_m^{(x)} \right) \text{ ソフトクラスタリングと連続的なマッピング}$$

# 線形回帰に基づく変換関数



$$\hat{y}_t = A_m \mathbf{x}_t + \mathbf{b}_m \quad \text{ハードクラスタリングと連続的かつ高精度な変換}$$

# 混合正規分布モデルに基づく変換関数



$$\hat{y}_t = \sum_{m=1}^M \gamma_{m,t}^{(x)} (A_m \mathbf{x}_t + \mathbf{b}_m) \quad \text{ソフトクラスタリングと連続的かつ高精度な変換}$$

# 混合正規分布モデル(GMM)

- 正規分布の重み付け和により多様な形状の*p.d.f.*をモデル化する.
- モデルパラメータは各分布における混合重み  $\alpha_m$ , 平均ベクトル  $\mu_m$ , 並びに共分散行列  $\Sigma_m$  である.

$$\text{制約条件: } \sum_{m=1}^M \alpha_m = 1, \quad \alpha_m \geq 0$$

## 1次元の場合

$$p.d.f.: P(x_t | \lambda) = \sum_{m=1}^M \alpha_m N(x_t; \mu_m, \Sigma_m)$$

$$m\text{番目の正規分布: } N(x_t; \mu_m, \sigma_m^2) = \frac{1}{\sqrt{2\pi\sigma_m^2}} \exp\left(-\frac{(x_t - \mu_m)^2}{2\sigma_m^2}\right)$$

## 多次元の場合

$$p.d.f.: P(\mathbf{x}_t | \lambda) = \sum_{m=1}^M \alpha_m N(\mathbf{x}_t; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m)$$

$m$ 番目の正規分布:

$$N(\mathbf{x}_t; \boldsymbol{\mu}_m, \boldsymbol{\Sigma}_m) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}_m|}} \exp\left(-\frac{1}{2} (\mathbf{x}_t - \boldsymbol{\mu}_m)^\top \boldsymbol{\Sigma}_m^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_m)\right)$$

# 結合確率密度のモデル化

[Kain and Macon]

- パラレルデータから結合ベクトルを作成し、結合*p.d.f.*をGMMでモデル化する。

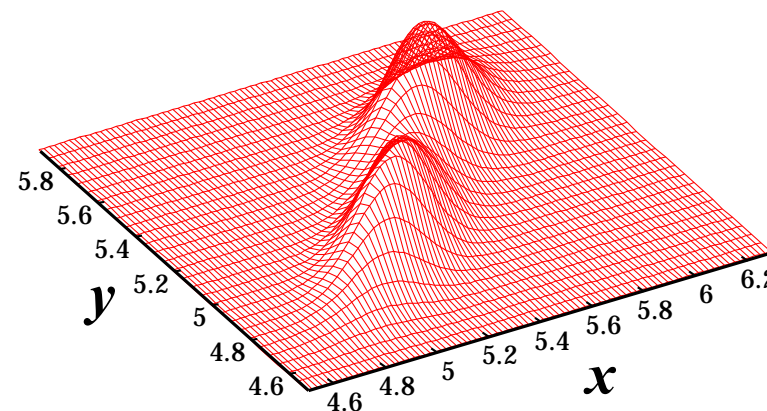
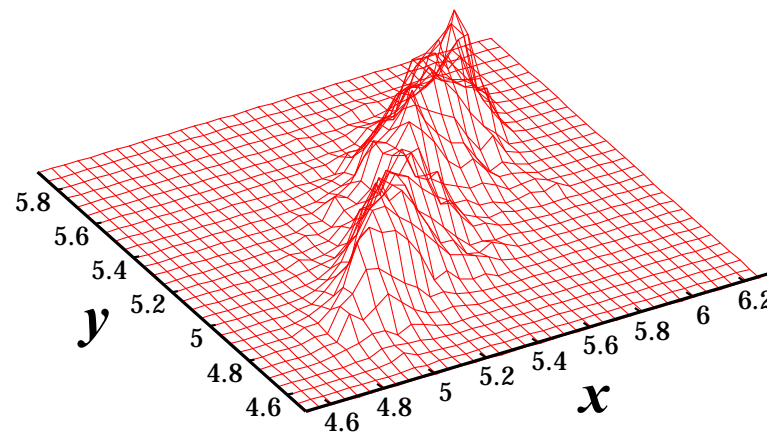
学習データのヒストグラム

結合ベクトル:  $\begin{bmatrix} \mathbf{x}_t \\ \mathbf{y}_t \end{bmatrix}$

GMM

$$p.d.f.: P(\mathbf{x}_t, \mathbf{y}_t | \lambda^{(x,y)}) = \sum_{m=1}^M P(m | \lambda^{(x,y)}) P(\mathbf{x}_t, \mathbf{y}_t | m, \lambda^{(x,y)})$$

$$= \sum_{m=1}^M \alpha_m N\left(\begin{bmatrix} \mathbf{x}_t \\ \mathbf{y}_t \end{bmatrix} \mid \boldsymbol{\mu}_m^{(x,y)}, \boldsymbol{\Sigma}_m^{(x,y)}\right)$$

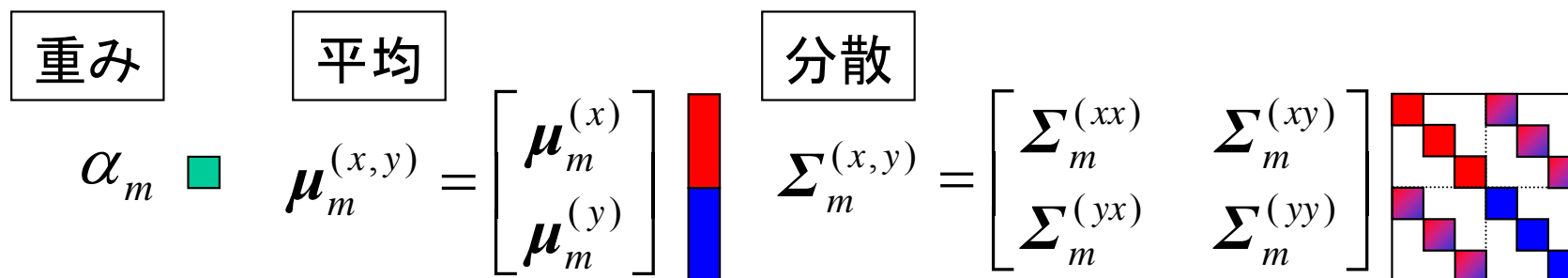


# GMMの構造

- $m$ 番目の混合分布パラメータ

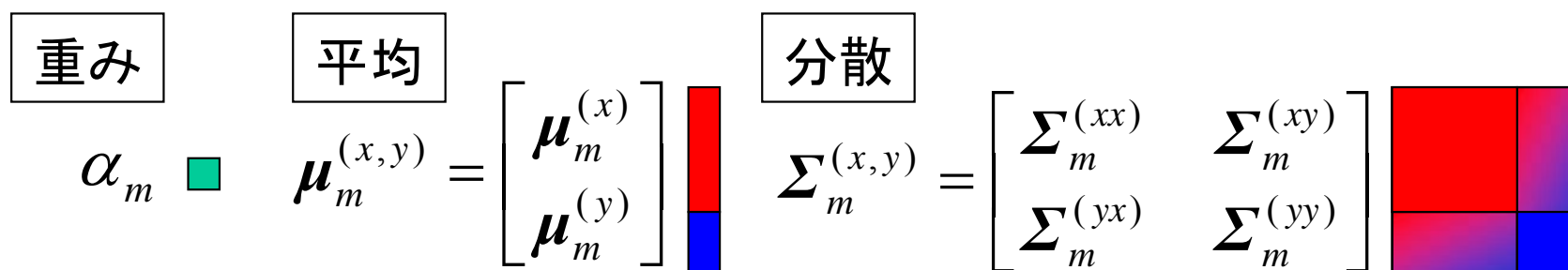
## 次元間の相関をモデル化しない場合

同一の特徴量パラメータ間で変換を行う際に限って用いることができる。



## 次元間の相関をモデル化する場合

異なる特徴量パラメータ間で変換を行う際には必要不可欠である。



# 条件付p.d.f.

・結合p.d.f.:  $P(\mathbf{x}_t, \mathbf{y}_t | \lambda^{(x,y)}) = \sum_{m=1}^M P(m | \lambda^{(x,y)}) P(\mathbf{x}_t | m, \lambda^{(x,y)}) P(\mathbf{y}_t | \mathbf{x}_t, m, \lambda^{(x,y)})$

周辺化すると  $P(\mathbf{x}_t | \lambda^{(x,y)}) = \sum_{m=1}^M P(m | \lambda^{(x,y)}) P(\mathbf{x}_t | m, \lambda^{(x,y)})$

・入力  $\mathbf{x}_t$  が与えられた際の条件付p.d.f.:

$$P(\mathbf{y}_t | \mathbf{x}_t, \lambda^{(x,y)}) = \frac{P(\mathbf{x}_t, \mathbf{y}_t | \lambda^{(x,y)})}{P(\mathbf{x}_t | \lambda^{(x,y)})}$$

$$= \sum_{m=1}^M \underbrace{P(m | \mathbf{x}_t, \lambda^{(x,y)})}_{\text{事後確率}} \underbrace{P(\mathbf{y}_t | \mathbf{x}_t, m, \lambda^{(x,y)})}_{\text{正規分布}}$$

GMMでモデル化される。



事後確率

$$\frac{P(m | \lambda^{(x,y)}) P(\mathbf{x}_t | m, \lambda^{(x,y)})}{\sum_{n=1}^M P(n | \lambda^{(x,y)}) P(\mathbf{x}_t | n, \lambda^{(x,y)})}$$

正規分布

$$\begin{aligned} \text{平均: } \mu_{m,t}^{(y|x)} &= \mu_m^{(y)} + \Sigma_m^{(yx)} \Sigma_m^{(xx)^{-1}} (\mathbf{x}_t - \mu_m^{(x)}) \\ &= \mathbf{A}_m \mathbf{x}_t + \mathbf{b}_m \end{aligned}$$

$$\text{共分散: } \Sigma_m^{(y|x)} = \Sigma_m^{(yy)} - \Sigma_m^{(yx)} \Sigma_m^{(xx)^{-1}} \Sigma_m^{(xy)}$$

# 最小平均自乗誤差推定

[Stylianou *et al.*]

- 入力が与えられた際の出力の条件付期待値へと変換する.

$$\begin{aligned}\hat{y}_t &= \int \underline{P(\mathbf{y}_t | \mathbf{x}_t, \lambda^{(x,y)})} \mathbf{y}_t d\mathbf{y}_t \\ &= \int \underline{\sum_{m=1}^M P(m | \mathbf{x}_t, \lambda^{(x,y)}) P(\mathbf{y}_t | \mathbf{x}_t, m, \lambda^{(x,y)})} \mathbf{y}_t d\mathbf{y}_t \\ &= \sum_{m=1}^M \underline{P(m | \mathbf{x}_t, \lambda^{(x,y)})} \int \underline{P(\mathbf{y}_t | \mathbf{x}_t, m, \lambda^{(x,y)})} \mathbf{y}_t d\mathbf{y}_t \\ &= \sum_{m=1}^M \underline{\gamma_{m,t}^{(x)}} \underline{(A_m \mathbf{x}_t + \mathbf{b}_m)}\end{aligned}$$

ソフトクラスタリング

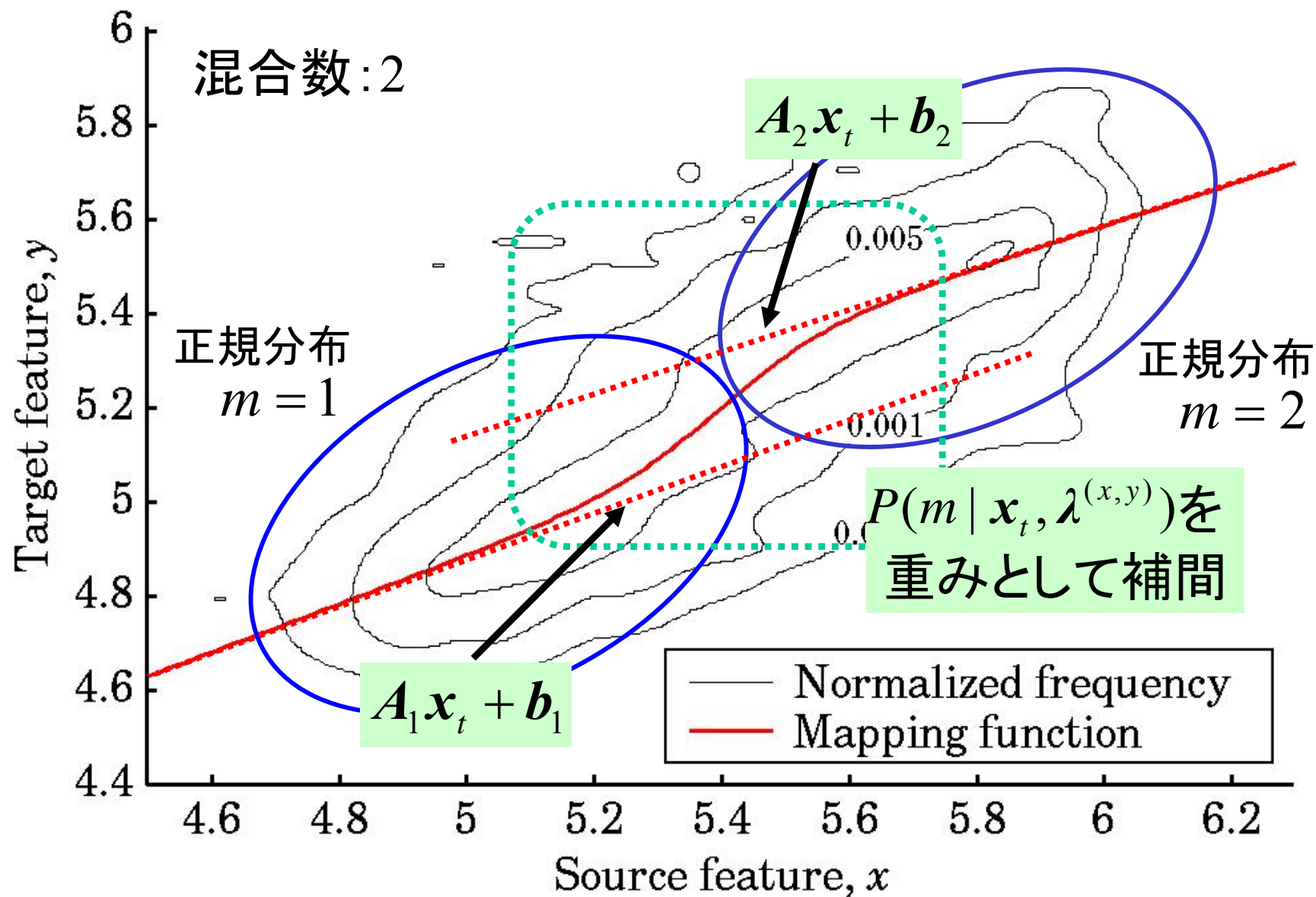
$$\begin{aligned}\gamma_{m,t}^{(x)} &= P(m | \mathbf{x}_t, \lambda^{(x,y)}) \\ &= \frac{\alpha_m N(\mathbf{x}_t; \boldsymbol{\mu}_m^{(x)}, \boldsymbol{\Sigma}_m^{(xx)})}{\sum_{n=1}^M \alpha_n N(\mathbf{x}_t; \boldsymbol{\mu}_n^{(x)}, \boldsymbol{\Sigma}_n^{(xx)})}\end{aligned}$$

入・出力の相関を考慮した連続的な変換

$$\begin{aligned}A_m &= \boldsymbol{\Sigma}_m^{(yx)} (\boldsymbol{\Sigma}_m^{(xx)})^{-1} \\ \mathbf{b}_m &= \boldsymbol{\mu}_m^{(y)} - A_m \boldsymbol{\mu}_m^{(x)}\end{aligned}$$



# GMMに基づく変換関数の解釈



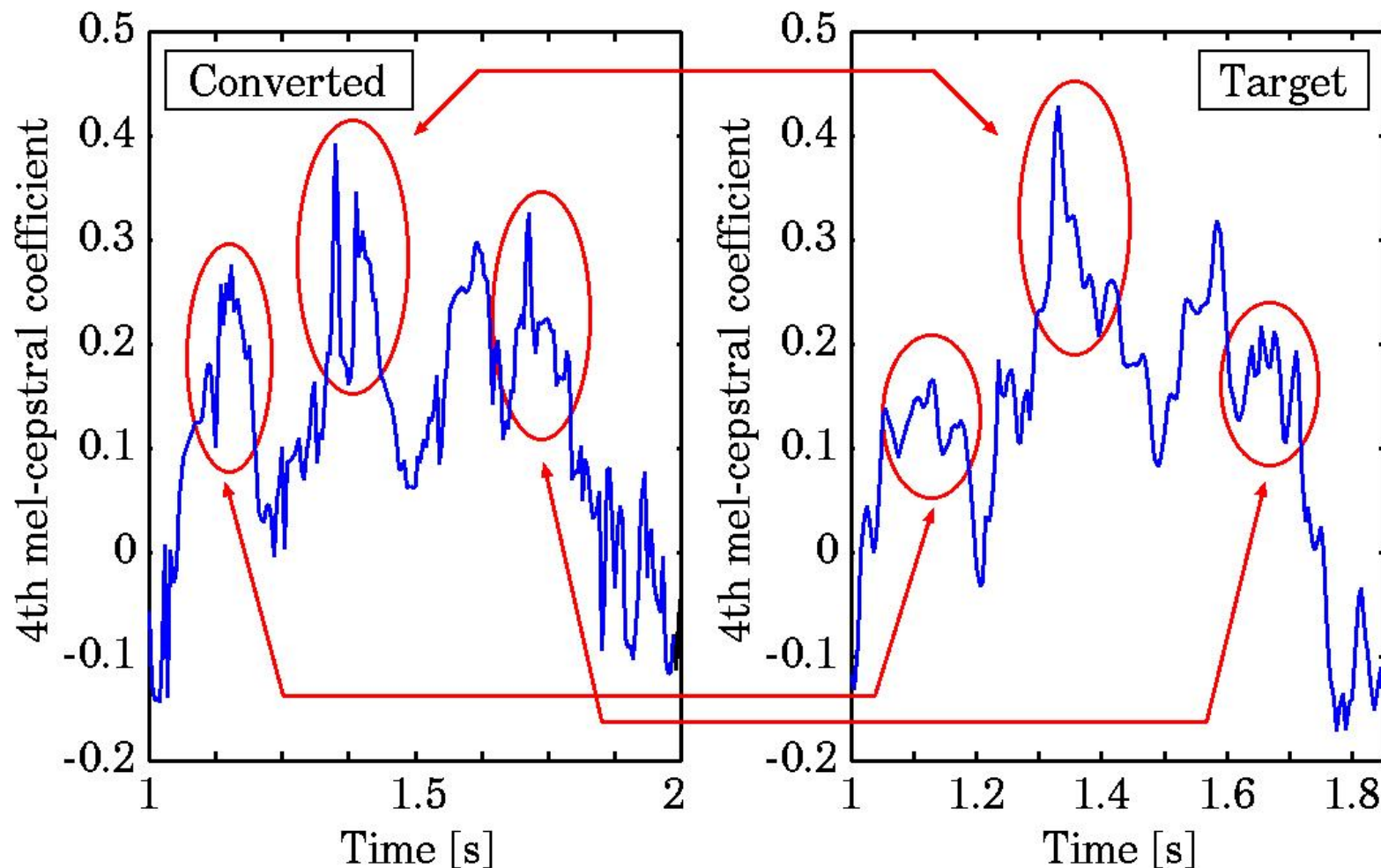
# 話者変換音声サンプル

- 話者: 4名 (男: bdl, 男: rms, 女: clb, 女: slt)
- 学習: 50文対 (完全自動学習)
- 変換法: 最小平均自乗誤差推定法 (フレームベース変換法)

		目標話者			
		bdl	rms	clb	slt
元話者	bdl				
	rms				
	clb				
	slt				

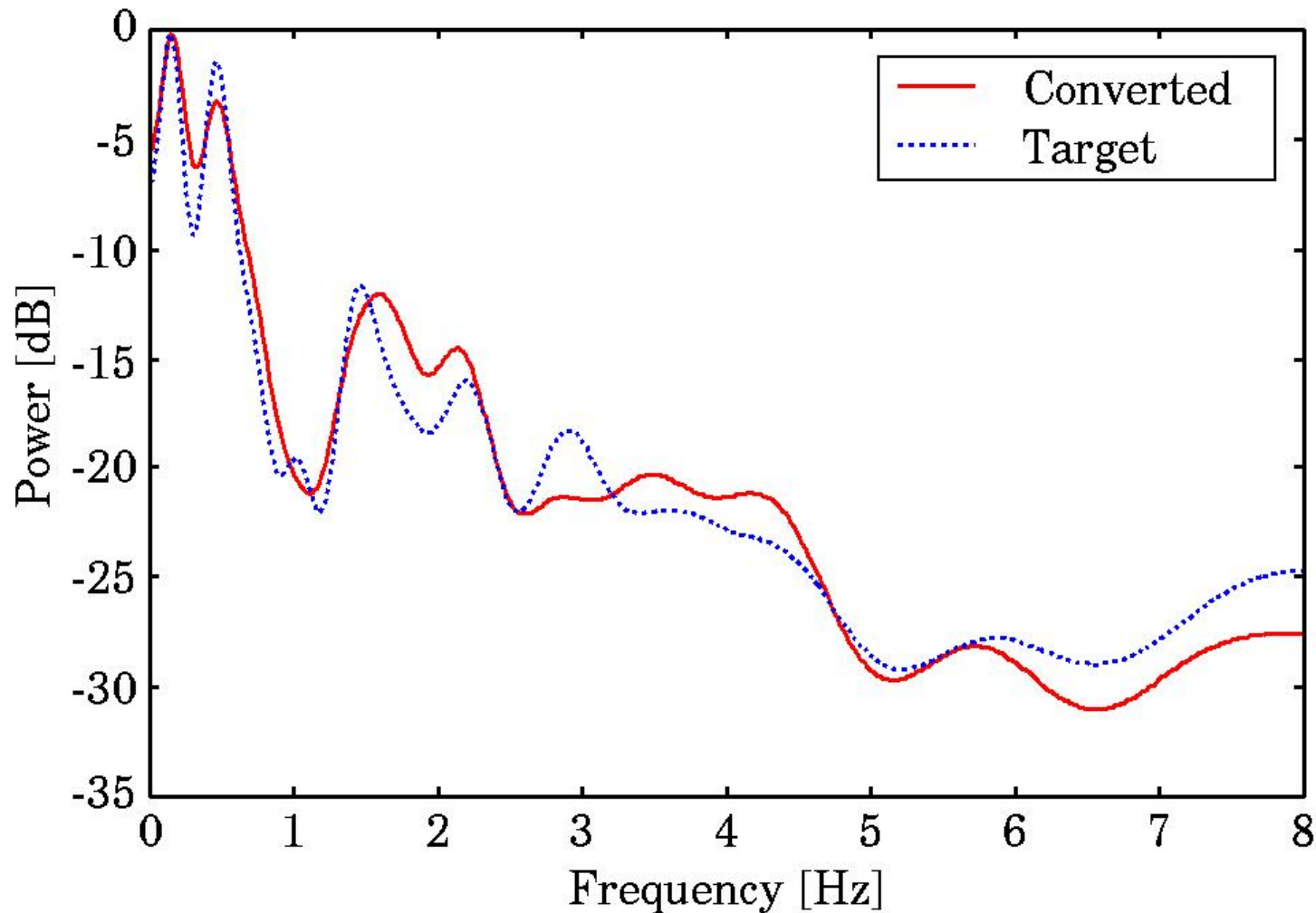
# 問題点1: 時間的依存関係の無視

- フレーム毎に独立して変換処理が行われるため、不適切なパラメータ遷移が発生する場合があります。



# 問題点2: 過剰な平滑化

- 汎化処理により詳細なスペクトル構造が消失し、変換音声の肉声感が大幅に損なわれる。



# 内容

---

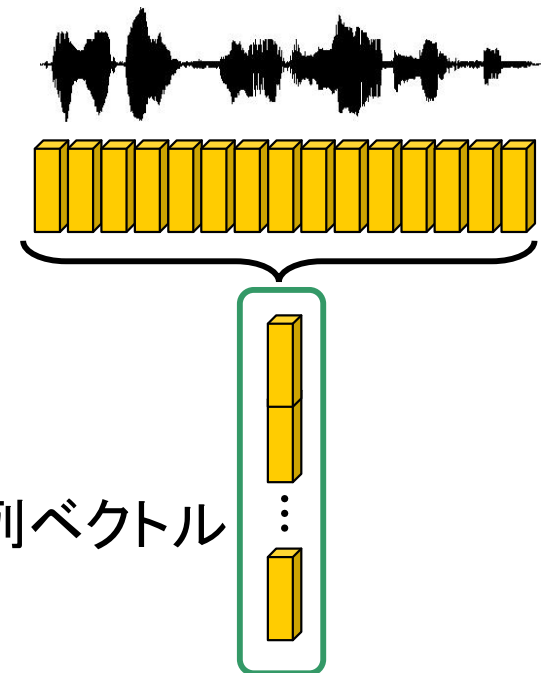
1. 音声変換のしくみ
2. 統計的手法による声質変換
  - 2.1. 基本的な枠組み
  - 2.2. フレームベース変換法
  - 2.3. 系列ベース変換法
3. 応用例

# 系列ベースの変換関数

- 時系列単位で変換処理を行う。

- 入力特徴量系列ベクトル:  $\mathbf{x} = [\mathbf{x}_1^T \mid \mathbf{x}_2^T \mid \dots \mid \mathbf{x}_T^T]^T$
- 出力特徴量系列ベクトル:  $\mathbf{y} = [\mathbf{y}_1^T \mid \mathbf{y}_2^T \mid \dots \mid \mathbf{y}_T^T]^T$
- 統計モデル:  $\lambda$
- 変換特徴量系列ベクトル:  $\hat{\mathbf{y}} = [\hat{\mathbf{y}}_1^T \mid \hat{\mathbf{y}}_2^T \mid \dots \mid \hat{\mathbf{y}}_T^T]^T$

$$\begin{aligned}\hat{\mathbf{y}} &= F_\lambda(\mathbf{x}) \\ &= \arg \max P(\mathbf{y} \mid \mathbf{x}, \lambda)\end{aligned}$$



系列ベクトル

# 最尤系列変換法

[Toda *et al.*]

- 系列に基づく特徴量を考慮した変換処理を行う.

## 1. 動的特徴量を考慮した系列ベースの最尤推定

- フレーム間相関を考慮した変換を行う.
- 問題点1 (時間的依存関係の無視)を解決できる.

## 2. 系列内変動の明示的なモデル化の導入

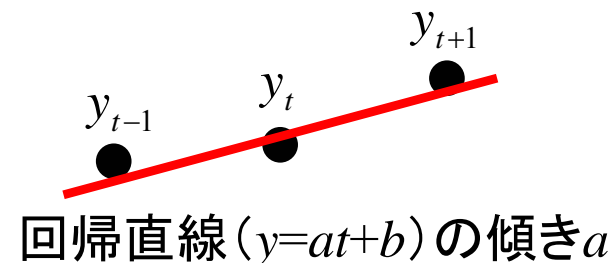
- 2次モーメントを考慮した変換を行う.
- 問題点2 (過剰な平滑化)の影響を大幅に抑えることができる.

# 1. 動的特徴量の導入

- 学習時には、静的特徴量のみでなく動的特徴量も含んだ結合*p.d.f.*をGMMでモデル化する。

動的特徴量(微分係数)の計算

例えば  $\Delta y_t = y_t - y_{t-1}$  や  $\Delta y_t = \frac{1}{2}(y_{t+1} - y_{t-1})$



静的・動的特徴量ベクトルに対する結合*p.d.f.*をモデル化する。

- 入力特徴量ベクトル:  $X_t = [x_t^T \ \Delta x_t^T]^T$
- 出力特徴量ベクトル:  $Y_t = [y_t^T \ \Delta y_t^T]^T$
- GMMによる結合*p.d.f.*のモデル化:

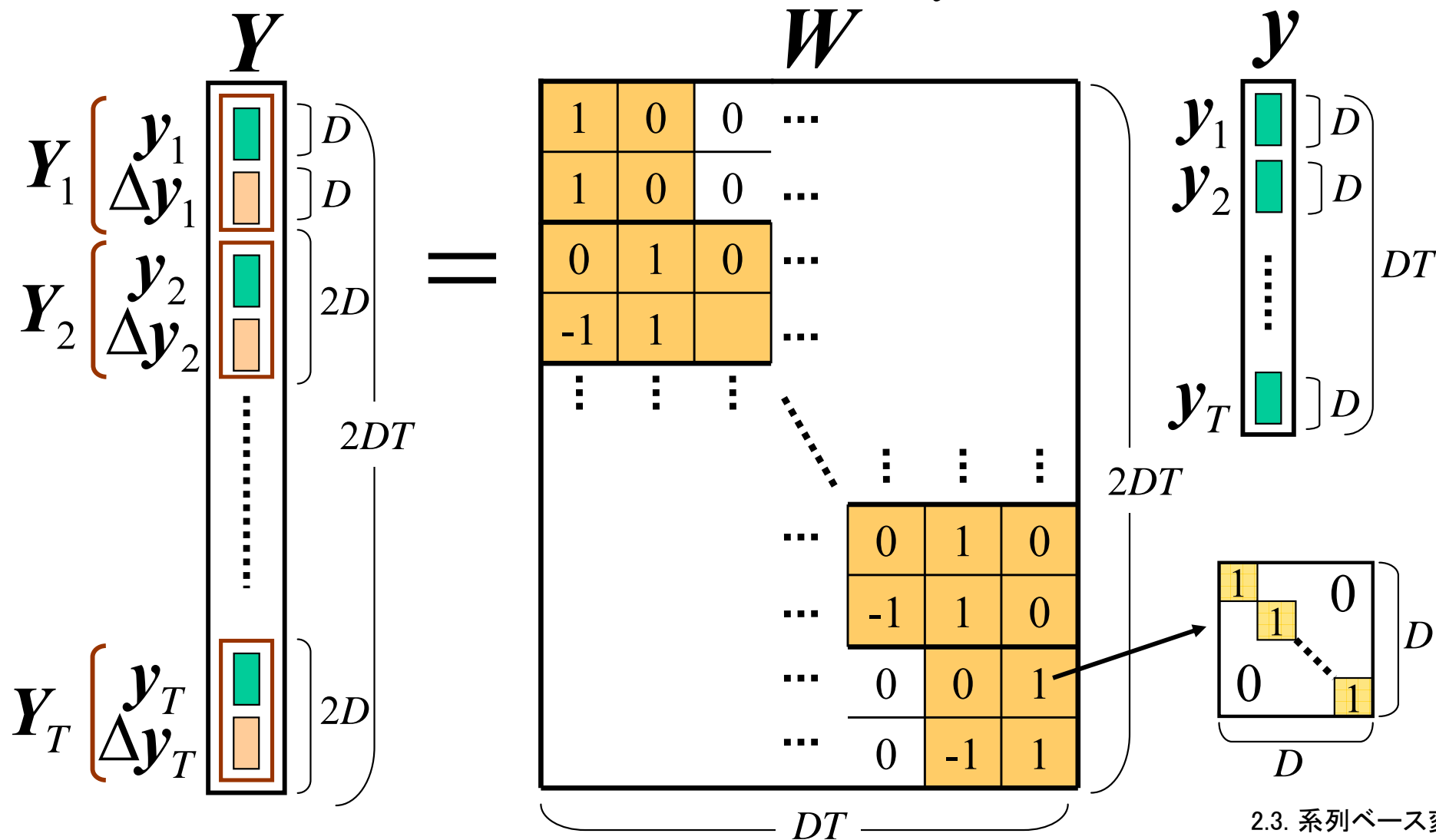
$$P(X_t, Y_t | \lambda^{(X,Y)}) = \sum_{m=1}^M \alpha_m N \left( \begin{bmatrix} X_t \\ Y_t \end{bmatrix}; \begin{bmatrix} \mu_m^{(X)} \\ \mu_m^{(Y)} \end{bmatrix}, \begin{bmatrix} \Sigma_m^{(XX)} & \Sigma_m^{(XY)} \\ \Sigma_m^{(YX)} & \Sigma_m^{(YY)} \end{bmatrix} \right)$$

- 変換時には、**静的・動的特徴量間の明示的な関係**を考慮して、条件付*p.d.f.*の尤度最大化に基づき入力特徴量ベクトルを変換する。



# 静的・動的特徴量の明示的な関係

- 静的・動的特徴量系列  $Y = [Y_1^T \ Y_2^T \ \dots \ Y_T^T]^T$  は静的特徴量系列  $y = [y_1^T \ y_2^T \ \dots \ y_T^T]^T$  の線形変換  $Y = Wy$  により表される。



# 動的特徴量を考慮した最尤系列変換

[Toda *et al.*]

- 静的・動的特徴量系列に対する条件付 *p.d.f.* の尤度を最大化する静的特徴量系列を推定する.

制約条件:  $Y = Wy$

$$\begin{aligned}\hat{y} &= \arg \max P(Y | X, \lambda^{(X,Y)}) \\ &= \arg \max P(Wy | X, \lambda^{(X,Y)})\end{aligned}$$

最尤な  $Y$  となる  
様に  $\mathcal{Y}$  を推定

静的特徴量系列  
に関する最大化

静的・動的特徴量  
系列に関する尤度

※HMM音声合成におけるパラメータ生成[Tokuda *et al.*]と同じ処理



# 系列単位の尤度関数

尤度関数

$$P(Y | X, \lambda^{(X,Y)}) = \prod_{t=1}^T \sum_{m=1}^M P(m | X_t, \lambda^{(X,Y)}) P(Y_t | X_t, m, \lambda^{(X,Y)})$$

$$= \sum_{\text{all } m} \underbrace{P(m | X, \lambda^{(X,Y)})}_{\text{事後確率}} \underbrace{P(Y | X, m, \lambda^{(X,Y)})}_{\text{正規分布}}$$

事後確率

$$P(m | X, \lambda^{(X,Y)}) = \prod_{t=1}^T P(m_t | X_t, \lambda^{(X,Y)})$$

分布系列:  $m = \{m_1, m_2, \dots, m_T\}$

正規分布

$$P(Y | X, m, \lambda^{(X,Y)}) = \prod_{t=1}^T P(Y_t | X_t, m_t, \lambda^{(X,Y)})$$

$$= P(Y | X, m, \lambda^{(X,Y)})$$

$$= N(Y; \mu_m^{(Y|X)}, \Sigma_m^{(Y|X)})$$

観測

ベクトル

$$Y = \begin{bmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_T \end{bmatrix}$$

平均

ベクトル

$$\mu_m^{(Y|X)} = \begin{bmatrix} \mu_{m_1}^{(Y|X)} \\ \mu_{m_2}^{(Y|X)} \\ \vdots \\ \mu_{m_T}^{(Y|X)} \end{bmatrix}$$

共分散  
行列

$$\Sigma_m^{(Y|X)} = \begin{bmatrix} \Sigma_{m_1}^{(Y|X)} & 0 & \dots & 0 \\ 0 & \Sigma_{m_2}^{(Y|X)} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \Sigma_{m_T}^{(Y|X)} \end{bmatrix}$$

# 最尤特徴量系列の導出

- 単一分布系列  $\mathbf{m} = \{m_1, m_2, \dots, m_T\}$  による近似を用いる場合:  
尤度関数

$$P(Y | X, \lambda^{(X,Y)}) = \sum_{\text{all } \mathbf{m}} P(\mathbf{m} | X, \lambda^{(X,Y)}) P(Y | X, \mathbf{m}, \lambda^{(X,Y)}) \\ \approx P(\mathbf{m} | X, \lambda^{(X,Y)}) P(Y | X, \mathbf{m}, \lambda^{(X,Y)})$$

準最適な分布系列の決定 (ハードクラスタリング)

$$\hat{\mathbf{m}} = \arg \max P(\mathbf{m} | X, \lambda^{(X,Y)})$$

出力静的特徴量系列の推定

$$\hat{\mathbf{y}} = \arg \max P(\hat{\mathbf{m}} | X, \lambda^{(X,Y)}) \underline{\underline{P(Y | X, \hat{\mathbf{m}}, \lambda^{(X,Y)})}} \\ = \arg \max \underline{\underline{-\frac{1}{2} (\mathbf{W}\mathbf{y} - \boldsymbol{\mu}_{\hat{\mathbf{m}}}^{(Y|X)})^T \boldsymbol{\Sigma}_{\hat{\mathbf{m}}}^{(Y|X)-1} (\mathbf{W}\mathbf{y} - \boldsymbol{\mu}_{\hat{\mathbf{m}}}^{(Y|X)})}} \\ = \left( \mathbf{W}^T \boldsymbol{\Sigma}_{\hat{\mathbf{m}}}^{(Y|X)-1} \mathbf{W} \right)^{-1} \mathbf{W}^T \boldsymbol{\Sigma}_{\hat{\mathbf{m}}}^{(Y|X)-1} \boldsymbol{\mu}_{\hat{\mathbf{m}}}^{(Y|X)}$$

※EMアルゴリズムを用いてソフトクラスタリングを行うことも可能である。

# トラジェクトリモデルとしての解釈

[Zen et al.]

静的特徴量系列ベクトル  $\mathbf{y}$  を確率変数としたモデルの尤度最大化

$$P(Y | X, \mathbf{m}, \lambda) = P(W\mathbf{y} | X, \mathbf{m}, \lambda)$$

正規分布

$$\left[ \begin{array}{l} \text{平均: } \boldsymbol{\mu}_m^{(Y|X)} \\ \text{共分散: } \boldsymbol{\Sigma}_m^{(Y|X)} \end{array} \right]$$

$$= \frac{1}{\sqrt{(2\pi)^{2DT} |\boldsymbol{\Sigma}_m^{(Y|X)}|}} \exp\left(-\frac{1}{2} (W\mathbf{y} - \boldsymbol{\mu}_m^{(Y|X)})^T \boldsymbol{\Sigma}_m^{(Y|X)-1} (W\mathbf{y} - \boldsymbol{\mu}_m^{(Y|X)})\right)$$

$$= Z_m \frac{1}{\sqrt{(2\pi)^{2DT} |\mathbf{P}_m|}} \exp\left(-\frac{1}{2} (\mathbf{y} - \bar{\mathbf{y}}_m)^T \mathbf{P}_m^{-1} (\mathbf{y} - \bar{\mathbf{y}}_m)\right)$$

$$= Z_m P(\mathbf{y} | X, \mathbf{m}, \lambda)$$

正規化項

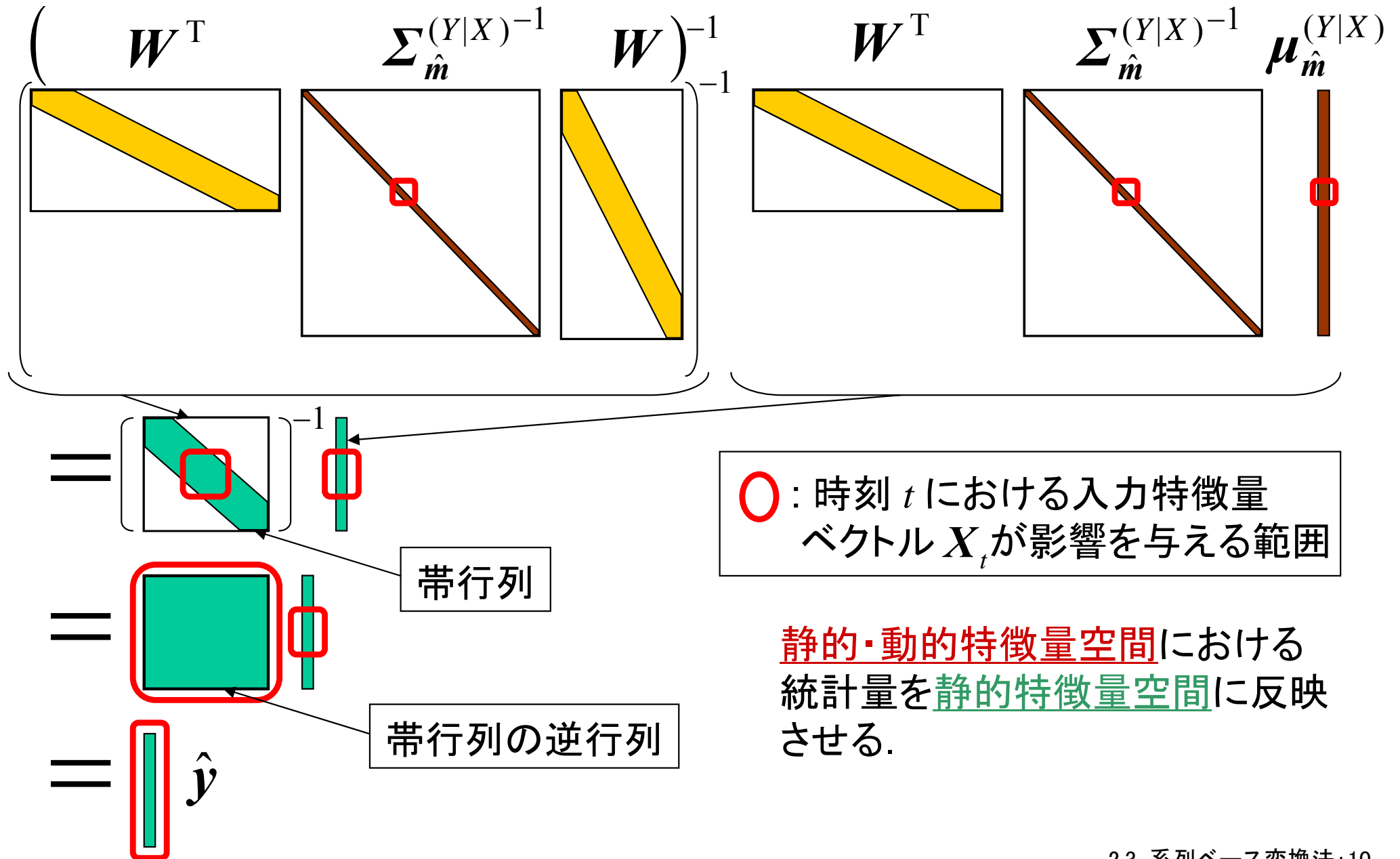
$$\frac{\sqrt{(2\pi)^{DT} |\mathbf{P}_m|}}{\sqrt{(2\pi)^{2DT} |\boldsymbol{\Sigma}_m^{(Y|X)}|}}$$

$$\exp\left(-\frac{1}{2} \left( \boldsymbol{\mu}_m^{(Y|X)T} \boldsymbol{\Sigma}_m^{(Y|X)-1} \boldsymbol{\mu}_m^{(Y|X)} - \bar{\mathbf{y}}_m^T \mathbf{P}_m^{-1} \bar{\mathbf{y}}_m \right)\right)$$

正規分布

$$\left[ \begin{array}{l} \text{平均: } \bar{\mathbf{y}}_m = \left( \mathbf{W}^T \boldsymbol{\Sigma}_m^{(Y|X)-1} \mathbf{W} \right)^{-1} \mathbf{W}^T \boldsymbol{\Sigma}_m^{(Y|X)-1} \boldsymbol{\mu}_m^{(Y|X)} \\ \text{共分散: } \mathbf{P}_m = \left( \mathbf{W}^T \boldsymbol{\Sigma}_m^{(Y|X)-1} \mathbf{W} \right)^{-1} \end{array} \right]$$

# 最尤特徴量系列の演算



# コレスキー分解を用いた解法

[Tokuda et al.]

生成パラメータ

$TD \times TD$ サイズの逆行列演算

$$\hat{y} = R_{\hat{m}}^{-1} r_{\hat{m}} \quad \text{ここで } R_{\hat{m}} = W^T \Sigma_{\hat{m}}^{(Y|X)^{-1}} W \quad \text{および } r_{\hat{m}} = W^T \Sigma_{\hat{m}}^{(Y|X)^{-1}} \mu_{\hat{m}}^{(Y|X)}$$

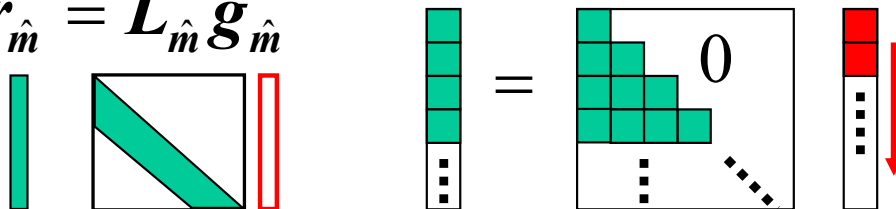
コレスキー分解

$$R_{\hat{m}} = L_{\hat{m}} L_{\hat{m}}^T$$

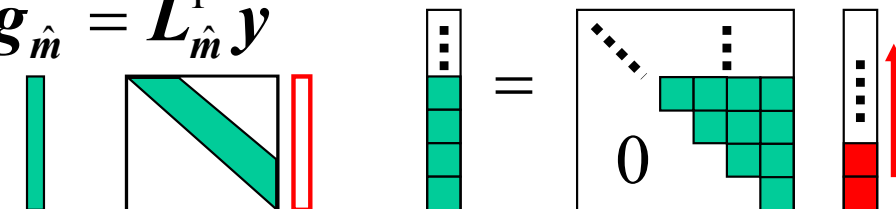
下三角行列とその転置の積に分解

$$r_{\hat{m}} = R_{\hat{m}} \hat{y} = L_{\hat{m}} L_{\hat{m}}^T \hat{y} \quad \text{ここで } L_{\hat{m}}^T \hat{y} = g_{\hat{m}} \quad \text{とすると } r_{\hat{m}} = L_{\hat{m}} g_{\hat{m}}$$

前進代入  $r_{\hat{m}} = L_{\hat{m}} g_{\hat{m}}$



後退代入  $g_{\hat{m}} = L_{\hat{m}}^T \hat{y}$



逆行列演算を必要とせず  
簡単な代入処理による計算

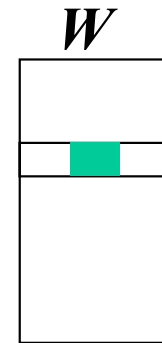


# 再帰式を用いた解法

[Tokuda et al.]

$$\hat{y} = \left( W^T \Sigma_{\hat{m}}^{(Y|X)^{-1}} W \right)^{-1} \left( W^T \Sigma_{\hat{m}}^{(Y|X)^{-1}} \mu_{\hat{m}}^{(Y|X)} \right)$$

$W$  を静的特徴量を求める行と  
動的特徴量を求める行に分解すると



$\Delta y_t$  を求める行

$$w_t = [0, \dots, 0, -I, I, 0, \dots, 0]$$

$$= \left( \underbrace{\Sigma_{\hat{m}}^{(y|X)^{-1}}}_{\text{静的特徴量に関する項}} + \sum_{t=1}^T \underbrace{w_t^T \Sigma_{\hat{m}_t}^{(\Delta y|X)^{-1}} w_t}_{\text{動的特徴量に関する項}} \right)^{-1} \left( \underbrace{\Sigma_{\hat{m}}^{(y|X)^{-1}} \mu_{\hat{m}}^{(y|X)}}_{\text{静的特徴量に関する項}} + \sum_{t=1}^T \underbrace{w_t^T \Sigma_{\hat{m}_t}^{(\Delta y|X)^{-1}} \mu_{\hat{m}_t}^{(\Delta y|X)}}_{\text{動的特徴量に関する項}} \right)$$

静的特徴量に関する項

動的特徴量に関する項

動的特徴量を考慮しない場合:

※静的・動的特徴量の相互共分散は無視

$$\begin{aligned} \hat{y}_{(0)} &= P_{(0)} r_{(0)} \\ &= \mu_{\hat{m}}^{(y|X)} \end{aligned} \quad \begin{cases} P_{(0)} = \Sigma_{\hat{m}}^{(y|X)} \\ r_{(0)} = \Sigma_{\hat{m}}^{(y|X)^{-1}} \mu_{\hat{m}}^{(y|X)} \end{cases}$$

← フレーム毎に独立に計算可能

フレーム  $t$  までの動的特徴量を考慮した場合:

$$\hat{y}_{(t)} = P_{(t)} r_{(t)} \quad \begin{cases} P_{(t)} = \left( P_{(t-1)}^{-1} + w_t^T \Sigma_{\hat{m}_t}^{(\Delta y|X)^{-1}} w_t \right)^{-1} \\ r_{(t)} = r_{(t-1)} + w_t^T \Sigma_{\hat{m}_t}^{(\Delta y|X)^{-1}} \mu_{\hat{m}_t}^{(\Delta y|X)} \end{cases}$$

← フレーム  $t-1$  の結果を用いて再帰的に計算可能

# 逆行列補題の適用

行列  $A$  とその逆行列  $A^{-1}$  が既知の時, 低ランクの変動  $\nu N \nu^T$  が行列  $A$  に加わった際のその逆行列  $(A + \nu N \nu^T)^{-1}$  は次式で計算可能

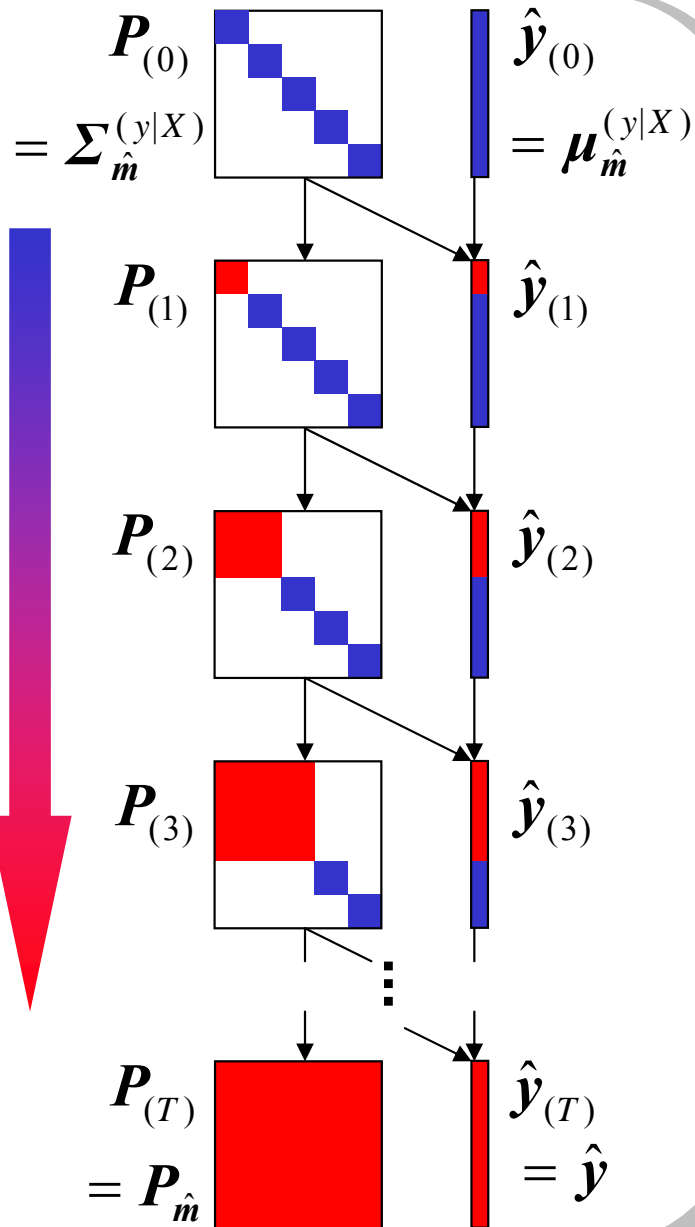
$$(A + \nu N \nu^T)^{-1} = A^{-1} - A^{-1} \nu (N^{-1} + \nu^T A^{-1} \nu)^{-1} \nu^T A^{-1}$$

$$\begin{aligned} P_{(t)} &= \left( P_{(t-1)}^{-1} + \mathbf{w}_t^T \Sigma_{\hat{m}_t}^{(\Delta y|X)^{-1}} \mathbf{w}_t \right)^{-1} = P_{(t-1)} - \underbrace{P_{(t-1)} \mathbf{w}_t^T \left( \Sigma_{\hat{m}_t}^{(\Delta y|X)} + \mathbf{w}_t P_{(t-1)} \mathbf{w}_t^T \right)^{-1} \mathbf{w}_t P_{(t-1)}}_{= \mathbf{k}_{(t)} \text{ とおく.}} \\ &= (I - \mathbf{k}_{(t)} \mathbf{w}_t) P_{(t-1)} \end{aligned}$$

$$\begin{aligned} \hat{\mathbf{y}}_{(t)} &= P_{(t)} \left( \mathbf{r}_{(t-1)} + \mathbf{w}_t^T \Sigma_{\hat{m}_t}^{(\Delta y|X)^{-1}} \boldsymbol{\mu}_{\hat{m}_t}^{(\Delta y|X)} \right) \\ &= (I - \mathbf{k}_{(t)} \mathbf{w}_t) P_{(t-1)} \left( \mathbf{r}_{(t-1)} + \mathbf{w}_t^T \Sigma_{\hat{m}_t}^{(\Delta y|X)^{-1}} \boldsymbol{\mu}_{\hat{m}_t}^{(\Delta y|X)} \right) \\ &= \underbrace{(I - \mathbf{k}_{(t)} \mathbf{w}_t)}_{= \mathbf{k}_{(t)} \boldsymbol{\mu}_{\hat{m}_t}^{(\Delta y|X)} \text{ になる.}} \hat{\mathbf{y}}_{(t-1)} + \underbrace{(I - \mathbf{k}_{(t)} \mathbf{w}_t) P_{(t-1)} \mathbf{w}_t^T \Sigma_{\hat{m}_t}^{(\Delta y|X)^{-1}} \boldsymbol{\mu}_{\hat{m}_t}^{(\Delta y|X)}}_{= \mathbf{k}_{(t)} \boldsymbol{\mu}_{\hat{m}_t}^{(\Delta y|X)} \text{ になる.}} \\ &= \hat{\mathbf{y}}_{(t-1)} + \mathbf{k}_{(t)} \left( \boldsymbol{\mu}_{\hat{m}_t}^{(\Delta y|X)} - \mathbf{w}_t \hat{\mathbf{y}}_{(t-1)} \right) \end{aligned}$$

# 更新式

順次動的特徴量の影響を反映



## $k_{(t)}$ (カルマンゲイン) の計算

$$k_{(t)} = P_{(t-1)} \mathbf{w}_t^T \left( \Sigma_{\hat{m}_t}^{(\Delta y|X)} + \mathbf{w}_t P_{(t-1)} \mathbf{w}_t^T \right)^{-1}$$

## $P_{(t)}$ (誤差の共分散) の更新

$$P_{(t)} = (I - k_{(t)} \mathbf{w}_t) P_{(t-1)}$$

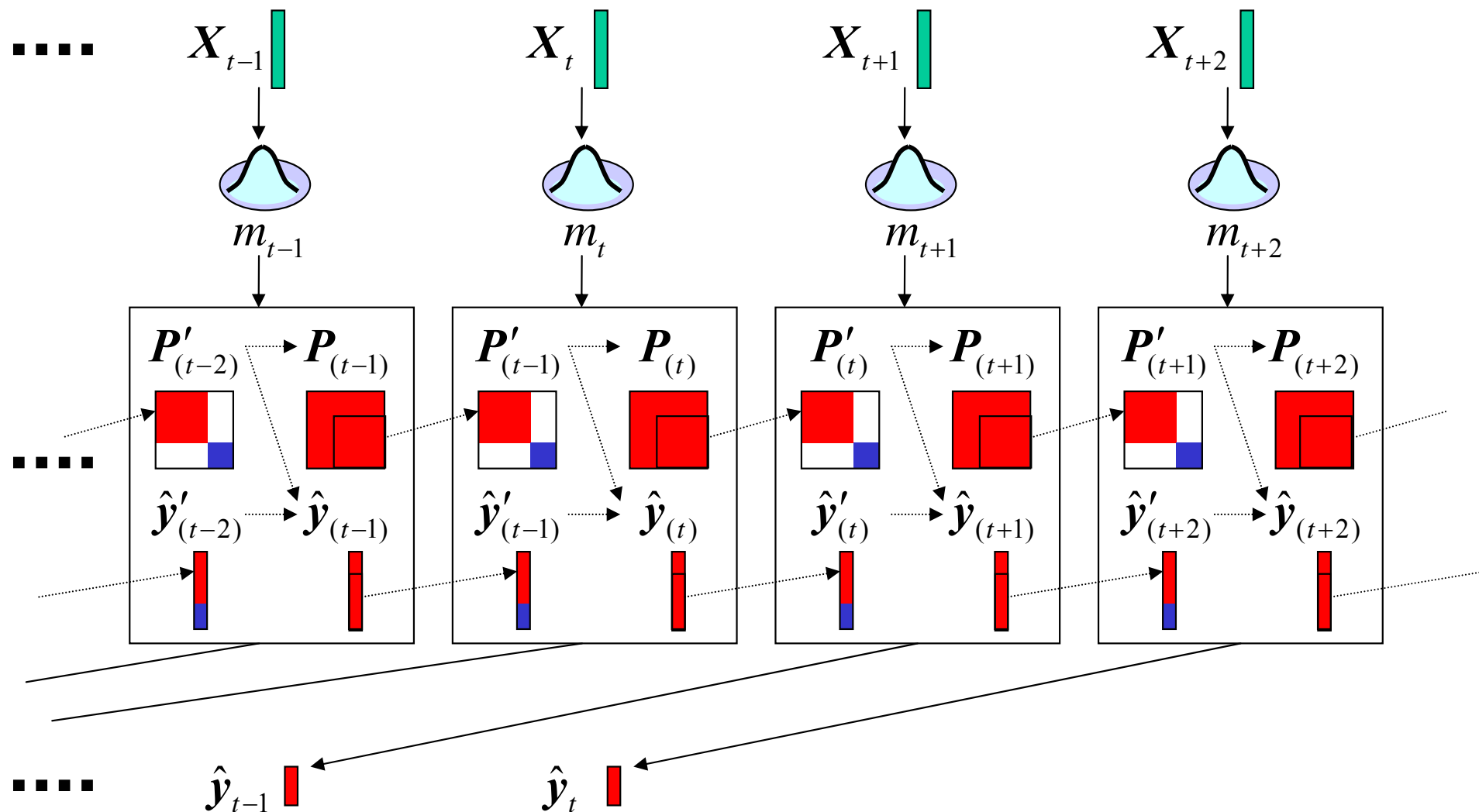
## $\hat{y}_{(t)}$ (推定状態) の更新

$$\hat{y}_{(t)} = \hat{y}_{(t-1)} + k_{(t)} \left( \mu_{\hat{m}_t}^{(\Delta y|X)} - \mathbf{w}_t \hat{y}_{(t-1)} \right)$$

# 短遅延変換処理

[Muramatsu *et al.*]

動的特徴量が与える影響の範囲を数フレームに限定することで  
短遅延変換処理が実現可能



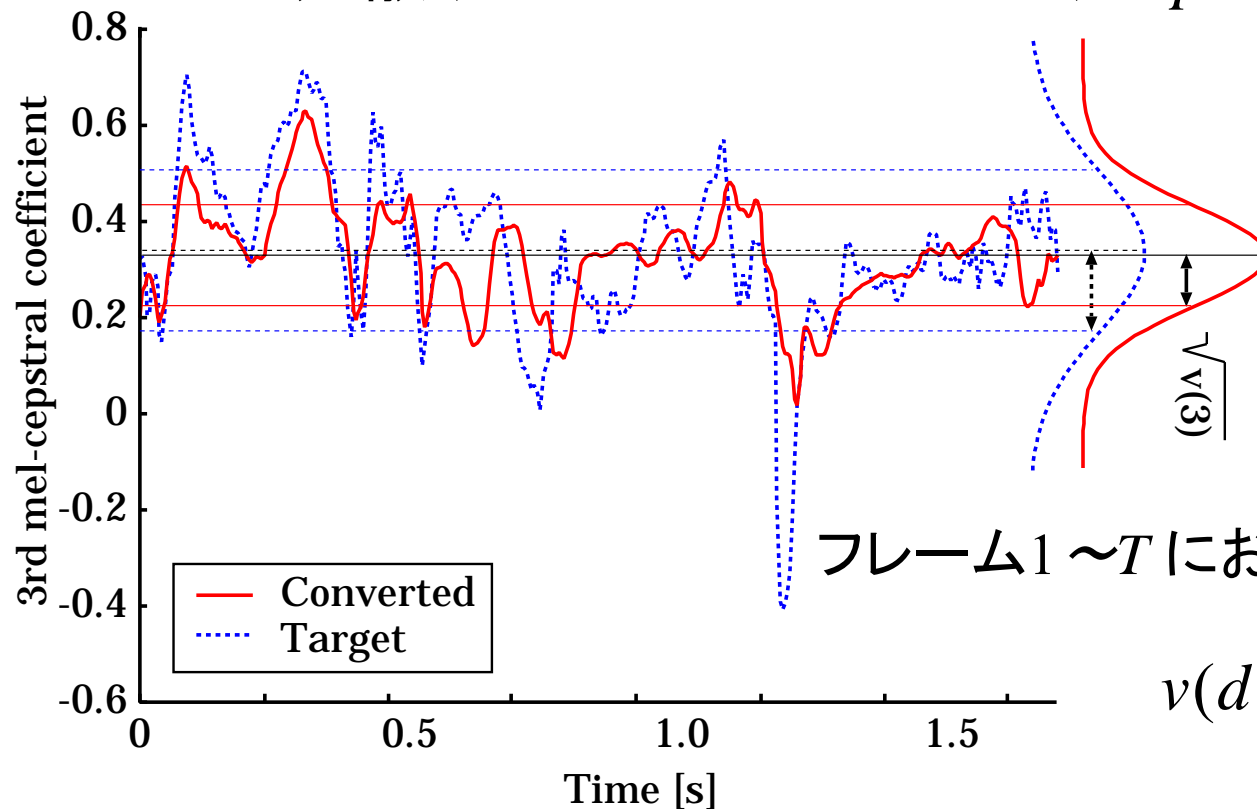
# 最尤系列変換法

[Toda *et al.*]

- 系列に基づく特徴量を考慮した変換処理を行う.
  1. 動的特徴量を考慮した系列ベースの最尤推定
    - フレーム間相関を考慮した変換を行う.
    - 問題点1 (時間的依存関係の無視)を解決できる.
  - 2. 系列内変動の明示的なモデル化の導入**
    - 2次モーメントを考慮した変換を行う.
    - 問題点2 (過剰な平滑化)の影響を大幅に抑えることができる.

## 2. 系列内変動モデルの導入

- 系列内(例えば一発話)に含まれる全フレームにわたって計算される分散(Global Variance: GV)の*p.d.f.*をモデル化する.



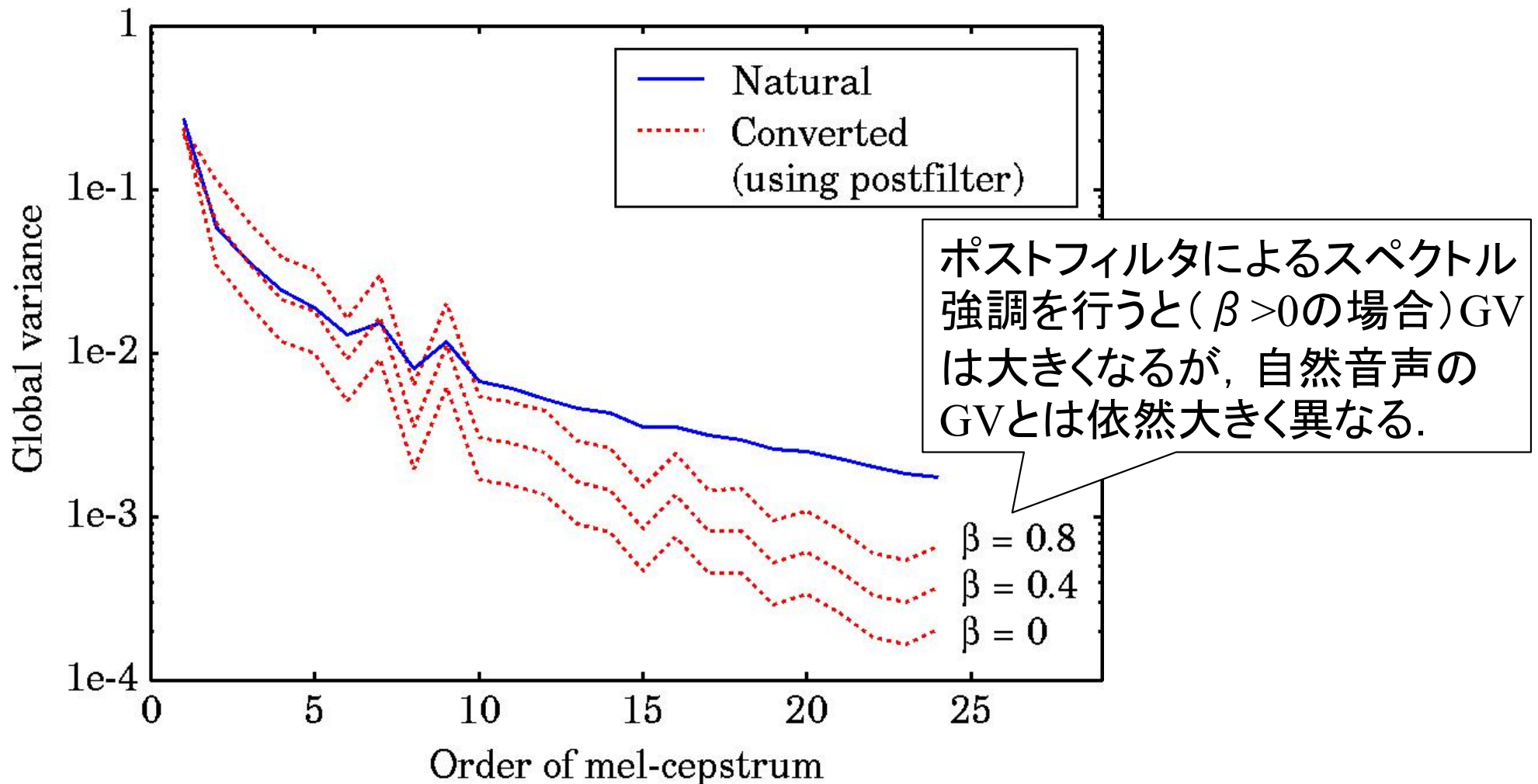
フレーム1 ~  $T$  において計算される  $d$ 次元目のGV:

$$v(d) = \frac{1}{T} \sum_{t=1}^T \left( y_t^{(d)} - \frac{1}{T} \sum_{\tau=1}^T y_{\tau}^{(d)} \right)^2$$

- 変換時には静的特徴量とGV間の明示的な関係を考慮して、条件付*p.d.f.*及びGVの*p.d.f.*の尤度最大化に基づき入力特徴量ベクトルを変換する.

# GVの統計的特徴

- 一般に、統計的変換処理で得られる変換特徴量系列のGVは減少する。



# GVを考慮した最尤系列変換法

[Toda et al.]

- GVの確率密度も考慮した尤度関数を用いる.

$$L(\mathbf{y}) = \log \left\{ P(W\mathbf{y} | \mathbf{X}, \lambda)^\omega \cdot P(\mathbf{v}(\mathbf{y}) | \lambda_v) \right\}$$

静的・動的特徴量  
系列に関する尤度

静的特徴量系列のGV  
 $\mathbf{v}(\mathbf{y})$  に関する尤度

$\omega$  : 2つの尤度のバランスを調節する重み

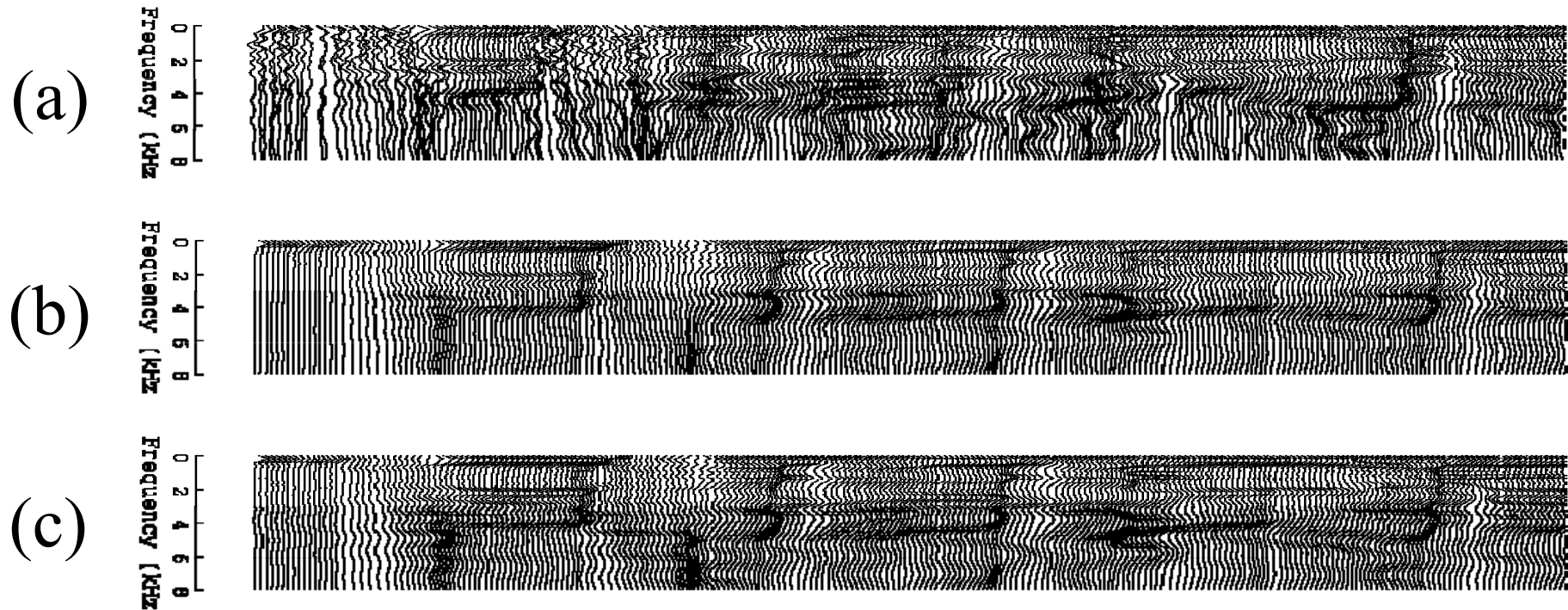
- 勾配法により尤度最大化を行う.

$$\hat{\mathbf{y}} = \arg \max L(\mathbf{y})$$

$$\text{最急降下法: } \mathbf{y}^{(i+1)\text{-th}} = \mathbf{y}^{(i)\text{-th}} + \alpha \cdot \left. \frac{\partial L(\mathbf{y})}{\partial \mathbf{y}} \right|_{\mathbf{y}=\mathbf{y}^{(i)\text{-th}}}$$































# GVを考慮する効果



- (a) 目標スペクトル
- (b) 変換スペクトル(GV未使用)
- (c) 変換スペクトル(GV使用)

# 話者変換音声サンプル

- 話者:4名(男:bdl, 男:rms, 女:clb, 女:slt)
- 学習:50文対(完全自動学習)
- 変換法
  - 左側:最小平均自乗誤差推定法(フレームベース変換法)
  - 右側:GVありの最尤系列変換法

		目標話者							
		bdl		rms		clb		slt	
元話者	bdl								
	rms								
	clb								
	slt								

# 内容

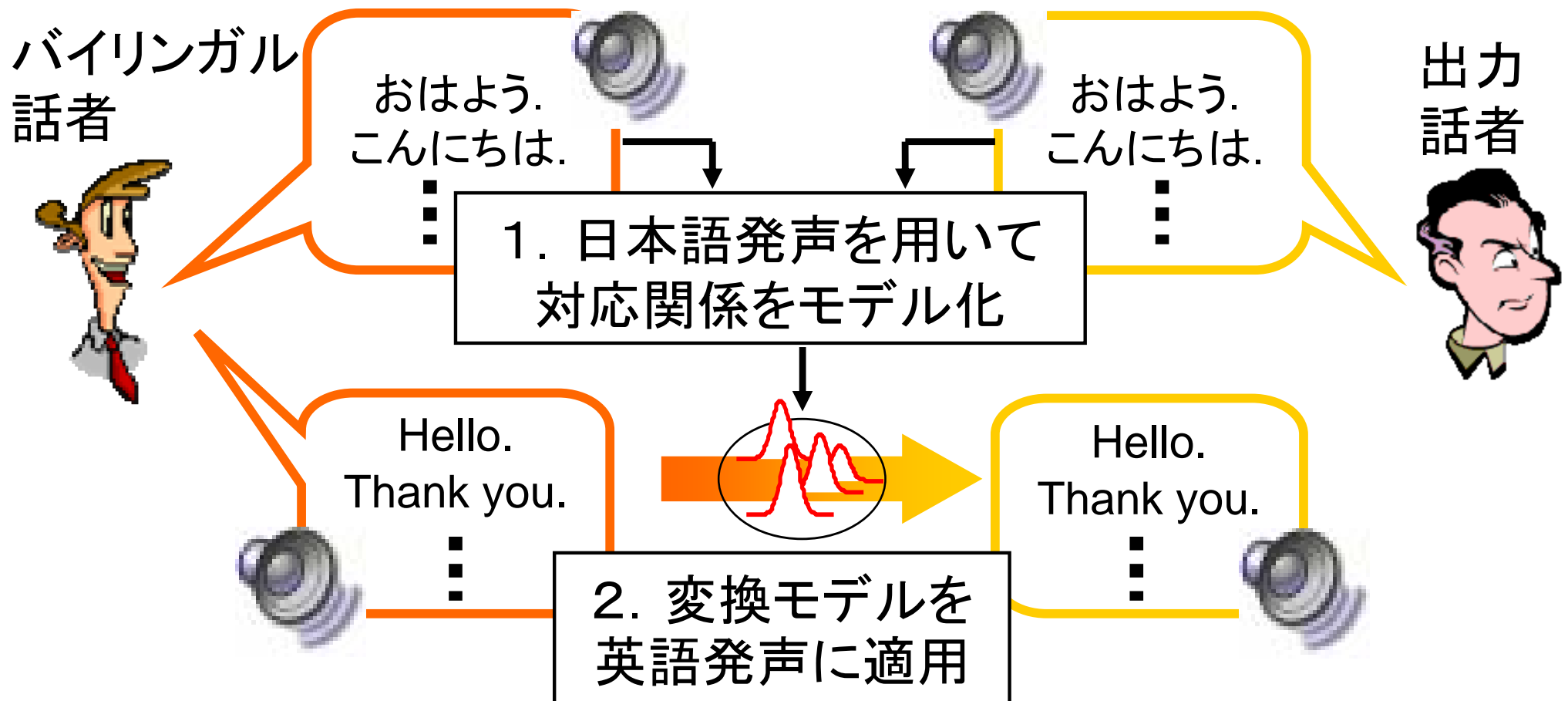
---

1. 音声変換のしくみ
2. 統計的手法による声質変換
  - 2.1. 基本的な枠組み
  - 2.2. フレームベース変換法
  - 2.3. 系列ベース変換法
3. 応用例

# 1. 異なる言語間における話者変換

[Abe et al.]

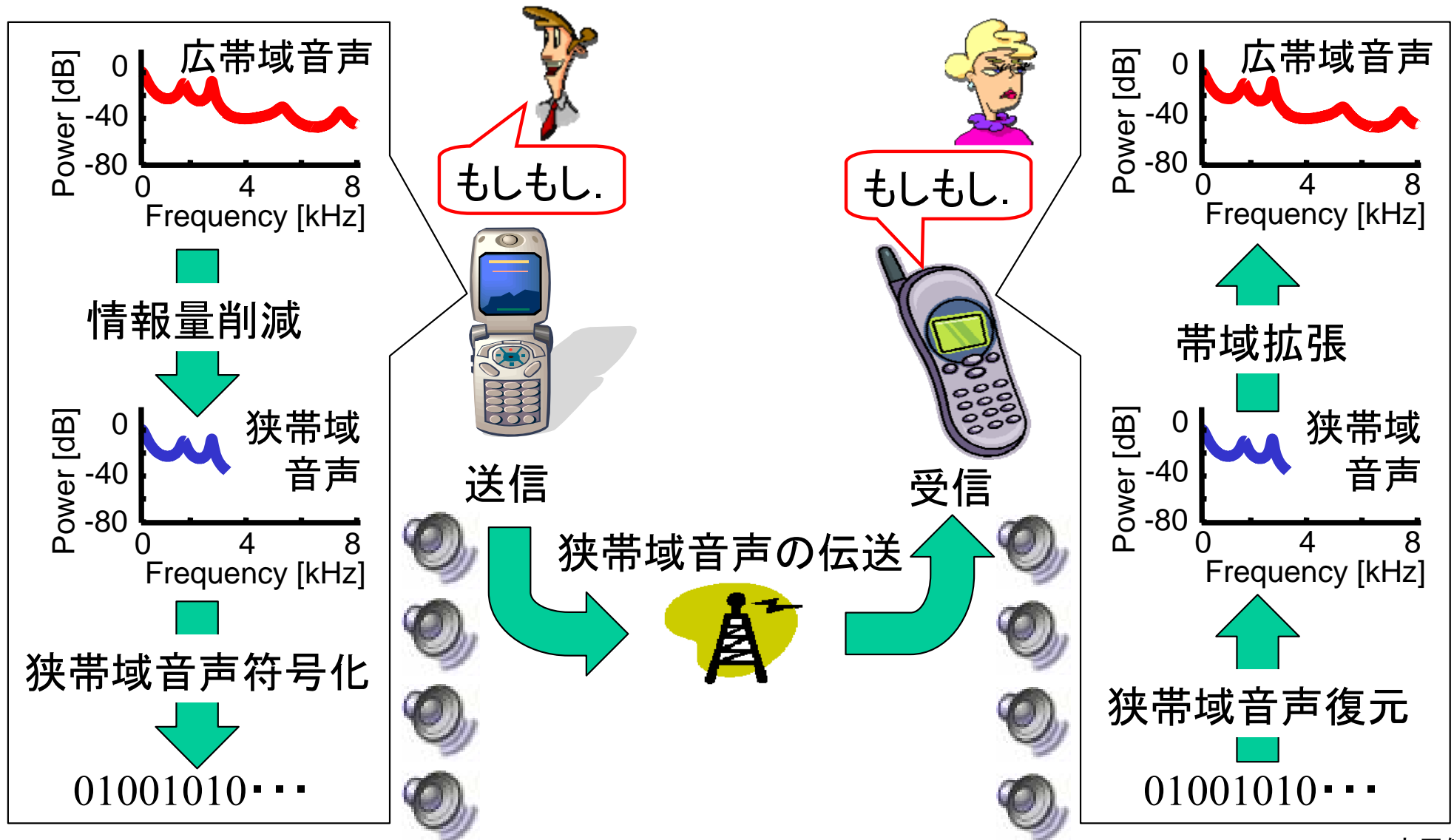
- 自身の声で他の言語を発声することができる。
- 自動音声翻訳, 映画の吹き替え, 外国語学習に有用である。



# 2. 携帯電話音声の帯域拡張

[Jax and Vary]

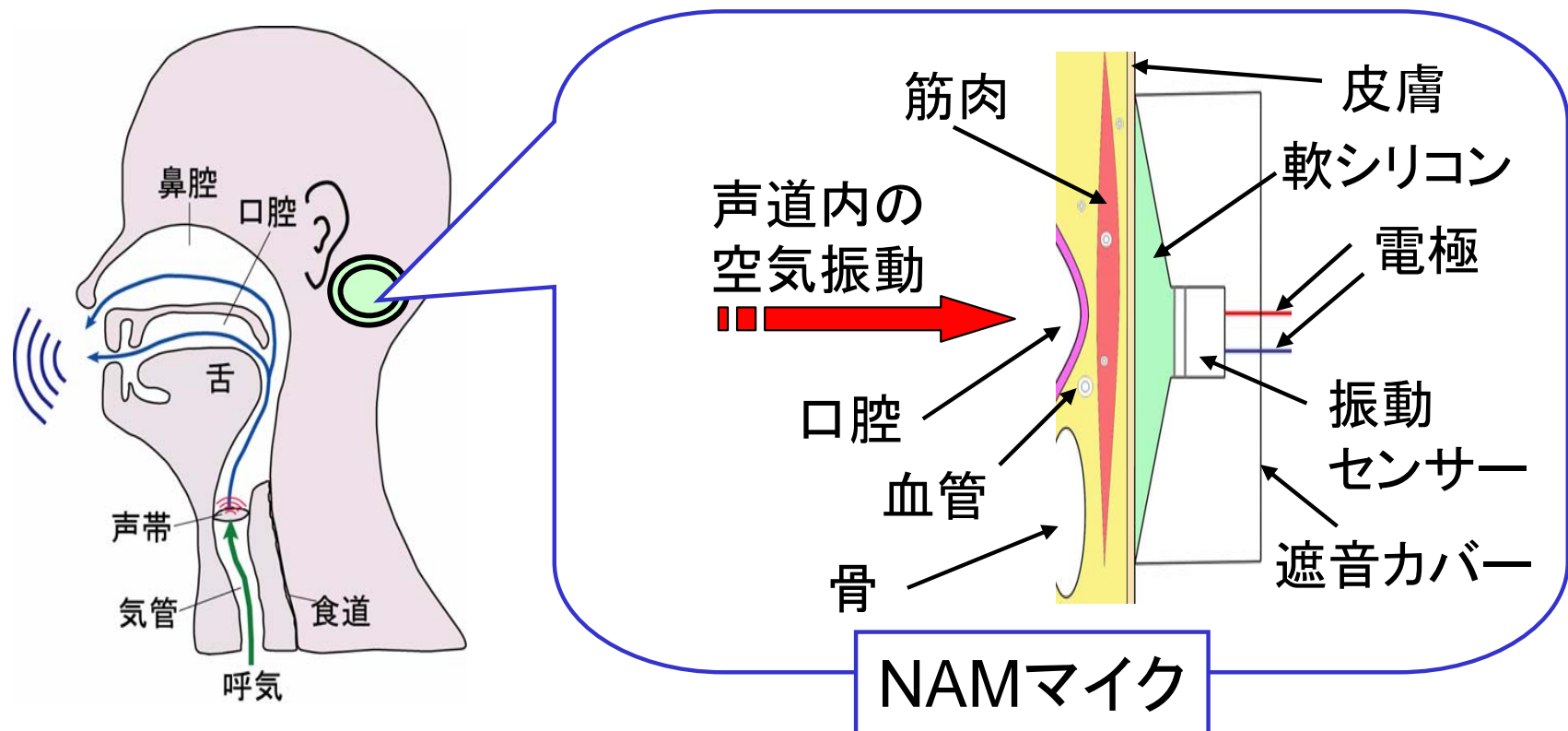
- 従来の伝送量で高品質な広帯域音声を受聴できる。



# 3. 肉伝導音声コミュニケーション

[中島 他]

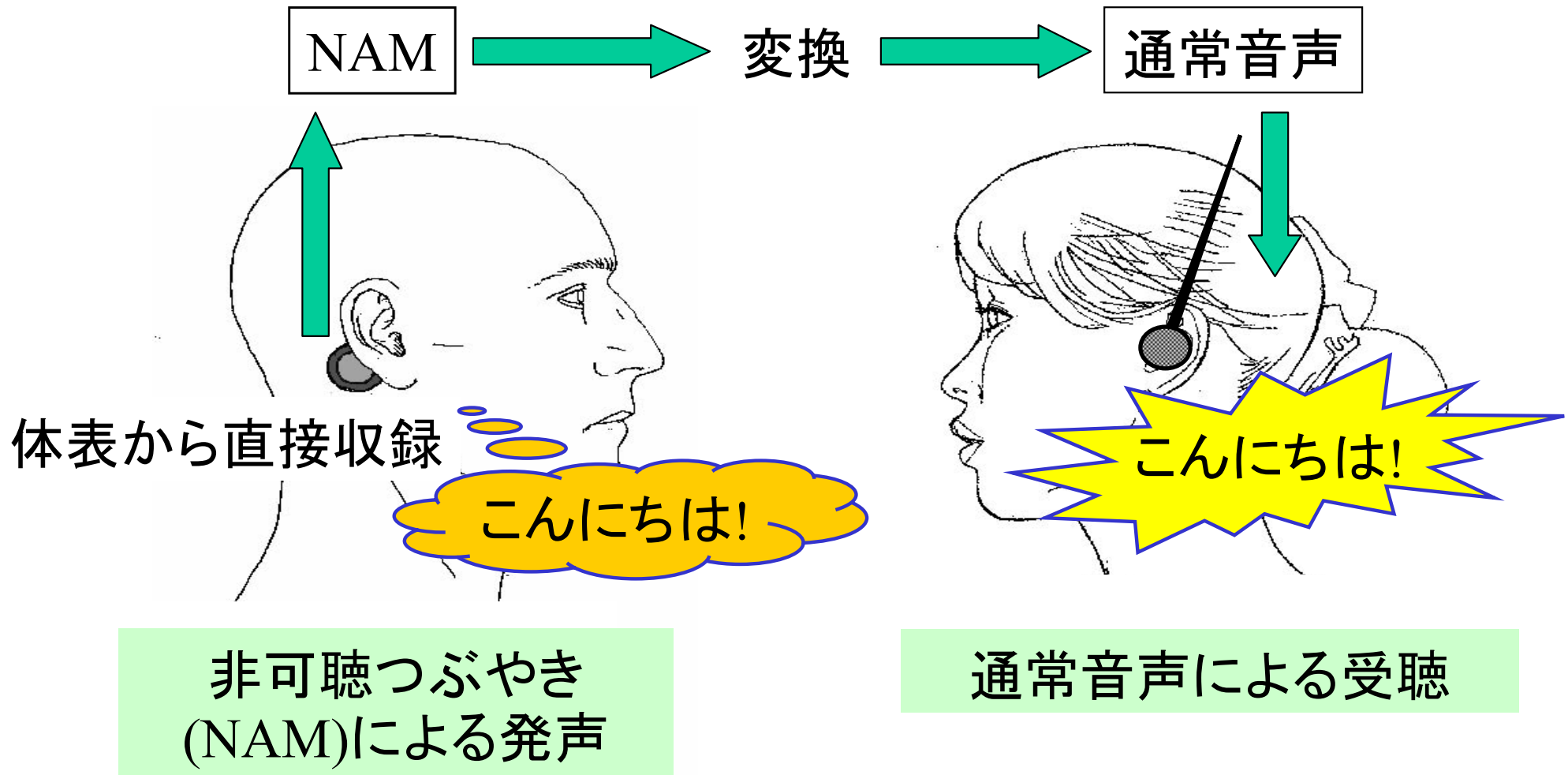
- NAMマイクを用いることで、軟組織を伝わる音(肉伝導音)を体表から収録し、それを音声コミュニケーションに用いる。



※非可聴つぶやき (NAM: Non-Audible Murmur) :  
周りに聞こえないほど小さなつぶやき声のようなものであり、  
呼気が調音されることで生成される無声音である。

# 応用例：無音声電話

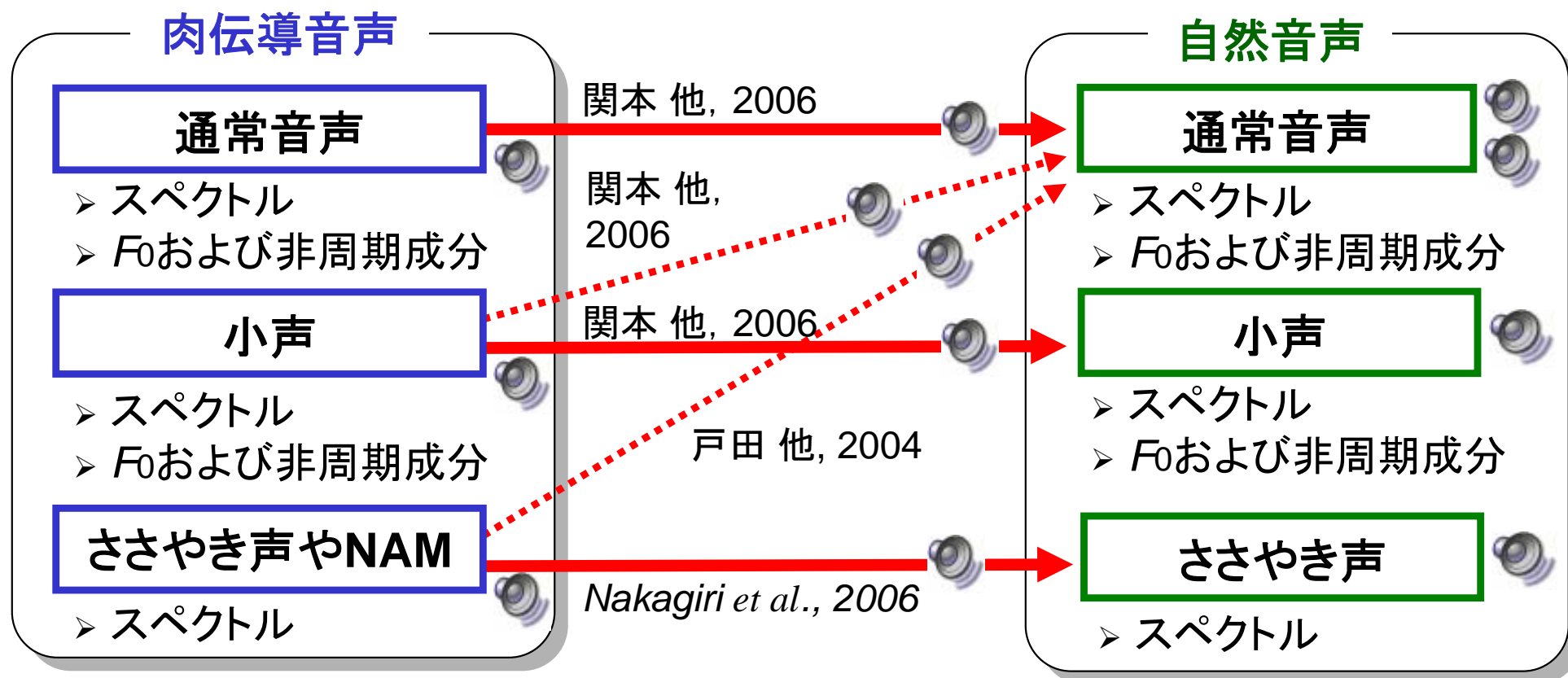
- 第三者には聞こえない音声コミュニケーションを実現する。



# 肉伝導音声の変換法

[Toda et al.]

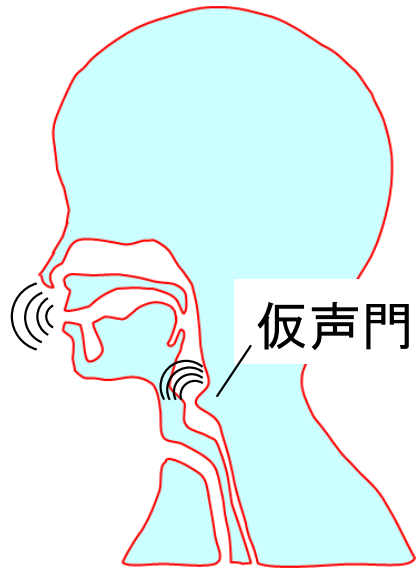
- NAMマイクは様々な肉伝導音声を収録することができ、各種肉伝導音声は様々な用途に使用できる。





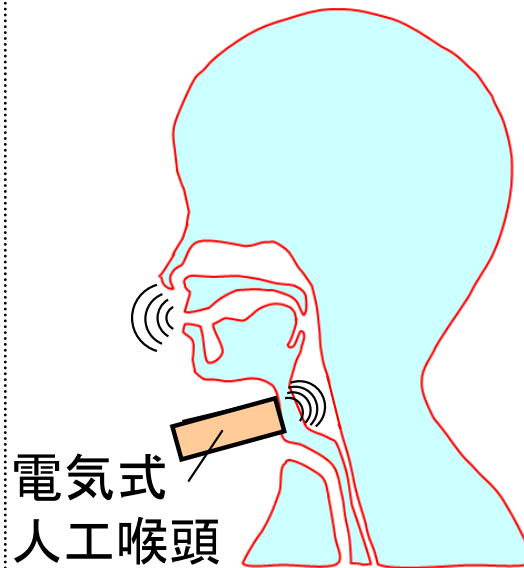
# 4. 発声障害者補助

## 食道発声



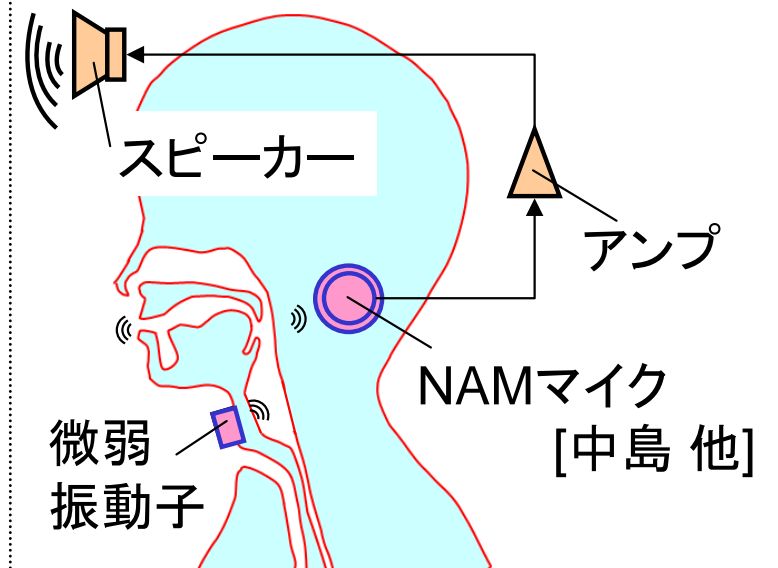
## 電気式人工喉頭を用いた発声

[橋場 他]



## 微弱振動子及びNAMマイクを用いた発声

[中村 他]

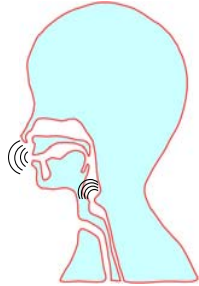


補助器具	不要	必要	必要
習得	困難	容易	比較的容易
伝達される音信号	無喉頭音声	無喉頭音声及び人工喉頭の音源信号	スピーカーから提示される無喉頭音声
無喉頭音声	食道音声	電気音声	肉伝導微弱電気音声

# 無喉頭音声の変換法

## 食道音声の変換

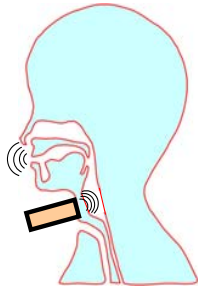
[Doi *et al.*, 2009–2010]



## 電気音声の変換

[Nakamura *et al.*, 2009–2010]

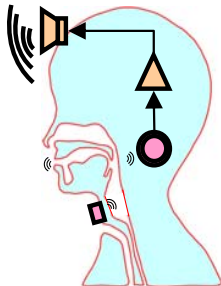
[Doi *et al.*, 2010]




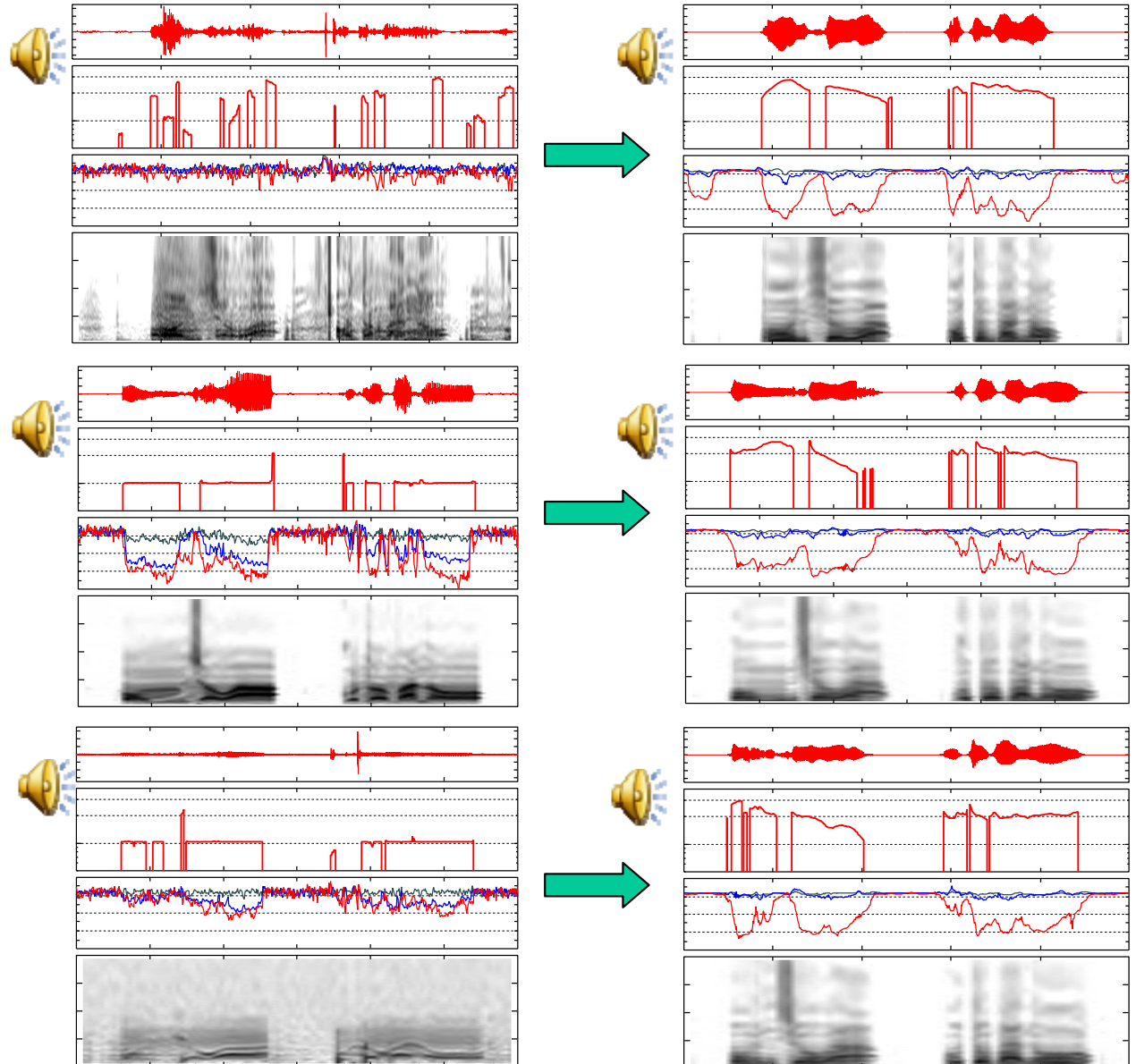
## 肉伝導微弱電気音声の変換

[Nakamura *et al.*, 2007–2010]

[Doi *et al.*, 2010]



目標話者: 



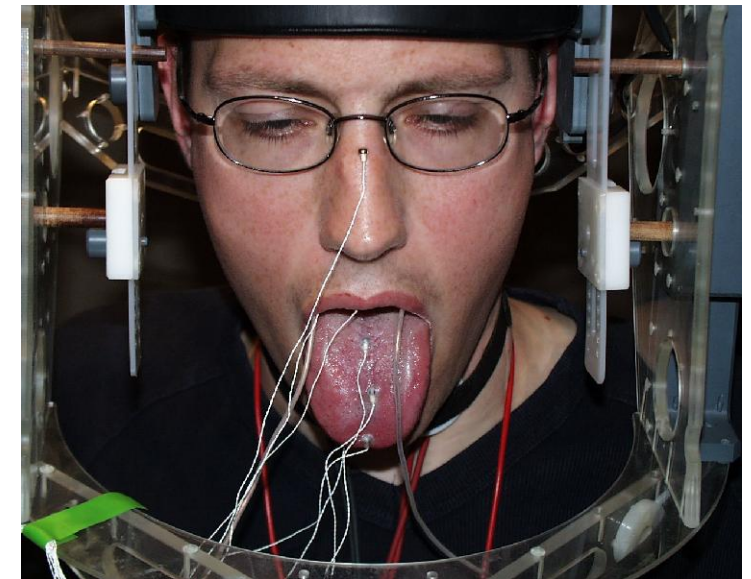
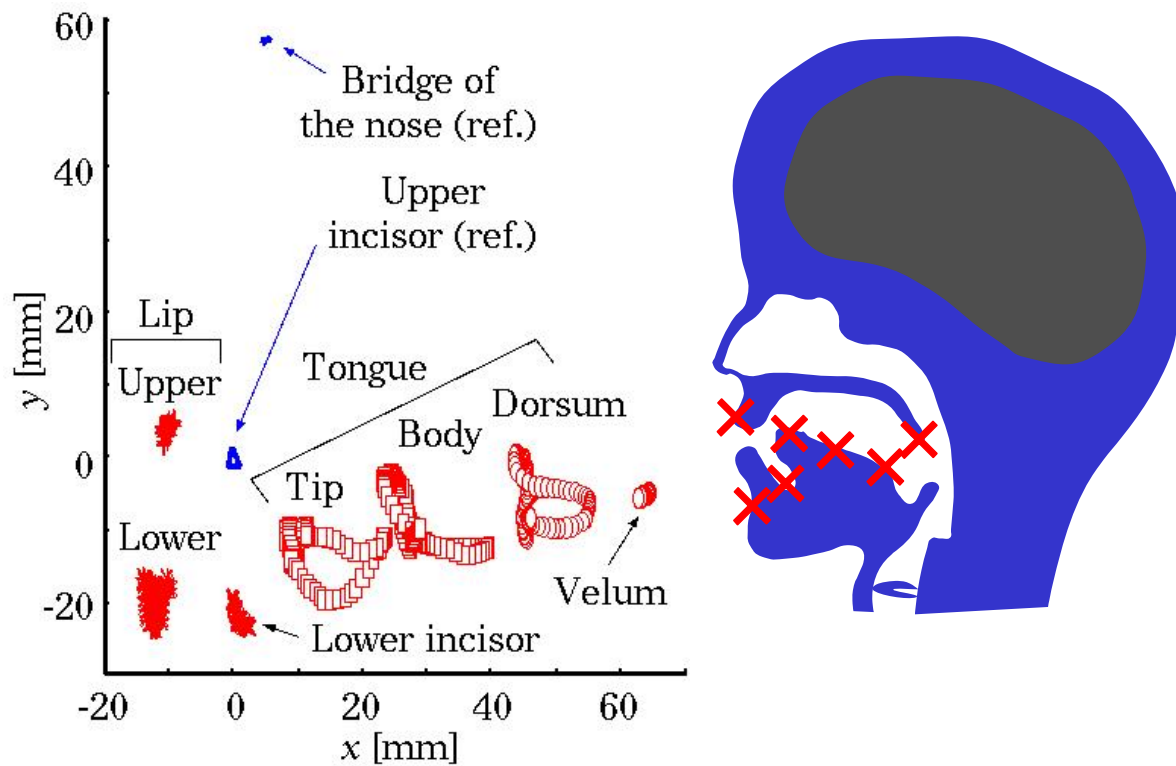
# 5. 調音運動制御による音声変換

[Toda *et al.*]

- 調音運動制御による音声変換を実現する.

手順1. 調音運動と音声信号を同時に収録する.

手順2. 調音運動と音声特徴量の対応関係をモデル化する.



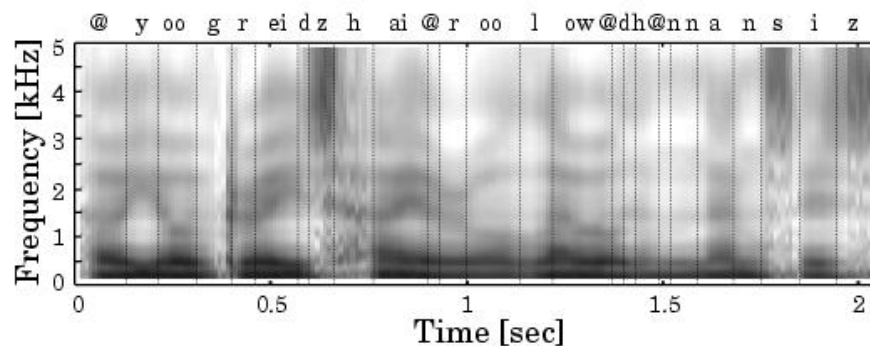
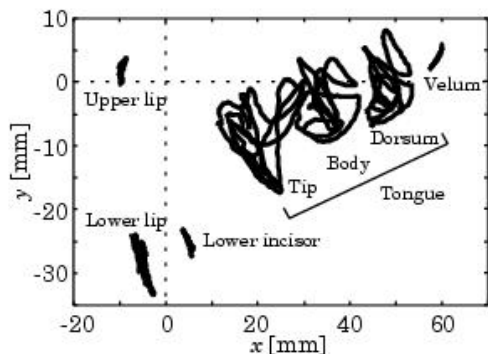
MOCHAデータベース  
Edinburgh大から公開

<http://www.cstr.ed.ac.uk/research/projects/artic/mocha.html>

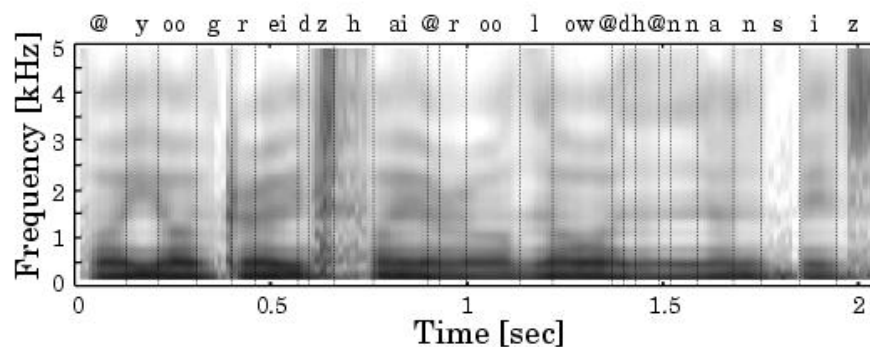
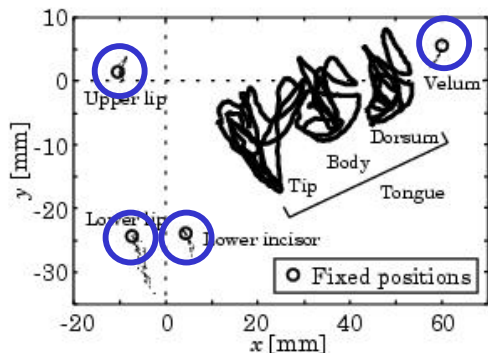
Electromagnetic articulograph (EMA)データ

# 調音運動制御による音声変換の一例

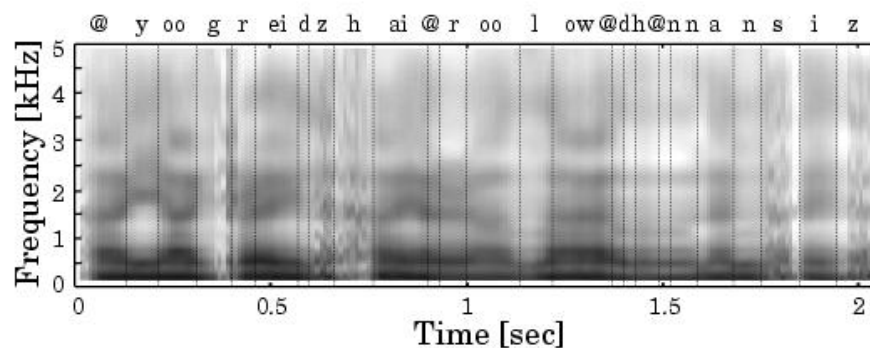
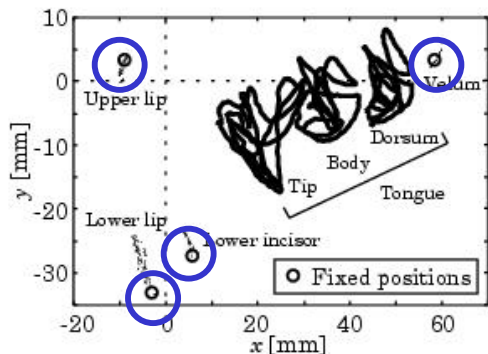
## 1. 自然な調音運動からの合成音



## 2. 口をすぼめたままの合成音



## 3. 口を開けたままの合成音





# 参考文献

- M. Abe, S. Nakamura, K. Shikano, H. Kuwabara, “Voice conversion through vector quantization,” *J. Acoust. Soc. Jpn. (E)*, vol. 11, no. 2, pp. 71–76, 1990.
- 中村哲, 鹿野清宏, “ファジィベクトル量子化を用いたスペクトログラムの正規化,” *日本音響学会誌*, vol. 45, no. 2, pp. 107–114, 1989.
- H. Matsumoto and Y. Yamashita, “Unsupervised speaker adaptation from short utterances based on a minimized fuzzy objective function,” *J. Acoust. Soc. Jpn. (E)*, vol. 14, no. 5, pp. 353–361, 1993.
- H. Valbret, E. Moulines, and J. P. Tubach, “Voice transformation using PSOLA technique,” *Speech Communication*, vol. 11, no. 2–3, pp. 175–187, 1992.
- Y. Stylianou, O. Cappe, E. Moulines, “Continuous probabilistic transform for voice conversion,” *IEEE Trans. Speech Audio Process.*, vol. 6, no. 2, pp. 131–142, 1998.
- A. Kain, M. W. Macon, “Spectral voice conversion for text-to-speech synthesis,” *Proc. ICASSP*, Seattle, USA, pp. 285–288, May 1998.
- T. Toda, A.W. Black, K. Tokuda, “Voice conversion based on maximum likelihood estimation of spectral parameter trajectory,” *IEEE Trans. Audio, Speech and Language Process.*, vol. 15, no. 8, pp. 2222–2235, 2007.
- Heiga Zen, Keiichi Tokuda, Alan W. Black, “Statistical parametric speech synthesis,” *Speech Communication*, vol.51, no.11, pp.1039–1154, 2009.
- K. Tokuda, T. Yoshimura, T. Masuko, T. Kobayashi, and T. Kitamura, “Speech parameter generation algorithms for HMM-based speech synthesis,” *Proc. ICASSP*, pp. 1315–1318, Istanbul, Turkey, June 2000.
- H. Zen, K. Tokuda, and T. Kitamura, “Reformulating the HMM as a trajectory model by imposing explicit relationships between static and dynamic feature vector sequences,” *Computer Speech and Language*, vol. 21, no. 1, pp. 153–173, 2007.
- 徳田恵一, 益子貴史, 小林隆夫, 今井 聖, “動的特徴を用いた HMMからの音声パラメータ生成アルゴリズム,” *日本音響学会誌*, vol.53, no.3, pp.192–200, Mar. 1997.
- T. Muramatsu, Y. Ohtani, T. Toda, H. Saruwatari, K. Shikano, “Low-delay voice conversion based on maximum likelihood estimation of spectral parameter trajectory,” *Proc. INTERSPEECH*, pp. 1076–1079, Brisbane, Australia, Sep. 2008.

# 参考文献

- M. Abe, K. Shikano, and H. Kuwabara, “Statistical analysis of bilingual speaker’s speech for cross-language voice conversion,” *J. Acoust. Soc. Am.*, vol. 90, no. 1, pp. 76–82, 1991.
- M. Mashimo, T. Toda, H. Kawanami, K. Shikano, N. Campbell, “Cross-language voice conversion evaluation using bilingual databases,” *IPSJ Journal*, vol. 43, no. 7, pp. 2177–2185, 2002.
- P. Jax, P. Vary, “On artificial bandwidth extension of telephone speech,” *Signal Processing*, vol. 83, pp. 1707–1719, 2003.
- 中島淑貴, 柏岡秀紀, ニック キャンベル, 鹿野清宏, “非可聴つぶやき認識,” *電子情報通信学会論文誌*, vol. J87-D-II, no. 9, pp. 1757–1764, 2004.
- T. Toda, K. Nakamura, H. Sekimoto, K. Shikano, “Voice conversion for various types of body transmitted speech,” *Proc. ICASSP*, pp. 3601–3604, Taipei, Taiwan, Apr. 2009.
- M. Nakagiri, T. Toda, H. Kashioka, and K. Shikano, “Improving body transmitted unvoiced speech with statistical voice conversion,” *Proc. INTERSPEECH*, pp. 2270–2273, Pittsburgh, USA, Sep. 2006.
- 中村圭吾, 戸田智基, 猿渡洋, 鹿野清宏, “肉伝導人工音声の変換に基づく喉頭全摘出者のための音声コミュニケーション支援システム,” *電子情報通信学会論文誌*, vol. J90-D, no. 3, pp. 780–787, 2007.
- H. Doi, K. Nakamura, T. Toda, H. Saruwatari, and K. Shikano, “Esophageal speech enhancement based on statistical voice conversion with Gaussian mixture models,” *IEICE Trans. on Inf. and Syst.*, vol. E93-D, no. 9, pp. 2472–2482, 2010.
- T. Toda, A.W. Black, K. Tokuda, “Mapping between articulatory movements and acoustic spectrum with a Gaussian mixture model,” *Speech Communication*, vol. 50, no. 3, pp. 215–227, Mar. 2008.