

汎用・識別的特徴量を用いた音声区間検出 Voice Activity Detection using GA-based Informative Feature

奥田 博也[†] 田村 哲嗣[‡] 速水 悟[‡]
Hiroya Okuda Satoshi Tamura Satoru Hayamizu

1. はじめに

近年音声認識技術は様々な場面で活躍している。例えば、カーナビゲーションシステムや携帯電話などで用いられている。しかし雑音環境下では、認識率が低下してしまう問題点があり、雑音に対し何らかの前処理を施す必要がある。その 1 つとして音声区間検出(Voice Activity Detection, VAD)が広く用いられている。VAD とは、入力された音声信号を音声区間と非音声区間に分類する手法である。検出された音声区間に対してのみ認識を行うことで、認識性能を改善できる。過去の研究では、MFCC[1]や 3 次キュムラント[2]を特徴量として VAD を行ない、雑音に対する頑健性を検証している。本稿では、VAD に汎用・識別的特徴量(GA-based Informative Feature, GIF)[3]を用いて、検出性能を検証し、考察する。

2. 汎用・識別的特徴量を用いた音声区間検出

2.1 クラスタ分析

後述する特徴量抽出では変換行列をクラスを基に学習で求める必要がある。このため、入力ベクトルについてクラスタ分析を行う。本来、VAD では識別するクラスは音声と非音声の 2 クラスだが、GIF への変換ではクラス数を拡張させる必要があるため、音声 3、非音声 2 の計 5 クラスとした。

2.2 汎用・識別的特徴量

汎用・識別的特徴量(GIF)は田村らによって提案された特徴量である[3]。入力ベクトルに対し二段階で変換を施すことで、最終的な識別用の特徴量を計算する。はじめに、それぞれのクラスごとに入力ベクトルが該当クラスに属するか否かを判別する二値写像を形成する。そしてこれらを統合することで第一の変換を作成する。次に、入力ベクトルに第一の変換を行って得られた中間ベクトルにおいて、クラスごとに平均を求め、分散が最大になるよう、第二の変換を作成する。同時に、特徴量次元間の無相関化と特徴量次元の削減を行う。これらの変換を得るために、遺伝的アルゴリズム(Genetic Algorithm, GA)を用いる。

N 次元の入力特徴量空間を \mathbf{X} 、入力ベクトルを $\mathbf{x} \in \mathbf{X}$ 、出力される M 次元の特徴量空間を \mathbf{Z} 、 \mathbf{x} から得られる出力ベクトルを $\mathbf{z} \in \mathbf{Z}$ とする。本稿では、 \mathbf{x} から \mathbf{z} を得る際に以下のように二段階の変換(線形変換)を行う。図 1 に概略を示す。

第一の変換では \mathbf{x} から C 次元の中間ベクトル \mathbf{y} を(1)式により計算する(C は識別すべきクラスの数)。

$$\mathbf{y} = \mathbf{A} \times (\mathbf{x}^T - \mathbf{1})^T \quad (1)$$

\mathbf{A} は $C \times (N + 1)$ 次元の変換行列である。次に第二の変換では、(2)式のように中間ベクトル \mathbf{y} から次式により識別特徴量 \mathbf{z} を計算する。

$$\mathbf{z} = \mathbf{B} \times \mathbf{y} \quad (2)$$

\mathbf{B} は $M \times C$ 次元の行列である($1 \leq M \leq C$)。本稿では、GIF への入力ベクトルに、フィルタバンク 24 次元を用いた。GIF を用いた VAD の概略を図 2 に示す。

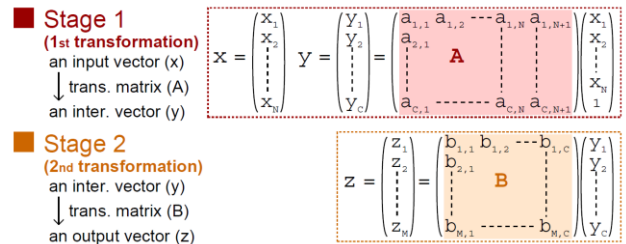


図 1 GIF の概略

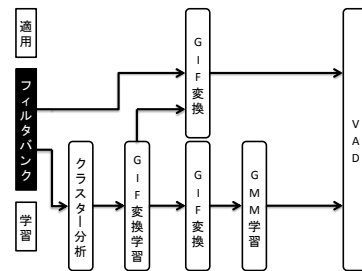


図 2 GIF を用いた VAD の概要

2.3 GMM を用いた VAD アルゴリズム

本稿では、GMM(Gaussian Mixture Model)を用いた VAD を実装した。はじめに、学習データから得られた GIF を用いてモデルの学習を行った。VAD を行う際にはフレームごとに対数尤度 $L(i)$ を計算し、閾値以上なら音声区間、閾値未満なら非音声区間とした。音声区間の尤度 $L_v(i)$ 、非音声区間の尤度 $L_o(i)$ は(3)、(4)式のように GMM を用いて計算する。 i はフレーム番号、 \mathbf{o}_i は i フレーム目で観測された特徴ベクトル、 $\boldsymbol{\pi}$ は混合重み、 K は混合数、 $\mathbf{x} = \mathbf{v}$ または \mathbf{o} である。

$$L_x(i) = \sum_{k=1}^K \pi_{k,x} N(\mathbf{o}_i | \boldsymbol{\mu}_{k,x}, \boldsymbol{\Sigma}_{k,x}) \quad (3)$$

$$L(i) = \log L_v(i) - \log L_o(i) \quad (4)$$

ただし、音声区間は 100ms 以上、非音声は 500ms 以上の連続した区間があるものとする。

CENSREC-1-C[4]のベースラインと同様に、各データについて初期閾値 THR_{int} を判別分析法[5]で決定する。その後、初期閾値を用いて(5)、(6)式のように閾値 THR を変動させる。

$$THR = THR_{int} + k \cdot \alpha \quad (5)$$

$$\alpha = \frac{L_h - L_l}{40} \quad (6)$$

[†] 岐阜大学大学院工学研究科

[‡] 岐阜大学工学部

ここで、 L_h は閾値以上の尤度の平均、 L_l は閾値未満の尤度の平均である。 k を変動させることで閾値を自動的に決定できる。本稿では $k = -40, -39, \dots, 39, 40$ とし、81通りの閾値についてVADを行った。

3. 評価実験

3.1 実験条件

CENSREC-1-C[4]を用いてVADを行った。CENSREC-1-CはVADの評価基盤であり、収録されている音声は男性話者、女性話者それぞれ52名の合計104名による連続数字読み上げである。クリーン音声と人工的に作成したシミュレーション雑音、実環境データが含まれている。本稿ではSubway, Babble, Car, Exhibition, Restaurant, Street, Airport, Stationのシミュレーション雑音を用い、SNRはそれぞれ20dB, 10dB, 5dB, 0dB, -5dBとした。

GIFとGMMの学習には、CENSREC-1-Cの男性26名、女性26名の計52名分のクリーン音声を用いた。識別は学習と異なる話者52名に対して行った。GIFの出力ベクトルは4次元で、GMMの混合数は音声、非音声それぞれ32である。VADの評価は10ミリ秒ごとにフレーム単位で行い、評価尺度は(7)、(8)式のFRR(False Rejection Rate)とFAR(False Acceptance Rate)を用いた。雑音別の評価はEER(Equal Error Rate)を用いた。EERはFRRとFARが等しくなる時の結果である。

$$FRR = \frac{N_{FR}}{N_s} \times 100 \quad (7)$$

$$FAR = \frac{N_{FA}}{N_{ns}} \times 100 \quad (8)$$

N_s は音声フレームの総数、 N_{FR} は非音声と検出された音声フレームの数、 N_{ns} は非音声フレームの総数、 N_{FA} は音声と検出された非音声フレームの数である。

3.2 実験結果

クリーン環境のGIFの特徴量を図3に示す。図3をみると得られた特徴量は音声区間で1次元目の値が高くなっており、VADに有効であると考えられる。GIFの尤度によるVADは対数パワーのみを用いたベースラインと同程度の識別性能であった。またGIF1次元目を用いたVADの追加実験として2.3節のように閾値を設けてVADを行ったところ、ほぼ同等の結果が得られた。1次元目を用いて閾値処理したときのSNR別の結果を図4に示す。また、表1は雑音別の識別結果である。数値は $100 - EER$ を表し、大きいほど性能が良い。Babble, Exhibition, Restaurant, Airport雑音下では、識別率の改善が見られた。GIFは比較的非常なこれらの雑音下においてVAD性能の向上が見られた。GIFは音声の周波数的な特徴を表現できていると考えられるため、単純なパワーよりも良くなったことが理由として挙げられる。

4. まとめ

本稿ではGIFを用いたVADを提案した。GIFの尤度によるVADとGIFの1次元目によるVADの性能を評価した。SNR別ではベースラインと同程度であるが、雑音の種類によっては性能の改善が見られた。今後は実環境での検出性能の検証、画像特徴量と統合したマルチモーダル音声区間検出への応用が挙げられる。

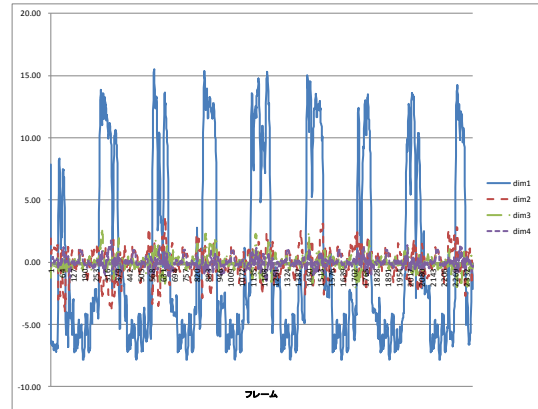


図3 GIF変換結果

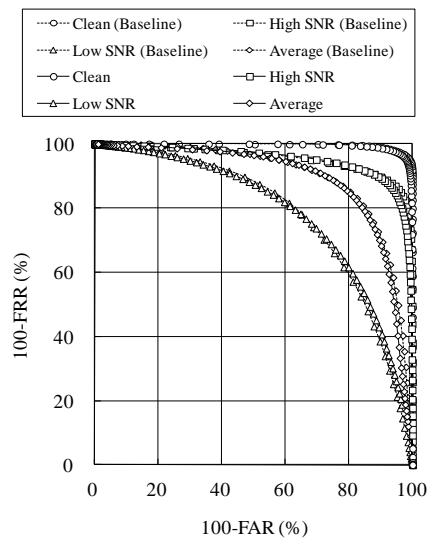


図4 SNR別VAD結果

表1 雑音別の識別結果(100-EER)[%]

	Subway	Babble	Car	Exhibition
Baseline	82.10	81.58	85.50	85.2
GIF	81.23	82.35	84.35	87.23
	Restaurant	Street	Airport	Station
Baseline	81.30	83.05	82.35	81.83
GIF	82.93	78.83	84.10	79.28

参考文献

- [1] 和田直哉, 早坂昇, 宮永善一, 畑岡信夫, "メルケプストラムを用いたロバスト音声区間検出", 電子情報通信学会技術研究報告, DSP, デジタル信号処理, Vol.103, No.146, pp.25-30, 2003.
- [2] 松田博義, 滝口哲也, 有木康夫, "3次元キュラント音声特徴を用いた音声区間検出", 電子情報通信学会技術研究報告, SP, 音声, Vol.106, No.263, pp.37-42, 2006.
- [3] 田村哲嗣, 田上陽嗣, 速水悟, "GIF-SP:汎用・識別的な特徴量を用いた音声認識性能の改善", 電子情報通信学会技術研究報告, SP2011-92, Vol.111, No.364, pp.119-124, 2008.
- [4] 北岡教秀, 山田武志, 柘植寛, 宮島千代美, 西浦敬信, 中山雅人, 傳田遊亀, 藤本雅清, 山本一公, 滝口哲也, 黒岩真吾, 武田一哉, 中村哲, "CENSREC-1-C:雑音下音声区間検出評価基盤の構築", 情報処理学会研究報告, 2006-SLP-63, Vol.2006, No.107, pp.1-6, 2006.
- [5] 大津展之, "判別および最小2乗基準に基づく自動しきい値選定法", 電子通信学会論文誌, Vol.J63-D, No.4, pp.349-356, 1980.