

# NMF と音響モデル併用型 DNN に基づく音声区間検出

太刀岡 勇気<sup>1,a)</sup>

**概要:** 音声区間検出を行う際には、パワーに基づく方法がよく使われる。しかしながらこの方法は高騒音下において性能の低下が著しいため、近年ではスペクトルの形状を考慮するような方法が提案されており、とりわけ深層神経回路網に基づく方法が性能が良いことが知られている。本報では、この方法の更なる改善を目的として、発話者の特徴や発話内容に応じた補助特徴量を用いる方法を提案する。特徴量として、非負値行列因子分解の活性化と音素ごとの事後確率を採用し、これらの有効性を車内環境での評価実験により確認した。

**キーワード:** DNN に基づく音声区間検出, 非負値行列因子分解, 音声認識, 補助特徴量

## DNN based voice activity detection with joint use of NMF and acoustic model

TACHIOKA YUUKI<sup>1,a)</sup>

**Abstract:** For voice activity detection, power-based methods are widely used; however, because these methods are susceptible to noise, recently, methods that consider the shape of spectrum have been proposed. In particular, deep neural network based methods have outperformed other methods. This paper aims to improve these methods by using auxiliary features that correspond to the speaker characteristics and the contents of the utterances. This paper proposes to use activation of non-negative matrix factorization and posterior probabilities of phonemes as an auxiliary feature and validates the effectiveness on the experiments in in-car environments.

**Keywords:** DNN-based voice activity detection, non-negative matrix factorization, automatic speech recognition, auxiliary features

### 1. はじめに

遠隔マイクでの音声操作や発話スイッチレス音声認識など、実環境で音声インターフェースを利用する機会が増加している。これに伴い、騒音環境下で対象話者が発話した区間を検出する音声区間検出技術が、重要性になってきている。以前は音声のパワーが騒音のそれよりも大きいことに依拠したパワーに基づく方法が使われていたが、高騒音下では性能が低下するという問題があった。この問題を解決するため、各周波数ビンでスペクトルのモデル化を行い、スペクトルの形状をとらえる尤度比検定に基づく方

法 [1], [2] が主流となった。また、事前に学習した音声モデルを用いる方法 [3] の有効性も知られている。近年では、Deep Neural Network (DNN) を使うことで、さらに性能が向上することが示されている [4]。

一方、DNN に基づく音声強調では、音声認識の情報などを補助的な特徴量として使う方法が提案されている [5], [6]。本報では、DNN 音声区間検出において、スペクトル特徴量に加えて、補助音声モデルより出力される補助特徴量を併用することで、音声区間検出の精度を向上させる方法を提案する。補助音声モデルとしては、非負値行列因子分解 (NMF) と音声認識用の音響モデルを利用する。前者では、NMF のアクティベーションを、後者では、音響モデルの音響スコアを補助特徴量として用いる。さらにこれらの特

<sup>1</sup> デンソーアイティラボラトリ  
東京都渋谷区渋谷 2-15-1 渋谷クロスタワー 28F 150-0002  
<sup>a)</sup> ytachioka@d-itlab.co.jp

微量を併用した際の性能も調査する。CENSREC-2 を用いた車内発話の音声を用いた実験により、提案法の有効性を示す。

## 2. Sohn の方法

ここでは DNN 以前の従来法として一般的な、周波数ごとのスペクトル特徴を利用して音声区間検出を行う Sohn の方法 [1] を概観する。短時間フーリエ変換により、観測音の FFT 係数  $\mathbf{X} \in \mathbb{C}^{F \times T}$  の時刻  $t$  における特徴量  $\mathbf{X}_t = \{X_{f=1, \dots, F}\} \in \mathbb{C}^F$  を求める。非音声区間  $H_N$  と音声区間  $H_S$  での音声と騒音の FFT 係数をそれぞれ  $\mathbf{S}_t = \{S_{f=1, \dots, F}\}$ ,  $\mathbf{N}_t = \{N_{f=1, \dots, F}\}$  とすると、観測音はそれぞれの区間で、

$$H_N: \mathbf{X}_t = \mathbf{N}_t, H_S: \mathbf{X}_t = \mathbf{N}_t + \mathbf{S}_t \quad (1)$$

のように表される。ここで  $H_N$ ,  $H_S$  において、それぞれの  $\mathbf{X}_t$  の確率密度関数が<sup>3</sup>、次式のように各次元で独立なガウス分布で表せると仮定する。

$$p(\mathbf{X}_t | H_N) = \prod_{f=1}^F \frac{1}{\pi \lambda_f^N} e^{-\frac{|X_f|^2}{\lambda_f^N}} \quad (2)$$

$$p(\mathbf{X}_t | H_S) = \prod_{f=1}^F \frac{1}{\pi [\lambda_f^N + \lambda_f^S]} e^{-\frac{|X_f|^2}{[\lambda_f^N + \lambda_f^S]}}$$

ここで  $\lambda_f^N$ ,  $\lambda_f^S$  は  $N_f$ ,  $S_f$  の分散を表す。すると  $f$  次元目の音声・非音声の尤度比は、式 (3) で表される。

$$\Lambda_f(X_f) = \frac{p(X_f | H_S)}{p(X_f | H_N)} = \frac{1}{1 + \xi_f} e^{\frac{\gamma_f \xi_f}{1 + \xi_f}} \quad (3)$$

$$\xi_f = \lambda_f^S / \lambda_f^N, \gamma_f = |X_f|^2 / \lambda_f^N$$

ここで  $\xi_f$ ,  $\gamma_f$  はそれぞれ事前、事後 SN 比と呼ばれる。それぞれの次元の尤度比の幾何平均により、音声・非音声を判断できる。

$$\log \Lambda(\mathbf{X}_t) = \frac{1}{F} \sum_{f=1}^F \log(\Lambda_f(X_f)) \underset{H_N}{\overset{H_S}{\gtrless}} \eta \quad (4)$$

$\log \Lambda(\mathbf{X}_t)$  が閾値  $\eta$  よりも大きければ時刻  $t$  は  $H_S$ , 小さければ  $H_N$  となる。ここで  $\lambda_f^N$  は観測された騒音の分散を集めた騒音モデルであり、事前に推定しておく。音声モデル  $\lambda_f^S$  を最尤基準により推定すると、最終的に音声・非音声の判別式は式 (5) のようになる。

$$\log \Lambda^{(ML)}(\mathbf{X}_t) = \frac{1}{F} \sum_{f=1}^F (\gamma_f - \log \gamma_f - 1) \quad (5)$$

## 3. DNN に基づく音声区間検出法

音声認識で DNN の有効性が示されるのとはほぼ同時に、音声区間検出においても DNN の有効性が確認された [4]。  $\mathbf{X}$  より得られるスペクトル特徴量を入力として、例えば

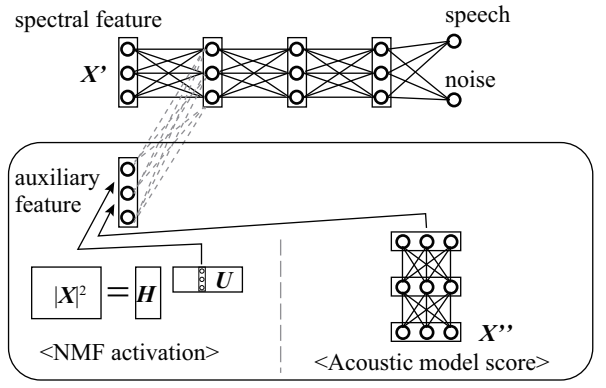


図 1 The proposed DNN-based voice activity detection (VAD) system using auxiliary speech models.

2 ノードの出力を設けて置き、学習データの音声/非音声の状態に対応して、一方のノードの出力が 1 になるように DNN を学習する。テスト時にはそれらの出力の softmax をとることで、音声の事後確率を算出することができる。式 (6) のように、スペクトル特徴量  $\mathbf{X}' \in \mathbb{R}^{F' \times T}$  を DNN に入力し (変換を  $f$  で表す)、出力値  $\mathbf{y} \in \mathbb{R}^{2 \times T}$  を得る。

$$\mathbf{y} = f(\mathbf{X}') \quad (6)$$

## 4. 補助特徴量の利用

音声強調の分野において、スペクトル特徴に加えて補助的な特徴量を用いることで、音声強調の性能が向上することが示されている [5], [6]。また音声認識の分野においても、DNN の音響モデルに対して話者性を表す特徴量を補助的に入力することで、音声認識の性能が向上することが知られている [7]。これらは補助特徴量を使うことで、DNN を環境に合わせて適応化しているととらえることができる。音声区間検出においても、このような手法が有効であると考えられる。音声区間検出を行う問題は、音声と騒音を区別する問題であるが、音声は多様性が大きくこの問題を直接解くことは難しい。そこで、補助特徴量として、音声のパターンを限定するような特徴量を用いれば、音声の多様性を縮小できる。例えば、音素を表す特徴量を用いれば、先の音声と騒音を区別する問題が、ある特定の音素と騒音を区別する問題に単純化でき、音声区間検出の性能が向上することが期待される。図 1 に提案のシステムを示す。NMF によるアクティベーション、もしくは音響モデルのスコアを補助特徴量として用いて、DNN により音声区間検出を行う。

### 4.1 NMF アクティベーション

騒音が混ざった音声  $\mathbf{X}$  のパワースペクトルを NMF によって、騒音と音声に分離する。

$$|\mathbf{X}|^2 \simeq \mathbf{H}\mathbf{U} = \mathbf{H}_s\mathbf{U}_s + \mathbf{H}_n\mathbf{U}_n \quad (7)$$

ここで  $\mathbf{H} \in \mathbb{R}_{\geq 0}^{F \times K}$  は  $K$  個の基底からなる基底行列、

$U \in \mathbb{R}_{\geq 0}^{K \times T}$  は、基底  $k$  の時刻  $t$  における活性化度  $U_{k,t}$  を表すアクティベーション行列である。基底を音声の基底  $H_s$  と騒音の基底  $H_n$  に分けると、それぞれのアクティベーションも  $U_s$  と  $U_n$  に分けられる。ここでは、この  $U$  もしくは  $U_s$  に着目する。NMF のアクティベーションは基底  $H$  が適切なものであるならば、発話に含まれる音声の特徴をよく表していると考えられる。実際 [8] では、音声の基底に対応するアクティベーションを利用して条件付き確率場で音声区間検出を行っている。そこで  $U$

$$y = f([X'; U]) \quad (8)$$

もしくは  $U_s$  を

$$y = f([X'; U_s]) \quad (9)$$

のように、補助特徴量として用いる。

#### 4.2 音声認識の音響モデルの音響スコア

文献 [5], [6] では、音声認識の結果をフィードバックすることで、音声強調の性能を向上させる方法が提案されている。ここにおいても、音声認識に用いる音響モデルにより、音素毎に属する確率を算出し、それを補助特徴量として用いる方法も考えられる。音響モデルによるスペクトル特徴量  $X'^1$  の音素毎の事後確率への変換を  $g$  とする\*2。音響モデルのスコアを正規化して\*3 DNN の入力とすると、式 (10) のようになる。

$$y = f\left(\left[X'; \frac{g(X')}{|g(X')|}\right]\right) \quad (10)$$

#### 4.3 結合型

結合型では上記の補助特徴量を結合した特徴量を使う。ただしそのままでは次元が高くなりすぎるので、必要に応じて主成分分析 (PCA) により次元削減行列  $P$  をかける。

$$y = f\left(\left[X'; P \left[U; \frac{g(X')}{|g(X')|}\right]\right]\right) \quad (11)$$

### 5. CENSREC-2 による実験

#### 5.1 実験条件

車内で実収録された音声データセットである CENSREC-2 を用いて\*4、音声区間検出のためのデータセットを構築した [9]。CENSREC-2 では発話毎にファイルが切り出されているが、これを連結して一人当たり 1 分程度の音声データを作成し、音声区間検出の実験を行った。一つの走行速度につき、話者数は学習セット 58 人、評価セット 15

\*1 式 (6) でのスペクトル特徴量と同じである必要はない。

\*2 GMM 音響モデル、DNN 音響モデルのいずれも使える。

\*3 必ずしも正規化する必要はないが、値のレンジの大きなスコアを扱う場合は何らかの配慮が必要となる。

\*4 音声区間の実験をおこなうためのデータセットとして、CENSREC-1C があるが、サンプリング周波数が 8kHz で実情と合わないため、独自のデータセットを構築した。

人である\*5。音声区間検出用の DNN の学習は、3 種の走行速度 (アイドリング (i.a), 低速 (市街地) 走行 (c.a), 高速走行 (e.a)) すべての音声を用いて行った。各走行速度において、4 種類の車内環境 (通常走行, エアコン On, オーディオ On, 窓開) をほぼ同じ割合で組み合わせた 12 種類の環境が存在するが、集計は走行速度別に行った。発話は数字 11 種類 (1~9, 0(まる), Z(ゼロ)) から構成される。CENSREC-2 には、音声区間の時間ラベルが含まれていないため、ラベル付けは接話マイク収録された音声を自動音声認識して行った。付属のスクリプトで「Condition 3」で音響モデルを適合学習し、そのモデルにより音声認識した際のアライメント結果の時間情報から、10ms 周期のフレーム単位で音声/非音声 2 値のラベル付けを行った。

表 1 に実験の設定を示す。音響特徴量は、音声区間検出用の DNN と音響スコア計算用の DNN には、0 次から 22 次のフィルターバンク (fbank) 特徴量を、前後 4 フレームコンテキスト拡張したものを用いた。NMF のアクティベーションはすべての基底に対する  $U$  (式 (8)) と音声のみの基底に対応する  $U_s$  (式 (9)) の 2 通りで実験した。また GMM による音響スコアの計算には、0 次から 12 次までの MFCC 特徴量とその動的特徴量を用いた。結合型の場合は、両補助特徴量を単純に結合したものと、NMF アクティベーションを PCA により 400 次元に圧縮したのちに特徴量結合したものを比較した。Sohn の方法では、平均的に最も良い音声区間検出精度が得られるときの閾値を、環境に共通で与えた。DNN の方法では、2 ノードの出力の softmax 値を取り、音声の事後確率が 0.5 を超えた場合に音声、それ以外は騒音として判定した。

図 2 に、高速走行時の近接マイクと遠隔マイクにより収録された音声の比較を示す。近接マイクにより収録された音声は音声区間を視察で与えることもできそうだが、遠隔マイクの方は全体が騒音に埋もれており、視察では音声区間を特定することは難しい。

表 1 Setup for the VAD system.

Sampling frequency	16 kHz
Window length	25 ms
Window shift	10 ms
Features	0-22th fbanks
Splice	9 frames
# NMF bases	50
# DNN output nodes	2
# DNN nodes per layer	1,000 nodes
DNN layer size	3 layers

\*5 CENSREC-2 の評価セットには接話マイクの音声がなく音声区間のラベリングが困難であったので、CENSREC-2 の学習セットを分割して、新たに学習セットと評価セットを構築した。

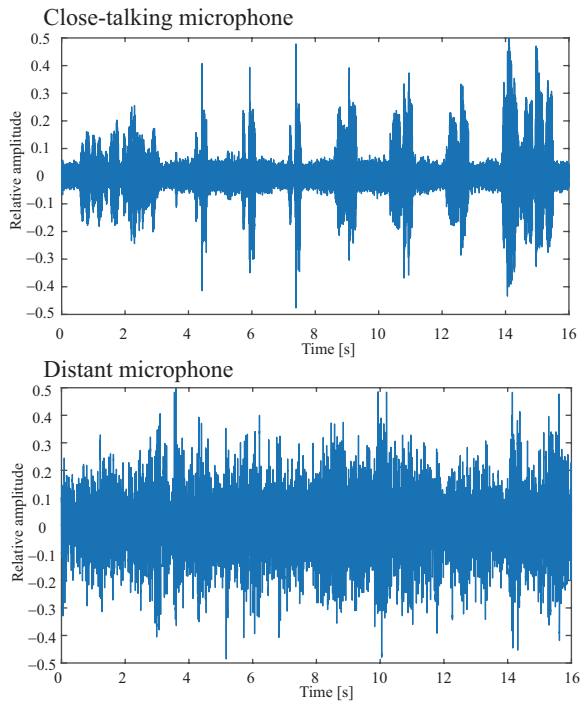


図 2 Waveforms recorded by cross-talking and distant microphones in a highly noisy environment.

## 5.2 ベースライン

表 2 には、フレーム単位での平均音声区間検出精度を示す。Sohn の方法のベースライン、MMSE-STSA により騒音抑圧した後に Sohn の方法を用いたもの、多層パーセプトロン (MLP) によるベースラインである。MLP(mel) がメルケプストラムを特徴量とした DNN に基づく手法のベースラインであるが、Sohn の方法に比べて非常に高い性能を示している。

## 5.3 NMF アクティベーション

表 2 には、NMF の音声の基底に対応するアクティベーション  $U_s$  のみ (speech activation) と全基底に対するアクティベーション  $U$  (all activation) を、特徴量として MLP により音声区間検出を行った結果も示している。MLP(mel) よりも性能が若干低いものの、音声区間検出はできており、アクティベーションに音声区間検出に有用な情報が含まれていることが確認できた。メルケプストラムに加えて補助特徴量としてアクティベーションを用いることで、補助特徴量を用いないものに比べて、どちらの場合も性能が向上した。騒音の基底に対するアクティベーションを用いても、精度は向上しないことが分かった。このことから、音声の基底に対するアクティベーションの有効性が示された。

## 5.4 音響スコア

表 3 には、音響モデル (GMM/DNN) の音響スコアを補助特徴量とした場合の結果を示す。5.3 節に示した NMF アクティベーションを用いた結果に比べ全体的に精度が高

表 2 Average frame-level VAD accuracy [%]. The performance of MLP was compared with that of the conventional Sohn's method. MLP used filterbank features with NMF activations.

	e_a	c_a	i_a
Sohn	52.08	63.34	63.23
Sohn (w MMSE-STSA)	62.25	60.81	65.06
MLP (speech activation)	76.00	84.36	90.00
MLP (all activation)	77.48	85.75	90.89
MLP (mel)	77.76	86.03	91.62
+ speech activation	79.37	<b>87.92</b>	<b>92.70</b>
+ all activation	<b>79.38</b>	87.79	92.64

表 3 Average frame-level VAD accuracy [%]. MLP used filterbank features with clean speech acoustic model (GMM/DNN) outputs.

	e_a	c_a	i_a
MLP (mel)	77.76	86.03	91.62
+ speech GMM	80.23	88.33	92.94
+ speech DNN	<b>81.70</b>	<b>90.14</b>	<b>94.28</b>

表 4 Average frame-level VAD accuracy [%]. MLP used filterbank features with both auxiliary features.

activation	GMM/DNN	PCA	e_a	c_a	i_a
speech	GMM	no	80.18	88.87	93.42
all	GMM	no	79.97	88.24	93.19
speech	DNN	no	81.51	<b>90.08</b>	94.11
all	DNN	no	81.49	89.91	94.00
speech	GMM	yes	80.01	88.47	93.30
all	GMM	yes	79.75	87.64	92.53
speech	DNN	yes	<b>81.66</b>	89.88	<b>94.38</b>
all	DNN	yes	81.44	89.51	94.10

く、GMM よりも DNN 音響モデルを用いた場合の方が有効性が高いことが分かった。

## 5.5 結合型

表 4 には結合型の音声区間検出性能を示す。全体的に音声の基底に対応するアクティベーションのみを使ったものの方がよく、DNNの方がGMMよりも優れていることがわかる。表 3 の結果と比べて性能の向上はほとんど見られなかった。

## 5.6 尤度比、事後確率の比較

図 3 には、図 2 の音声を与えたときの、式 (5) で計算される Sohn の方法の対数尤度比  $\log \Lambda$  と、提案の DNN の音声区間検出モデルに DNN の音響モデルの音響スコアを補助特徴量として与えた際に算出された音声の事後確率を示す。どちらの方法も発話を取り逃してはいないものの、Sohn の方法の結果が非常に変動が大きく、始末端検出の

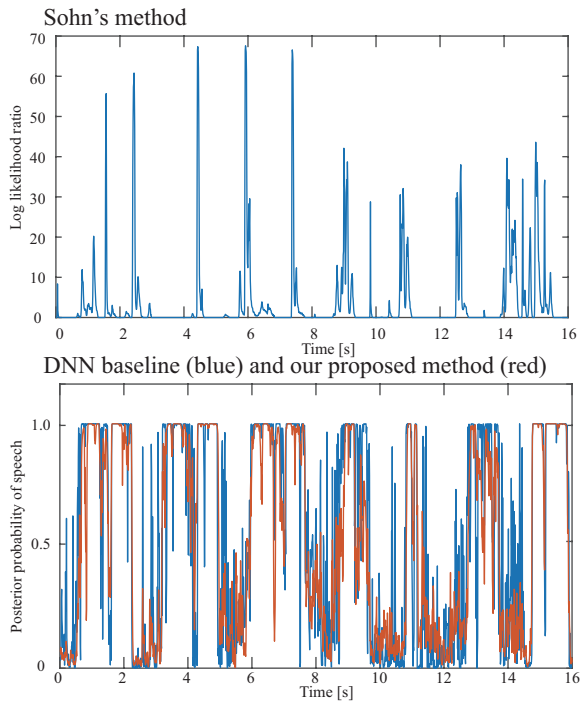


図 3 Log likelihood ratio of Sohn’s method,  $\log \Lambda$ , and the speech posterior probability of DNN baseline and the proposed method.

性能は、閾値  $\eta$  による影響を大きく受けることがわかる。これに対して DNN の結果は変動も少なく、閾値処理も容易である。青線の音響特徴量のみを用いた DNN よりも、提案の補助特徴量を併用した DNN の方がより、発話区間においては値が安定し、無音区間においては誤検出を抑制できていることがわかる。

### 5.7 スムージングの必要性

表 5 には、フレーム単位の音声区間検出結果に対して、隣接数フレームをまたいでスムージングした場合の性能を示す。DNN に基づく手法の場合に、顕著に性能が向上している。Sohn の方法では HMM hangover といったスムージング手法がすでに入っているが、DNN は入力特徴量の隣接コンテキストを利用することで暗に与えているだけなので、性能が向上した。

結合型に関しては、スムージングなしではあまり性能が向上しなかったものの、表 6 に示す通り、スムージングをすることで性能が大幅に改善した。これは 2 つの異なる特徴量を使うことでより安定的に音声区間が検出できるようになったためと考えられる。

### 5.8 音声認識での評価

CENSREC-2 付属のスクリプトにより学習した GMM 音響モデルにより、音声認識実験を行った。表 7 に単語正解率を示す。こちらも DNN に基づく音声区間検出を行った場合の性能が高く、補助特徴量の利用によりさらに性能が

表 5 Average frame-level VAD accuracy [%] with smoothing.

	e_a	c_a	i_a
Sohn	52.14	63.66	63.46
MLP (speech activation)	78.62	87.21	91.68
MLP (all activation)	80.14	88.47	92.76
MLP (mel)	80.91	89.02	94.03
+ speech activation	81.90	89.93	94.02
+ all activation	82.13	90.25	94.39
+ speech GMM	82.25	88.33	92.94
+ speech DNN	<b>82.68</b>	<b>91.19</b>	<b>94.91</b>

表 6 Average frame-level VAD accuracy [%] with smoothing. MLP used both auxiliary features.

activation	GMM/DNN	PCA	e_a	c_a	i_a
speech	GMM	no	81.87	90.77	94.26
all	GMM	no	82.12	90.16	94.28
speech	DNN	no	82.80	<b>91.14</b>	94.72
all	DNN	no	82.86	91.05	94.74
speech	GMM	yes	82.08	90.81	94.54
all	GMM	yes	82.32	89.96	94.26
speech	DNN	yes	<b>83.40</b>	<b>91.14</b>	<b>95.22</b>
all	DNN	yes	83.16	91.01	94.91

表 7 Word accuracy [%] of automatic speech recognition for the detected speech.

	e_a	c_a	i_a
Sohn	18.37	40.79	44.70
MLP (all activation)	67.76	73.99	82.71
MLP (mel)	69.33	78.30	87.25
+ speech activation	72.70	78.87	87.37
+ all activation	73.00	79.15	87.06
+ speech GMM	73.75	80.57	88.35
+ speech DNN	72.70	80.28	89.03
+ speech activation + GMM	73.15	80.11	88.72
+ all activation + GMM	74.79	81.19	89.21
+ speech activation + DNN	<b>76.51</b>	80.85	<b>89.88</b>
+ all activation + DNN	74.57	<b>81.98</b>	89.27

向上している。音声区間検出での取りこぼしは音声認識性能の低下に直結するため、高精度に音声区間検出を行う重要性が示された。また結合型により特に高騒音下での音声区間検出精度が大幅に改善された。

## 6. まとめ

DNN による音声区間検出の性能を向上させるために、補助音声モデルによる特徴量を併用する方法を提案した。CENSREC-2 による音声区間検出実験を行ったところ、従来のパワーに基づく方法よりも DNN に基づく方法の性能が顕著に高いことが分かった。また、NMF のアクティベーションや GMM/DNN 音響モデルのスコアを補助特徴量と

した実験により、補助特徴量の利用が有効であることを確認した。加えて音声認識実験においても、提案法の有効性を確認した。またこれら補助特徴量を結合して使うことにより、性能が改善し、特に音声認識性能の大幅な改善が見られた。

#### 参考文献

- [1] Sohn, J., Kim, N. S. and Sung, W.: A Statistical Model-based Voice Activity Detection, *IEEE Signal Processing Letters*, Vol. 6, pp. 1–3 (1999).
- [2] 太刀岡勇氣, 花沢利行, 成田知宏, 石井 純: 音声と騒音の密度比推定を用いた音声区間検出法, *電気学会論文誌 C*, Vol. 133, pp. 1549–1555 (2013).
- [3] Fujimoto, M. and Ishizuka, K.: Noise Robust Voice Activity Detection Based on Switching Kalman Filter, *IEICE Transactions on Information and Systems*, Vol. E91-D, pp. 467–477 (2008).
- [4] Zhang, X.-L. and Wu, J.: Deep Belief Networks Based Voice Activity Detection, *IEEE Transactions on Audio, Speech, and Language Processing*, Vol. 21, No. 4, pp. 1–14 (2013).
- [5] Sohrab, F. and Erdogan, H.: Recognize and Separate Approach for Speech Denoising Using Nonnegative Matrix Factorization, *Proceedings of EUSIPCO* (2015).
- [6] Kinoshita, K., Delcroix, M., Ogawa, A. and Nakatani, T.: Text-informed Speech Enhancement with Deep Neural Networks, *Proceedings of INTERSPEECH*, pp. 1760–1764 (2015).
- [7] Delcroix, M., Kinoshita, K., Hori, T. and Nakatani, T.: Context Adaptive Deep Neural Networks for Fast Acoustic Model Adaptation, *Proceedings of ICASSP*, pp. 4535–4539 (2015).
- [8] Teng, P. and Jia, Y.: Voice Activity Detection via Noise Reducing Using Non-Negative Sparse Coding, *IEEE Signal Processing Letters*, Vol. 20, No. 5, pp. 475–478 (2013).
- [9] Takeda, K., Fujimura, H., Itou, K., Kawaguchi, N., Matsubara, S. and Itakura, F.: Construction and Evaluation a Large In-Car Speech Corpus, *IEICE Transactions on Information and Systems*, Vol. E88-D, No. 3, pp. 553–561 (2005).