

スペシャルセッション  
「音楽情報処理と機械学習」

歌声情報処理：声質変換モデル

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8月9日

# 機械学習を用いた 研究の進展

# 歌声／音声変換・合成の定式化

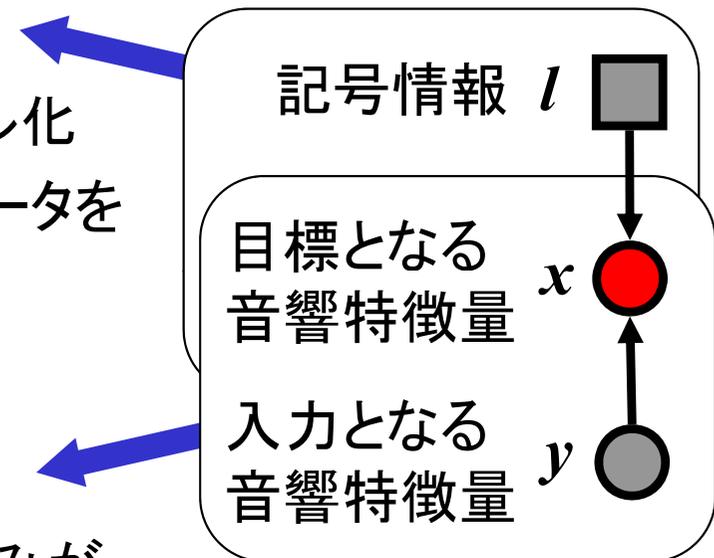
## • 確率モデルによる音響特徴量のモデル化

### – 歌声合成, テキスト音声合成

- 確率密度関数 (p.d.f.)  $P(x|l)$  のモデル化
- 記号情報が付与された歌声／音声データを用いて**自動的に学習**

### – 歌声声質変換, 声質変換

- $P(x|y)$  のモデル化
- 記号情報は同一で所望の声質成分のみが異なる歌声／音声データを用いて**自動的に学習**



HMM, GMM, ニューラルネット,  
トラジェクトリモデル, GP,  
最尤推定, ベイズ推定,  
EMアルゴリズム, 変分法, ...



モデル



データ

# 確率分布からの時系列データ生成法

[徳田 他, 1995]

静的特徴量と**動的特徴量** (各時刻における時間方向の変化量) の関係を導入することで、適切に遷移する時系列データを生成

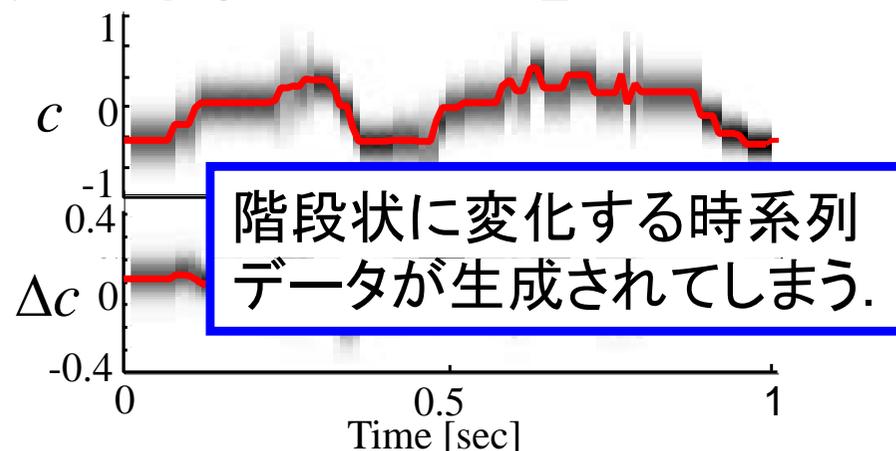
$c$ : 静的特徴量系列

$o$ : 静的・動的特徴量系列

$q$ :  $o$  をモデル化する分布系列

$$\bar{o}_q = \arg \max_o P(o | q, \lambda) \quad \text{静的・動的特徴量}$$

= 平均ベクトル系列  $\text{系列の} p.d.f.$

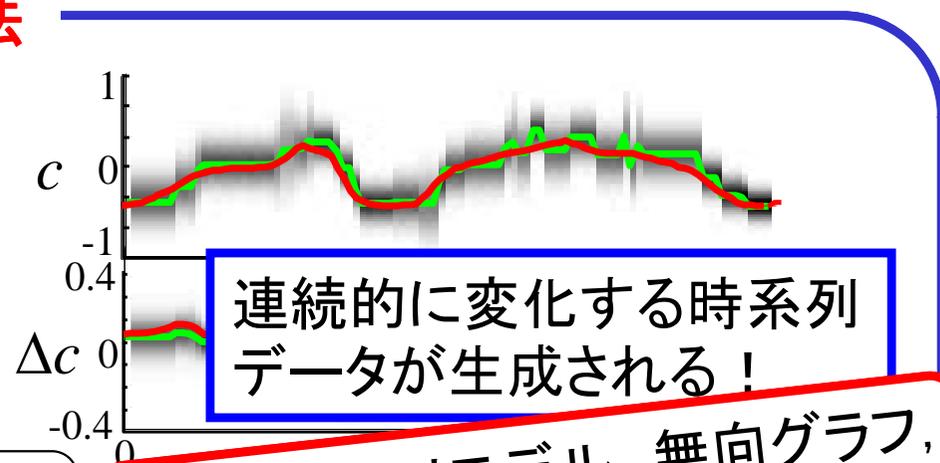


## 動的特徴量を考慮した時系列データ生成法

$$\bar{c}_q = \arg \max_c P(o | q, \lambda)$$
$$= \arg \max_c P(c | q, \lambda) P(\Delta c | q, \lambda)$$

静的特徴量  $\text{系列の} p.d.f.$       動的特徴量  $\text{系列の} p.d.f.$

静的特徴量系列を線形変換  $\Delta c = f_{\Delta}(c)$  した変数の  $p.d.f.$



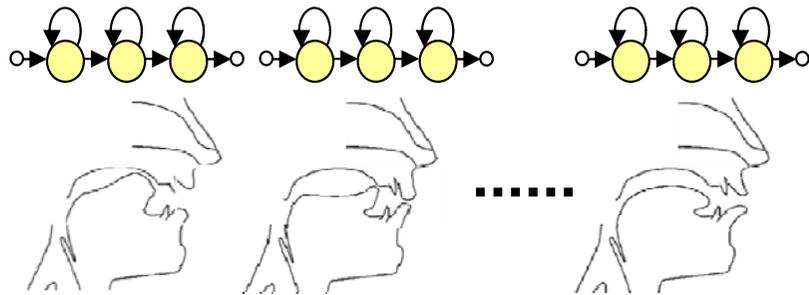
トラジェクトリモデル, 無向グラフ, 線形動的システム, ...

# 変換・合成：もっと音を良くしたい！

2004年当初・・・，GMMに基づく声質変換およびHMM音声合成による合成音声はこもり感が強く，自然性が低かった・・・

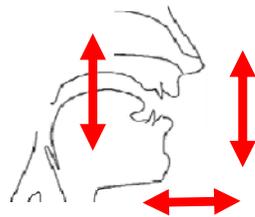


人がHMM尤度最大化基準で発声したらどうなるかな？



口・舌を動きが小さくなりそう・・・  
はっきりと話さなくなりそう・・・

各音韻を発声する際の**平均的**な  
口・舌の形を使って発声する。

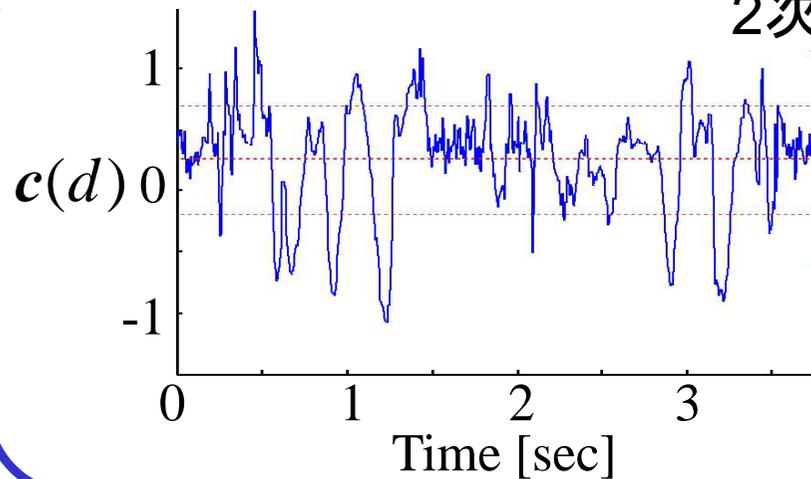


発話全体を通して口・舌をどの程度動かすかという**変動量**という基準も使っている気がする！

# 系列内変動 (GV) を考慮した生成法

[Toda *et al.*, 2004]

## GV (global variance)



2次モーメント (= 静的特徴量系列の**非線形変換**)

$$\mathbf{v}_c = [v_{c(1)}, \dots, v_{c(D)}]^T = f_v(\mathbf{c})$$

$$v_{c(d)} = \sum_{t=1}^T \left( c_t(d) - \frac{1}{T} \sum_{\tau=1}^T c_\tau(d) \right)^2$$

データからGVのp.d.f.も学習 

## GVを考慮した時系列データ生成法

$$\bar{c}_q = \arg \max_c \underbrace{P(\mathbf{c} | \mathbf{q}, \lambda)}_{\text{静的特徴量系列のp.d.f.}} \underbrace{P(\Delta \mathbf{c} | \mathbf{q}, \lambda)}_{\text{動的特徴量系列のp.d.f.}} \underbrace{P(\mathbf{v}_c | \lambda_v)}_{\text{系列内変動のp.d.f.}}$$

静的特徴量  
系列  $c$  の関数

静的特徴量 動的特徴量系列のp.d.f. 系列内変動のp.d.f.

系列のp.d.f. ※線形変換  $\Delta \mathbf{c} = f_\Delta(\mathbf{c})$  ※非線形変換  $\mathbf{v}_c = f_v(\mathbf{c})$

HMM音声合成によるサンプル音声

GVなし: 

GVあり: 

PoE, 正則化, 特徴空間,  
非線形写像, ...

# 変換：学習処理を無くしたい！

2005年当初・・・，声質変換の障害者補助応用を考え始める・・・



目標とする声質の音声が存在しない・・・どうしよう？

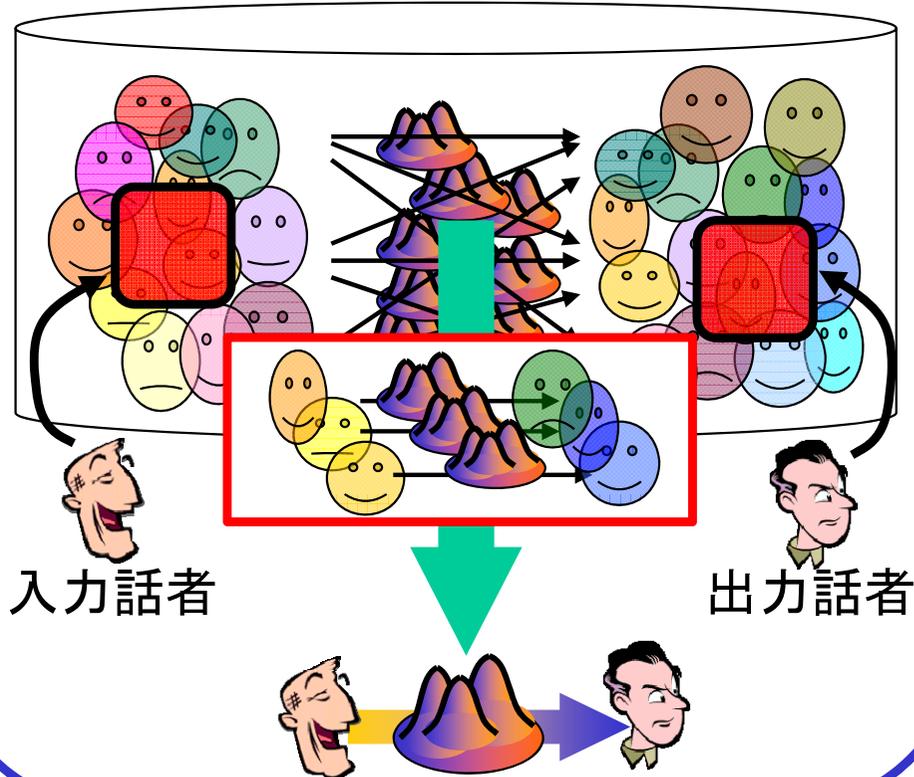
他人の変換モデルで代用できるのでは・・・

多人数集めれば似た声の人はいるのでは・・・

声質コントローラみたいな仕組みを作れないか・・・

音声合成／認識分野における数多くの研究成果を活用できないか・・・

他人の声を上手く混ぜることで別の人の声を創り出そう！



# 固有声変換

[Toda et al., 2006]

学習データ

多数パラレルデータセット

- 参照話者  $x$
- 多数事前収録話者  $y^{(1)}, \dots, y^{(S)}$

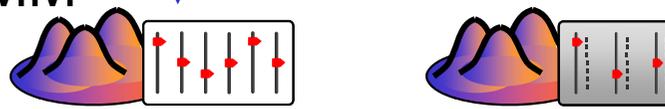
声質表現語スコア

- 各事前収録話者用  $w_c^{(1)}, \dots, w_c^{(S)}$

結合 *p.d.f.*

$$P(x, y^{(s)} | w^{(s)})$$

固有声GMM

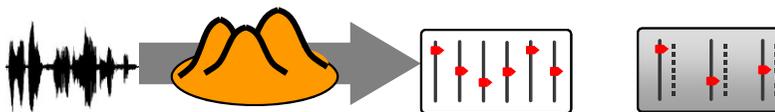


周辺 *p.d.f.*

$$P(y^{(s)} | w^{(s)})$$

声質分析

話者性パラメータ推定



声質表現語  
スコア推定

条件付 *p.d.f.*

$$P(x | y^{(s)}, w^{(s)})$$

$$P(y^{(s)} | x, w^{(s)})$$

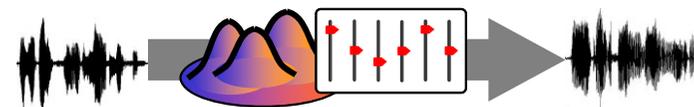
$$\int P(y^{(s)} | x, w^{(s)}) P(x | y^{(s)}, w^{(s)}) dx$$

声質変換

多対一変換

一對多変換

多対多変換



周辺化, 部分空間,  
回帰, カーネル法, ...

設定

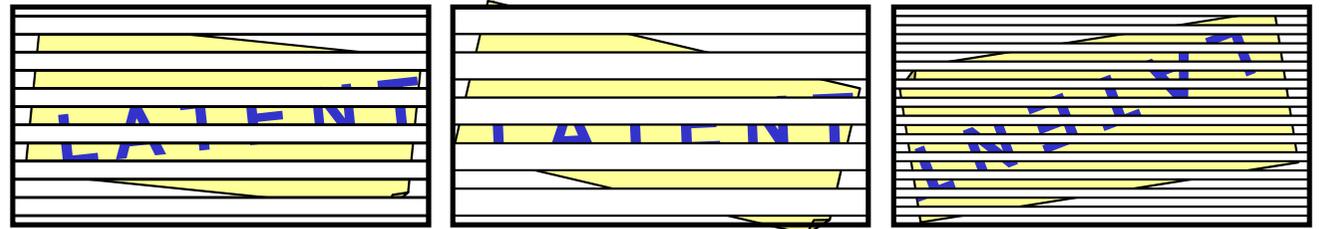
# 分析：統計処理で誤差を減らしたい！

2006年夏・・・，声質変換の分析エラーが気になり，分析合成も統計的にできるのではないかと思い始める・・・

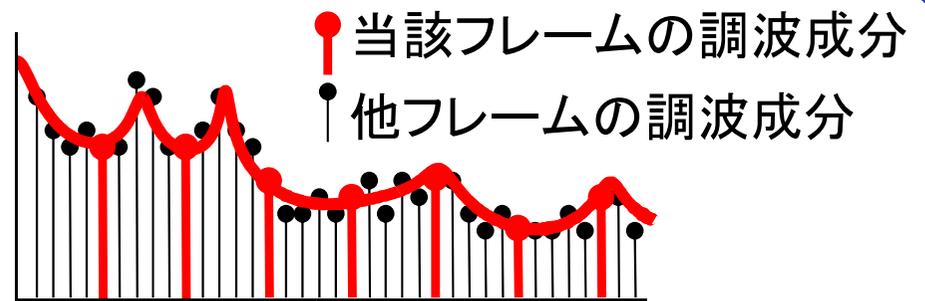
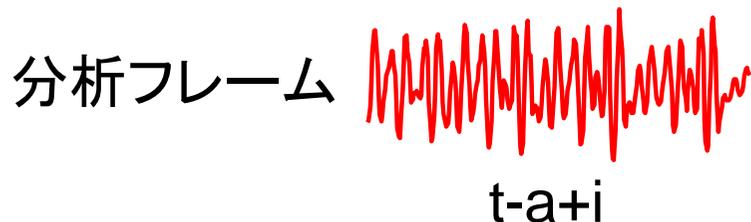


スペクトル分析は見えない箇所を見る問題・・・人は障害物の向こうにある物をどうやってイメージしているんだろう？

別の箇所が見えている他のサンプルも使えば推測できるのでは？



他フレームから得られる情報も使おう！



類似した音素環境を持つ他のフレーム

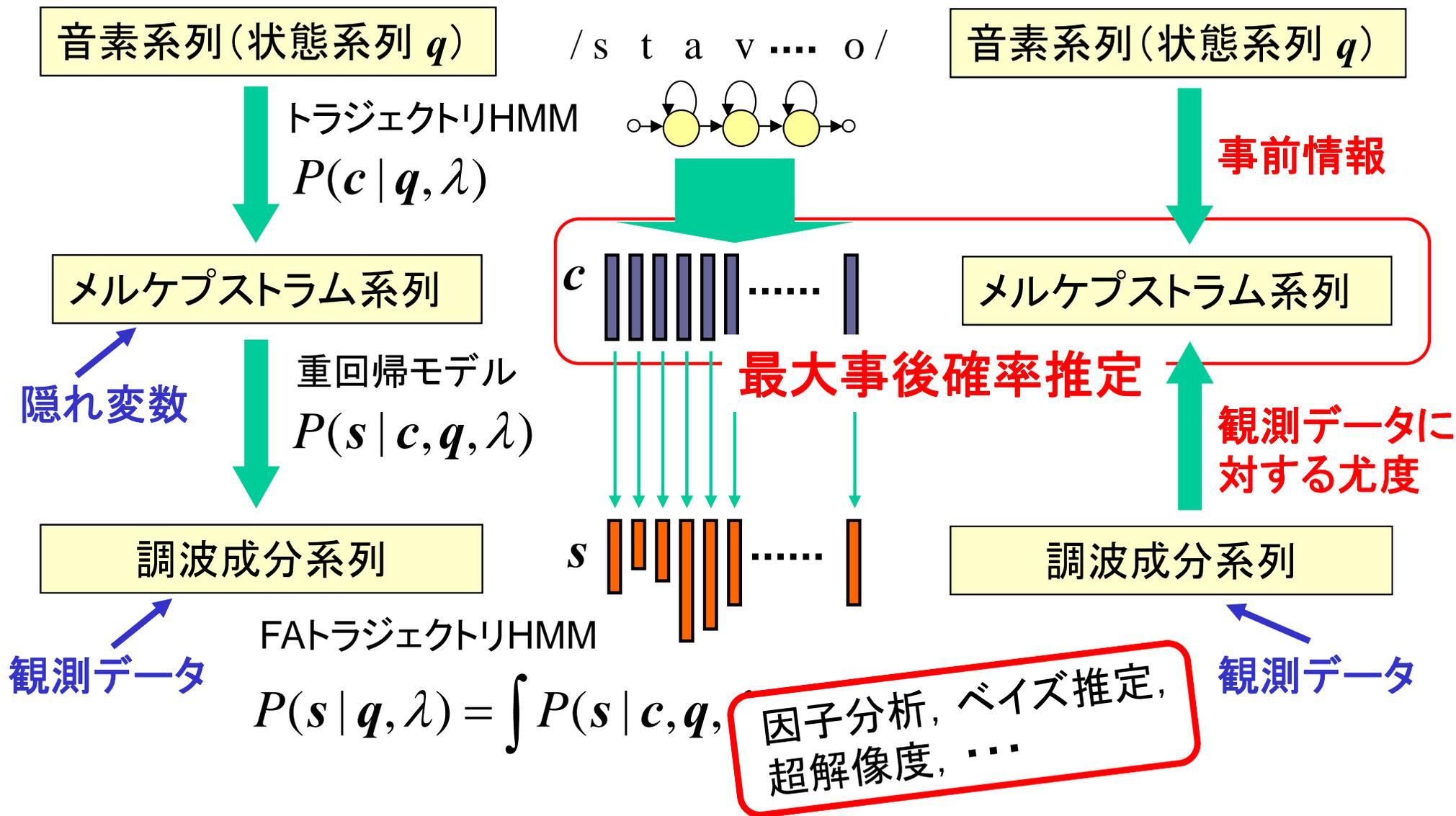


# STAVOCO (STAtistical VOCOrder) のスペクトル分析

[Toda et al., 2007]

## 学習処理

## 推定処理



# 機械学習の 効果と限界

# なぜ機械学習を使うのか？

- **こんなにも良いことがある！**

- 解くべき問題を数理的に記述することで、**高い汎用性と実用性**が得られる！
- 容易に**技術を共有**することができ、色々な問題へと適用できる！
- **データドリブン**な枠組みであるため、システムを**自動的に構築**できる！
- **統計処理**を用いることで、物理的に不可能なことも可能とする技術を生み出すことができる！

**人の能力 + 機械学習 = 世界が変わる！**

# 機械学習を適用する際に重要なことは？

- **適用する問題に対して特化させよう！**

- 応用先の処理(例えば合成)の特徴, モデル化対象となる信号の特徴, 物理的な生成過程などを考慮して, モデルを構築する.
- (機械学習のみに限らず)従来研究の知見を活用する.

- **解決すべき課題は多い！**

- 主観的・知覚的な評価との対応を上手にとるには？
- 如何に対象となるデータを大量に手に入れるか？
  - システム使用時に得られるデータを大量に集める必要あり...

今後の展開および  
実用化に向けて

# 今後の展開および実用化に向けて

- **歌声・音楽情報処理における非言語／パラ言語情報のモデル化技術をさらに発展させていきたい！**
  - **イメージを具現化できる歌声／音声／楽器音合成・変換技術**
    - 全ての音が合成音や変換音で作られた芸術作品(映画など)も可能に？
  - 人の主観値に基づくデータを自動的に獲得する仕組みの構築
    - 機械学習が苦手なところは人の力を借りる枠組みへ
  - 創造力をかきたてる技術を構築し, これまでには存在しなかった歌声／音声／楽器音データの作成を支援
    - 「人によるデータの創出」と「機械学習」のポジティブスパイラルへ

**人と機械学習が協力しあう仕組みが重要！  
2次創作文化を大いに活用！**

スペシャルセッション  
「音楽情報処理と機械学習」

続き：機械学習の効果と限界

歌声情報処理：声質変換モデル

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8月10日

# なぜ機械学習を使うのか？

- **こんなにも良いことがある！**
  - 解くべき問題を数理的に記述することで、**高い汎用性と実用性**が得られる！
  - **技術の共有**が容易であり、色々な問題へと適用できる！
  - **データドリブン**な枠組みであるため、システムを**自動的に構築**できる！
    - まだ発見されていない知見も見えてくるかも・・・
  - **統計処理**を用いることで、物理的に不可能なことも可能とする技術を生み出すことができる！
    - 声を真似る／混ぜる／創る／操るという処理を実現できる！

**人の能力 + 機械学習 = 世界が変わる！**

# 機械学習を適用する際に重要なことは？

- **適用する問題に対して特化させよう！**

- 応用先の処理(例えば合成)の特徴, モデル化対象となる信号の特徴, 物理的な生成過程などを考慮して, モデルを構築する.
  - 何を誤差(ゆらぎ)とみなしてモデル化するのかを考える.
- (機械学習のみに限らず)従来研究の知見を活用する.

- **解決すべき課題は多い！**

- **感性**: 主観的・知覚的な評価との対応を上手にとるには？
- **データ**: 如何に対象となるデータを大量に手に入れるか？
  - 未知のデータを創出するのは苦手(想像力に乏しい)・・・

✓ **人と機械学習が協力しあう仕組みが重要！**

「人によるデータの創出」と「機械学習」のポジティブスパイラルへ

✓ **二次創作文化を大いに活用！**