

大規模英語学習者を対象とした 音声の構造的表象に基づく発音分類とその応用

鎌田 圭[†] 高澤 真章^{††} 竹内 京子^{†††} 朝川 智[†] 峯松 信明[†] 牧野 武彦[†] 広瀬 啓吉^{††}
[†] 東京大学大学院新領域創成科学研究科 ^{††} 東京大学工学部情報工学科 ^{†††} 東京大学大学院総合文化研究科

[‡] 中央大学経済学部 ^{‡‡} 東京大学大学院情報理工学研究所

1 はじめに

音声の音響的特徴は、話者や音響機器などの非言語的要因によって不可避免的に変動し、同じ発話内容でも観測される音響事象は異なる。従来の音声システムでは、基本的に、大量の音声を用いて個々の言語事象を統計的にモデル化することで対処してきた。これは、不可避な音声歪みを内包した特徴量を利用して行うことに相当し、不一致問題の解決には至らず、話者を選ぶシステムとなっていた。音声認識の分野では、種々の適応技術（例えば話者適応）が用いられているが、本稿で検討するような発音学習に適用すると、下手な発音に対して高いスコアを与えるような適応が免れない。これは、話者の違いと発音の上手下手とがスペクトル包絡という同一物理現象に基づくため、両者を切り分ける術を持っていなかったことが根本原因である。

近年、非言語的要因に起因する静的な音響歪みを完全に排除して音声を表象する「音声の構造的表象」が提案された [1]。音声事象そのものではなく、事象間の相対的な関係のみを抽出/表象する。即ち、発話に含まれる有限個の事象から、全ての事象間差異（距離）を計算する。この際、一対一対応を満たす線形/非線形変換に対する変換不変量となる距離尺度 [2] を採択することで、非言語的変換に対して完全不変の音声表象となる。この音声表象は構造音韻論の数学的/物理的実装に相当し、外国語学習者の発音状態を非言語的要因の影響を受けずに記述することが可能である [3]。

提案表象を用いると、例えば子供の英語発音（細い声）と大人の英語発音（太い声）とを直接比較し、子供が母語話者であった場合、大人（学習者）の発音のどこを矯正すべきか、との教示を得ることができる [4]。また、複数の教師の中から、自分のお好みの教師を選ぶことが可能となる。図1は提案表象を用いて作成したデモシステム・インタフェースである。例えば、憧れの映画スターの音声データを使えば、その憧れに近づく最も効率的な発音矯正法を学習者別に提示することが可能である。更に本表象を用いることで、非言語的要因に影響を受けずに、学習者群の発音（方言）分類が可能となる。地球上には20億の英語学習者がいると言われるが、学習過程における彼らの時間的変化を考慮すると、 $20 \text{億} \times N$ 個となる発音を分類し、英語発音の世界地図を作ることも不可能ではない [5]。このようなインタフェースや試みは「複数話者の音声混ぜ合わせ」を基本とする従来法では不可能であり、本表象によって初めて可能となるアプリケーションである。

本稿ではこの発音分類について、模擬英語学習者96名に対して行った発音分類の実効性検証、及び、3年間に渡るオープンキャンパスを通して収集した柏市民564名の英語発音を分類することにより、柏に存在する英語方言の実験的定義を試みたので報告する。

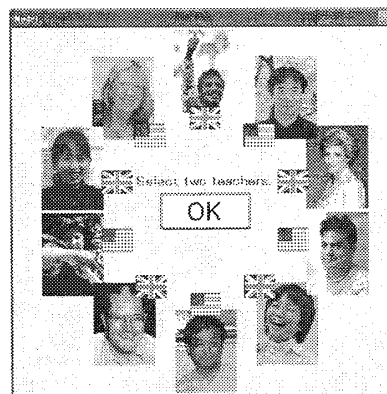


図1: Select your favorite teachers!

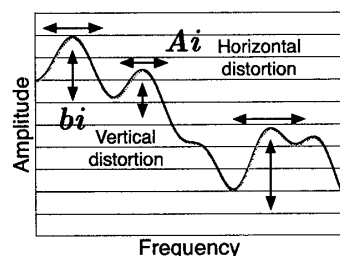


図2: スペクトルに対する水平/垂直方向の音響歪み

2 音声（発音）の構造的表象

静的な非言語的音響歪みは、乗算性歪みと線形変換性歪みに大別される。乗算性歪みは、マイクなどの伝送特性や話者の声道形状の違いの一部に相当し、ケプストラムベクトル c に対するベクトル b の加算 $c' = c + b$ となる。線形変換性歪みは、声道長の差異や聴覚特性の差異が挙げられ、行列 A の乗算 $c' = Ac$ で精度良くモデル化される。以上より、不可避かつ静的な非言語的音響歪みは、アフィン変換でモデル化できる。

図2はこの2種類の歪み（乗算性歪み=垂直、線形変換性歪み=水平）を示している。個々の音声事象を分布として捉え、全ての二事象間距離を計算すると（距離行列）、それは、事象群が成す幾何学構造を規定することになる。バタチャリヤ距離はアフィン変換のみならず、非線形変換も含め不変性を有するため、これを用いることで、静的な非言語的音響歪みに一切不変な音声表象が得られる。この構造を音響的普遍構造と呼んでいる。

3 日本人英語学習者の模擬音声を用いた発音分類

3.1 日本人英語学習者の模擬音声の作成

日本語・米語（英語）双方が話せる日本人話者12名（男性6名、女性6名）の米語11単母音と日本語5母音を収録した。収録では米語/bVt/、日本語/bVto/を5回発声させ、母音部分のみを切り出して実験に用い

