

音声認識のための周辺特徴の検討

福田 隆 瀧川 正史 新田 恒雄

豊橋技術科学大学大学院工学研究科 知識情報工学専攻
〒441-8580 愛知県豊橋市天伯町雲雀ヶ丘 1-1 E-mail : nitta@utkie.tut.ac.jp

あらし 本報告では周波数領域およびケプストラム領域において、周辺特徴の抽出機構を特徴抽出器へ組み込むことを試みる。最初に、周波数領域における周辺特徴抽出法を説明し、周辺特徴を組み込んだ特徴抽出器が従来の特徴抽出器に比べて、高い認識性能を示すことを **HMM** に基づく単語認識実験から述べる。次に、これらの実験結果からケプストラム領域の周辺特徴抽出法を提案し、“**MFCC** と周辺特徴”のパラメータセットを“**MFCC** と動的特徴”のパラメータセットと性能比較した結果から、周辺特徴として Δ_1 に Δ_q パラメータを併用することの有効性を指摘する。

キーワード 音声認識, 特徴抽出, ケプストラム分析, 直交基底, 写像演算子, 周辺特徴

Peripheral Features for Speech Recognition

Takashi FUKUDA, Masashi TAKIGAWA, and Tsuneo NITTA

Department of Knowledge-based Information Engineering, Graduate School of Engineering,
Toyohashi University of Technology, 1-1 Hibiriga-oka, Tempaku, Toyohashi, 441-8580 JAPAN

E-mail : nitta@utkie.tut.ac.jp

Abstract This paper describes an attempt to incorporate the mechanism of extracting the peripheral features into a feature extractor in frequency and quefrequency domain. Firstly, the method of extracting peripheral features in frequency domain is shown, and then the feature extractor that combines the peripheral features with **MFCC** improves the performance in comparison with the standard feature extractor in word recognition experiments with an **HMM**-based ASR system. Secondary, the method of extracting the peripheral feature in quefrequency domain is also provided. The proposed parameter set of **MFCC** with peripheral features shows significant improvements in comparison with the standard parameter set of **MFCC** with dynamic features in experiments.

key words Speech Recognition, Feature Extraction, Cepstrum Analysis, Orthogonal Basis, Mapping Operator, Peripheral Feature

1. はじめに

時間-スペクトル(TS)パターン $x(t, f)$ は長らく、自動音声認識における音響特徴として使用されてきた。近年、 Δ ケプストラム、 Δ パワーなどの動的特徴が ASR に導入され[1, 2]、MFCC と動的特徴のセットが広く用いられている。動的特徴は時間軸に沿った TS パターン上の点 $x(t, f)$ の周辺特徴を表現している。しかし、TS パターン上の $n \times n$ 近傍の点からはより多くの周辺特徴を得ることができると考えられる。

先行研究では[3]、複合音響特徴平面(MAFP)に基づいた特徴抽出を音素分類に適用し、性能を大きく改善することを示した。この方法では、 $x(t, f)$ を要素とする集合 X を局所的な空間写像演算子 $G_m (G_m \in G)$ により、さまざまな AFP (音響特徴平面 : Acoustic Feature Planes) $Y_m = y_m(t, f)$, $m=1, 2, \dots, M$ に写像した。

$$G_m : X \rightarrow Y_m \quad (1)$$

G_m としては、まず TS パターン上の 3×3 近傍の直交基底 $\{\Phi_1, \Phi_2, \dots, \Phi_9\}$ を音声データから直接抽出する。次に、直交基底を対称化して、モデル化したものを使用している。今回は、この手法を周波数領域およびケプストラム領域における周辺特徴抽出に適用することを試みる[4]。周波数領域においては、最初に直交基底を抽出し、それを写像演算子として捉え、TS パターンを周辺特徴空間に写像する。一方、ケプストラム領域においても同様の直交基底を抽出できる。抽出した直交基底、すなわち写像演算子は単純化かつ対称化した後、HMM に基づく ASR システムの特徴抽出器に組み込む。最後に、今回提案する特徴パラメータセットを不特定話者単語認識実験を通して、MFCC と動的特徴の標準的なパラメータセットと比較する。

以下、2. で TS パターン上の 7×7 近傍の幾何学的な特徴(構造特徴)について述べ、TS パターンから得られる周辺特徴を用いた認識実験結果を示す。3. では TS パターンに代えて、TQ パターン上の 7×3 近傍の直交基底を観測し、TQ パターンから抽出される周辺特徴を用いた認識実験結果を述べる。

2. 周波数領域における周辺特徴

2.1 TS パターン上の周辺特徴

TS パターン上にはさまざまな幾何学的構造が観測される。図 1 に TS パターン $x(t, f)$, $j=1, 2, \dots, 24$ の 7×7 ブロックにおける上位 9 個、 $\Phi_1, \Phi_2, \dots, \Phi_9$ の直交基底を示す。図では、黒と白の正方形はそれぞれ正と負の値を表し、正方形の大きさは振幅を表す。 7×7 直交基底は 2.3.1 で示す音声データから KL 変換(KLT : Karhunen-Loeve transform)を用いて抽出した。

これらの直交基底は空間写像演算子として観察すると、まず Φ_1 は平滑化演算子と見なせる。この演算子は MFCC の濃度情報を別に使用するなら、特徴抽出効果はないと考えられる。 Φ_2, Φ_3 はそれぞれ時間軸(Δ_t)と周波数軸(Δ_f)に関しての 1 次微分演算子、 Φ_4, Φ_5 は時間軸($\Delta_t \Delta_t$)と周波数軸($\Delta_f \Delta_f$)についての 2 次微分演算子、また Φ_6, Φ_7, Φ_8 は TS パターン上の尾根や谷を表す部分空間である。

時間-周波数空間演算子 $\{\Phi_m\} (\Phi_m \in \Phi)$ を用いると、TS パターン $x(t, f)$ からさまざまな周辺特徴 $Y_m = y_m(t, f)$, $m=1, 2, \dots, M$ を抽出することができる。周辺特徴の構成要素 $y_m(t, f)$ は 7×7 近傍の $x(t, f)$ と $\Phi_m = \phi_m(t, f)$ により次式から計算される。

$$y_m(t, f) = \sum_{i=-3}^3 \sum_{j=-3}^3 x(t+i, f+j) \phi_m(i, j) \quad (2)$$

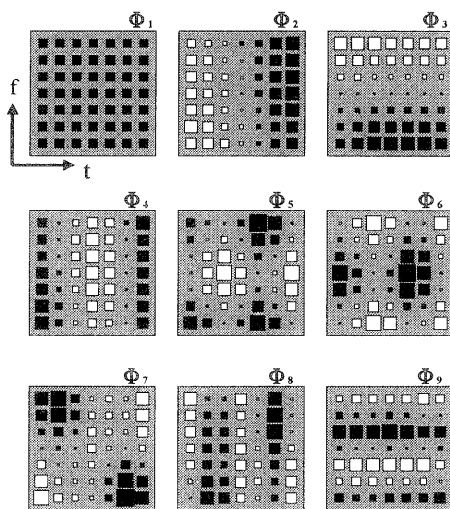


図 1 TS パターン上の 7×7 直交基底

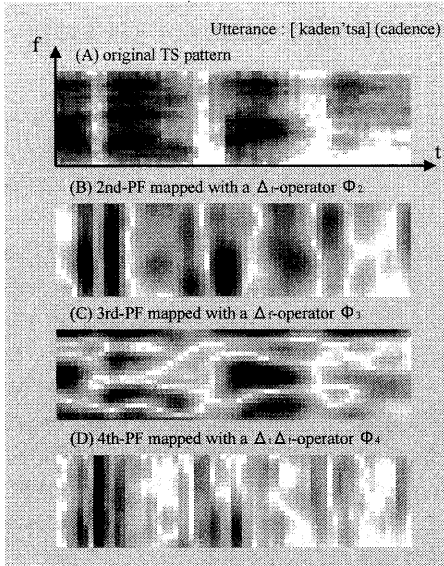


図2 周波数領域の時間 - スペクトルパターンと周辺特徴

図2は発話[kaden'tsa](カデンツァ)の上位3つの周辺特徴の例である。図において、(A)はTSパターン、(B)、(C)および(D)はそれぞれ Δ_t 演算子 Φ_2 を用いて写像した2番目のPF(PF: Peripheral Feature), Δ_t 演算子 Φ_3 による3番目のPF, そして $\Delta_t \Delta_t$ 演算子 Φ_4 による4番目のPFである。 $y_m(t, f)$ の正の値は正の傾斜を表しており、負の値は負の傾斜を示す。たとえば、破裂音の明瞭な時間変化は2番目のPF上で正負の値の対で表現され、同様に、定常部の明瞭なスペクトルピークは3番目のPF上で正負の値の対で表現される。なお、図では絶対値の大きさが示されている。

2.2 特徴抽出器への組み込み

この節では周辺特徴の抽出方法、および周辺特徴をMFCCと結合する方法を述べる。図3-Aは連続HMMに基づくASRシステムで標準的に用いられている特徴パラメータの抽出方法を示している。特徴抽出器では、入力音声を16KHzでサンプリングし、窓掛けした音声(25msハミング窓)を10ms毎に512点FFTした後、パワースペクトルをメルスケール化した24チャンネルのBPFを通して得る。その後、BPF出力を12次元のケプストラム係数(MFCC)にDCT変換した後、

対数パワー($\Delta P, \Delta \Delta P$)、および24次元の動的特徴($\Delta_t, \Delta_t \Delta_t$)を含む38次元の特徴パラメータを抽出する。

図3-BはTSパターン上のさまざまな周辺特徴を抽出し、それらを統合、圧縮する手順を示す。まず、24チャンネルのBPF出力 $x(t, f)$ から式(2)により周辺特徴 $y_m(t, f)$, $m=1, 2, \dots, 8$; $f=1, 2, \dots, 20$ を抽出する。続いて、それぞれの周辺特徴をDCTによりケプストラム $c(m, q)$, $m=1, 2, \dots, 8$; $q=1, 2, \dots, 10$ に変換する。最後に80次元の $c(m, q)$ を次式を用いてKL変換し、周辺特徴ベクトル $z(k)$, $k=1, 2, \dots, 24$ に圧縮する。

$$z(k) = \sum_{m=1}^8 \sum_{q=1}^{10} c(m, q) \phi_k(m, q) \quad k=1, 2, \dots, 24 \quad (3)$$

ここで、 $\phi_k(m, q)$ はKLTのためのk番目の固有ベクトルである。24次元の周辺特徴 $z(k)$ を、静的特徴を主として表すMFCC(12MFCC+ $\Delta p, \Delta \Delta p$)と結合する(計38次元)。

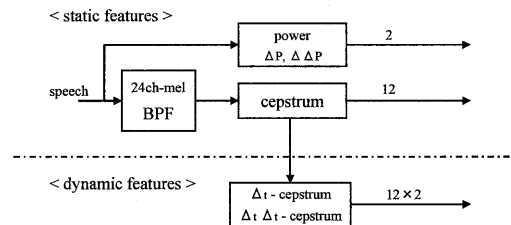


図3-A MFCCと動的特徴: Baseline

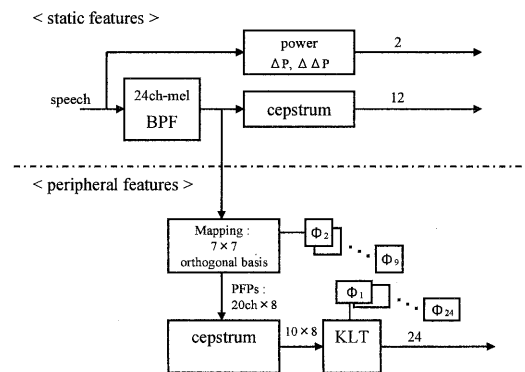


図3-B MFCCと周辺特徴(TS): PF-KL(7×7直交基底)

2.3 評価実験

2.3.1 音声試料

以下に示す4種類のデータセットを使用した。

- D1.** 音響モデル学習データセット：日本音響学会 (ASJ) 研究用連続音声データベース、(16KHz, 16bit)の男性話者30名(4503文)
- D2.** 評価データセット：東北大-松下単語音声データベース。男性話者10名(981単語)。サンプリング周波数は24KHzから16KHzへ変換。
- D3.** 7×7直交基底抽出用データセット：日本音響学会、新聞記事読み上げ音声コーパス (ASJ-JNUS, 16KHz, 16bit)の男性話者53名(2662文)
- D4.** KL変換用固有ベクトル学習データセット：直交基底抽出に用いたデータASJ-JNUSの内、前出とは異なる男性話者53名(5569文)

2.3.2 実験仕様

5状態3ループ、43音素の日本語 monophone-HMMをD1データセットを使用して設計した。HMMは出力確率をガウス混合分布で表現し、共分散行列を対角化している(混合数=1, 2, 4)。学習の後、D2のデータセットを用いて不特定話者単語認識実験を行った。

2.3.3 実験仕様

表1はBaselineの特徴抽出器と周辺特徴を組み込んだ特徴抽出器(PF-KL)の単語認識率である。周辺特徴を用いたPF-KLの認識率は、動的特徴のみを組み込んでいるBaselineの特徴抽出器よりも高い。特に、PF-KLの混合分布1のモデルの認識率が、同じ次元数のBaselineモデルの認識率よりきわめて高いことが注目される。この結果は、周辺特徴を付加したMFCCがロバスタな特徴パラメータセットを構成していることを推測させる。

表1 特徴抽出器の比較結果

method	word correct rate [%]		
	mix. = 1	mix. = 2	mix. = 4
Baseline	93.58	93.99	95.82
PF-KL	97.25	97.15	97.15

3. ケプストラム領域における周辺特徴

3.1 TQパターン上の周辺特徴

2節では、さまざまな幾何学的構造がTSパターン上で観測されることを示した。ここでは、TSパターンの代わりにTQパターン上の周辺特徴を観測する。図4はTQパターン $c(t, q)$, $j=1, 2, \dots, 12$ の7×3ブロックにおける上位9個、 $\Phi_1, \Phi_2, \dots, \Phi_9$ の直交基底を示している。直交基底の表現方法は図1と同様である。空間演算子の観点から、 Φ_1 は平滑化演算子と見なせる。 Φ_2, Φ_3 はそれぞれケフレンシー軸に沿った1次と2次の微分演算子(Δ_q 演算子, $\Delta_q \Delta_q$ 演算子)、 Φ_4, Φ_7, Φ_9 は時間軸に沿った1次, 2次および3次の微分演算子(Δ_t 演算子, $\Delta_t \Delta_t$ 演算子および $\Delta_t \Delta_t \Delta_t$ 演算子)、そして Φ_5, Φ_6 および Φ_8 はTQパターン上の谷や尾根を表す部分空間である。

TQパターン上の直交基底はTSパターン上の直交基底と異なり、ケフレンシー軸に沿った変化を示している基底(Φ_2, Φ_3)が時間軸に沿った変化を示している基底(Φ_4, Φ_7)よりも高い寄与を示した。

TQ空間演算子、すなわち写像演算子 Φ_m はTQパターン $c(t, q)$ を周辺特徴 $Y_m = y_m(t, q)$, $m = 1, 2, \dots, M$ に写像する。周辺特徴の構成要素 $y_m(t, q)$ は7×3近傍の $c(t, q)$ と $\Phi_m = \phi_m(t, q)$ により式(4)で計算される。

$$y_m(t, q) = \sum_{i=-3}^3 \sum_{j=-1}^1 c(t+i, q+j) \phi_m(i, j) \quad (4)$$

図5に発話[kaden'tsa](カデンツァ)の上位3つ

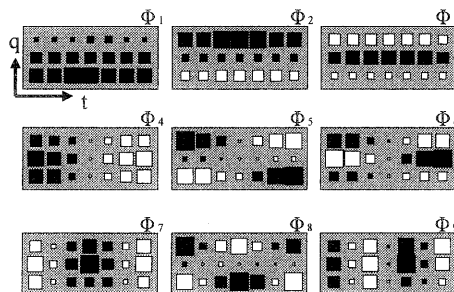


図4 TQパターン上の7×3直交基底

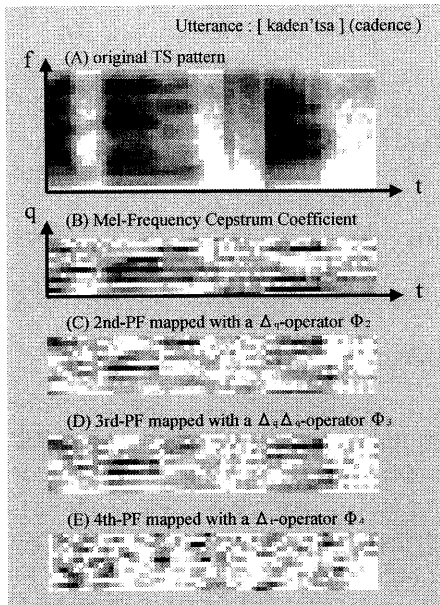


図5 ケプストラム領域の周辺特徴

の周辺特徴の例を示す。写像演算子 $\{\Phi_m\}$, $m=2, 3, 4$ は、単純化・対称化した後、TQ パターンに適用した。この場合、3つの写像演算子はそれぞれ Δ_q 演算子、 $\Delta_q \Delta_q$ 演算子および Δ_t 演算子に相当する。図において、(A)はTSパターン、(B)はMFCC、(C)、(D)および(E)はそれぞれ、 Δ_q 演算子 Φ_2 、 $\Delta_q \Delta_q$ 演算子 Φ_3 、 Δ_t 演算子 Φ_4 により抽出した周辺特徴である。周辺特徴と MFCC のパターンは絶対値で表示されている。図の(C)、(D)、(E)に示した周辺特徴には、写像演算子の性質を反映してケフレンシー方向、および時間方向の特徴的なパターンが現われている。

3.2 特徴抽出器への組み込み

3.1節では、TQパターン上の 7×3 直交基底を調べ、上位の主要な基底にはケフレンシー軸に沿う微分演算子 (Φ_2 , Φ_3) とみられる基底が存在することを見出した。一方、MFCCに基づく標準的な特徴抽出器はケフレンシー方向の変化を表す特徴 (Δ_q パラメータ) を含んでいない。図6はケプストラム領域の周辺特徴抽出過程を示す。図はMFCC および動的特徴の抽出過程を含んでいる。

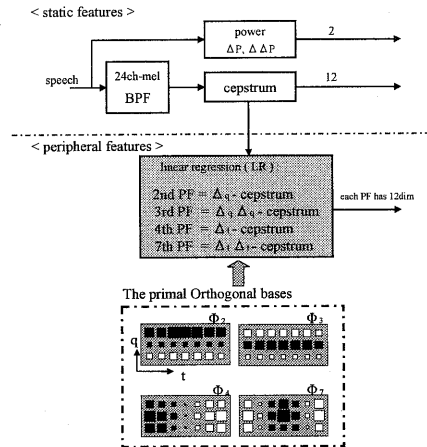


図6 MFCCと周辺特徴 (TQ)

図では、4つの周辺特徴 (Δ_q -, $\Delta_q \Delta_q$ -, Δ_t -, $\Delta_t \Delta_t$ -ケプストラム) を与える空間演算子を 1×3 ブロックの演算子 (Δ_q , $\Delta_q \Delta_q$) と 7×1 ブロックの演算子 (Δ_t , $\Delta_t \Delta_t$) の形式に単純化している。また、微分演算は線形回帰計算に置き換えている。周辺特徴は MFCC 特徴 ($12\text{MFCC} + \Delta P + \Delta \Delta P$) と合わせて、38次元の特徴パラメータセットとして使用される。

3.3 評価実験

3.3.1 実験仕様

実験に使用した音声試料は2.3.1節と同じである。また音響モデルも、2.3.2節に説明したと同じものを用いて、不特定話者単語認識実験を行う。特徴パラメータは Δ_t , $\Delta_t \Delta_t$, Δ_q および $\Delta_q \Delta_q$ を MFCC とさまざまな組み合わせで結合して評価する。

3.3.2 実験結果

表2および図7は種々の周辺特徴を MFCC に付加したときの実験結果である。表2の特徴パラメータセットの項では、MFCCの表記を省略し、付加した周辺特徴のみを記した。実験結果は以下のことを示している。

- 周辺特徴 Δ_q ケプストラムは Δ_t ケプストラムと同等の改善を与える。
- Δ_t ケプストラムと Δ_q ケプストラムの両方を付加したとき、最も高い認識率を得る。

なお、50次元の特徴である”MFCC+ Δt + Δq + $\Delta q \Delta q + \Delta P + \Delta \Delta P$ ”パラメータセットについても実験を行ったが、”MFCC+ Δt + Δq + $\Delta P + \Delta \Delta P$ ”のセットよりも認識率が低かった。これは学習データセットが小さいためと考えられ、今後、再検討が必要である。

ケフレンシー軸に沿った動的特徴量を組み込んだ特徴パラメータセットは、動的特徴のみを付加した特徴パラメータと比較して顕著な改善を示した。この結果は、ケフレンシー軸に沿った動的特徴量を標準的なMFCCパラメータセットへ付加することの重要性を示している。

表2 各特徴パラメータセットの認識率

特徴パラメータセット	word correct rate [%]		
	混合数		
	mix.=1	mix.=2	mix.=4
MFCC	81.96	81.86	81.04
Δt	91.54	92.25	92.86
Δq	93.68	94.29	93.58
$\Delta t, \Delta t \Delta t$	93.37	93.48	94.09
$\Delta q, \Delta q \Delta q$	93.17	94.19	94.19
$\Delta t, \Delta t \Delta t, \Delta P, \Delta \Delta P$	93.58	93.99	95.82
$\Delta q, \Delta q \Delta q, \Delta P, \Delta \Delta P$	95.21	95.21	95.31
$\Delta t, \Delta q, \Delta P, \Delta \Delta P$	97.96	98.37	98.47

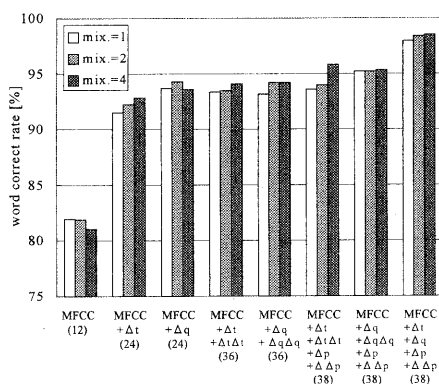


図7 MFCCパラメータセットの比較

提案した特徴抽出器の性能が標準的な特徴抽出器よりも高い理由としては、以下の点が考えられる。 Δt -パラメータは時間軸に沿った変動を表現するため子音の認識には有効に働くが、一方、定常的な母音等の認識には寄与しない。他方、 Δq -パラメータはケフレンシー方向の変動を表現するため、主に母音認識に大きく寄与する。こうした仮説は今後、実験を通して実証する必要がある。

4. まとめ

さまざまな幾何学的構造をASRシステムの特徴抽出器に組み込む方法を提案した。周辺特徴に関わる特徴抽出器の設計法としては、音声の直交基底を観測し、主要な基底を単純化・対称化した形で特徴抽出器に組み込む方法を示した。周辺特徴とMFCCを結合した特徴パラメータセットは、HMMに基づく単語認識実験において標準的な特徴セット(MFCC+動的特徴)と比較して顕著な改善を示した。

参考文献

- [1] K. Elenius and M. Blomberg, "Effect of emphasizing transitional or stationary parts of the speech signal in a discrete utterance recognition system", IEEE Proc. ICASSP'82, pp. 535-538 (1982).
- [2] S. Furui, "Speaker-independent isolated word recognition using dynamic features of speech spectrum", IEEE Trans. Acoust. Speech Signal Process. ASSP-34, pp. 522-59 (1986).
- [3] T. Nitta, "Feature extraction for speech recognition based on orthogonal acoustic-feature planes and LDA," Proc. IEEE ICASSP'99, Phoenix, Vol.1, pp.421-424 (1999-3).
- [4] T. Nitta, M. Takigawa, and T. Fukuda, "A Novel Feature Extraction Using Multiple Acoustic Feature Planes for HMM-based Speech Recognition" Proc. ICSLP'00, Vol. 1, pp. 385-388 (2000).