Deep Learning に基づくクリーン音声状態識別による雑音環境下音声認識*

☆柏木陽佑, 齋藤大輔, 峯松信明, 広瀬啓吉 (東大)

1 はじめに

近年,クリーン環境では自動音声認識の精度が大 きく向上し実用段階に入ってきている。しかし,ノイ ジー環境では未だにその性能は十分だとは言えない。 学習時と利用時において様々な環境ミスマッチが生じ るが,耐雑音性を高めることは非常に重要な課題で ある [1].

特徴量強調は耐雑音性を確保するためのフロントエ ンド処理であり、SPLICE [2] や Denoising AutoEncoder (DAE) [3], その他にも各種手法 [4–6] が提案 されている.これらは、ノイジー音声特徴量からク リーン音声特徴量を推定することによりミスマッチ を低減する.

SPLICE は 2 つの段階に分けることができる.ま ず,観測ノイジー音声特徴量空間をガウス混合モデル (Gaussian Mixture Model; GMM)でモデル化し,入 力フレームに対する,GMM コンポーネントの事後 確率を計算する.次に,計算した事後確率を重みとし て線形変換の足し合わせでクリーン音声特徴量を推 定する.しかし,ノイジー音声特徴量領域の分割とク リーン音声特徴量領域における分割は一致しないこ とが多く,局所領域での線形性の仮定が正しいとは限 らない.また,ノイジー音声特徴量領域の分割は雑音 の種類に大きく依存することも問題である.

そこで、クリーン音声特徴量領域で分割を行う手 法として REgularized piecewise linear mapping with DIscriminative region weighting And Long-span features (REDIAL) が提案されている [7–9]. REDIAL は GMM と線形判別分析 (Linear Discriminant Analysis; LDA) によりノイジー音声特徴量からクリーン音 声状態を識別的に推定する. REDIAL は特に HMM を雑音音声込みで学習する multi-condition で高い認 識性能を示すことが実験的に示されている. しかし, LDA は線形変換であるため、ノイジー音声特徴量と クリーン音声状態の複雑な関係を適切にモデル化で きているとは言えない.

一方,DAE はニューラルネットによって観測ノイ ジー音声特徴量からクリーン音声特徴量を非線形かつ 直接的に推定する.DAE を多層にした Deep Denoising AutoEncoder (DDAE) は特にノイズクローズド 環境で高い精度でクリーン音声を推定することが報 告されている.しかし、ノイズオープン環境では、過 学習の影響で大きく性能が低下する.

そこで、本提案手法は Deep Neural Network (DNN) [10] をクリーン音声状態の事後確率推定にの み用い、区分的線形変換によってクリーン音声特徴量 を推定する.まず、クリーン音声特徴量空間を GMM でモデル化し、クリーン音声状態ラベルを得る.ク リーン音声とノイジー音声は時間対応の取れている パラレルデータであるため、このクリーン音声状態ラ ベルを DNN により観測ノイジー音声特徴量より推定 を行い、クリーン音声状態に対する事後確率を得る. その後、REDIAL と同様に、観測ノイジー音声特徴 量から重み付きの線形変換によりクリーン音声特徴 量を推定する.

本稿の構成を示す.まず、2節で提案手法の定式化 を行う.3節で従来手法との比較を行い、4節で実験 を示す.最後に5節でまとめる.

2 定式化

特徴量強調は雑音環境下での音声認識精度を向上 させるため観測されたノイジー音声特徴量から対応 するクリーン音声特徴量を推定する.本研究では,観 測特徴量が得られた際のクリーン音声特徴量状態に 対する事後確率を DNN を用いて推定し,クリーン音 声特徴量を推定した事後確率を重み付けとした区分 的線形変換により求める.

DNN を採用する理由は、従来の区分的線形変換手 法は常に変換が線形で表現できるような領域に分割で きているとは限らないためである。例えば、SPLICE はノイジー音声特徴量空間を GMM でモデル化し、ノ イジー音声特徴量の状態に対する事後確率を重みとし て用いる。しかし、この重みの代わりに、クリーン音 声特徴量空間を GMM でモデル化して正解のクリー ン音声特徴量状態を用いて推定した事後確率を利用 すると, 高い精度でクリーン音声特徴量が推定するこ とができる。オラクルの結果と比較して SPLICE の 結果が大きく低下するのは、クリーン音声特徴量状 態とノイジー音声特徴量状態のミスマッチの影響だ と考えられる. REDIAL は LDA によってこの影響を 取り除くことを行っているが、LDA はあくまで線形 の写像なため限界がある.それに対して DNN は非線 形変換であるため、より正確にクリーン音声特徴量 状態を推定することができると考えられる.時刻 t に おける N 次元のクリーン音声特徴量 x_t とノイジー 音声特徴量 y_t のパラレルデータ { (x_t, y_t) } を考える. Fig. 1 に学習段階の流れを示す.まず、 クリーン音声 特徴量の確率密度関数 $p(\mathbf{x})$ を GMM で学習する.

$$p(\boldsymbol{x}_t) = \sum_k p(k)p(\boldsymbol{x}_t|k)$$
(1)

$$p(\boldsymbol{x}_t|k) = \mathcal{N}(\boldsymbol{x}_t; \boldsymbol{\mu}_k^{\boldsymbol{x}}, \boldsymbol{\sigma}_k^{\boldsymbol{x}})$$
(2)
$$p(k) = \pi_k^{\boldsymbol{x}}$$
(3)

これを用いて,GMM のコンポーネントインデクス に対する事後確率を,

$$p(k|\boldsymbol{x}_t) = \frac{p(\boldsymbol{x}_t|k)p(k)}{\sum_{k'} p(\boldsymbol{x}_t|k')p(k')}$$
(4)

と表すことができる.

一方,評価段階ではクリーン音声特徴量 *x*t は観測 できないため観測されたノイジー音声特徴量 *yt* から

^{*} Piecewise linear transformation based on discriminative approach using deep learning and its use for noise robust automatic speech recognition by Y.Kashiwagi D.Saito, N.Minematsu, and K.Hirose (The University of Tokyo)



Fig. 1 The training phrase of our proposed method.

Evaluation data

Deep Neural Network 🦟



Fig. 2 The testing phase of our proposed method.

クリーン音声状態に対する事後確率 $p(k|\mathbf{y}_t)$ を推定す る必要がある (Fig. 2). そこで、識別モデルとして DNN をクリーン音声特徴量状態と観測特徴量を用い て学習する.

$$p(k|\boldsymbol{y}_t) \simeq p(k|\boldsymbol{d}_t) = \operatorname{softmax}_k(\boldsymbol{V}\boldsymbol{h}^{(n)}(\boldsymbol{d}_t) + \boldsymbol{c})(5)$$

$$h^{(n)}(d_t) = \sigma(W^{(n)}h^{(n-1)}(d_t) + b^{(n)})$$
(6)

$$h^{(1)}(d_t) = \sigma(W^{(1)}d_t + b^{(1)})$$
(7)

ここで、 d_t は前後フレームを連結した入力特徴量ベクトルであり、

$$\boldsymbol{d}_{t} = [\boldsymbol{y}_{t-s}^{\top}, \dots, \boldsymbol{y}_{t-1}^{\top}, \boldsymbol{y}_{t}^{\top}, \boldsymbol{y}_{t+1}^{\top}, \dots, \boldsymbol{y}_{t+s}^{\top}]^{\top}$$
(8)

となる.また, σ はベクターシグモイド関数であり, $V, W^{(n)} \geq c, b^{(n)}$ はニューラルネットの重みとバ イアスのパラメータ, $h^{(n)}(y)$ は n 番目の隠れ層の出 力ベクトルである.なお,事前学習として各層の初期 値を Restricted Boltzmann Machine (RBM) で求め る [10].以上により,観測特徴量からクリーン音声特 徴量状態インデクスに対する事後確率 $p(k|y_t)$ を推定 することが可能となる.

評価段階では、クリーン音声特徴量 x_t を事後確率 $p(k|y_t)$ を重みとして用いる区分的線形変換によって 推定する.

$$\hat{\boldsymbol{x}}_t = \sum_k p(k|\boldsymbol{y}_t) \boldsymbol{A}_k \boldsymbol{e}_t \tag{9}$$

ここで、 A_k は、コンポーネントkにおける線形変換 行列であり、 e_t は前後フレームを連結した拡張行列

Table 1 Performance gap between ordinary SPLICE and the oracle SPLICE. (Word Error Rate (%))

	SPLICE	SPLICE
		(Oracle)
clean	0.57	0.69
SNR 20	1.08	0.76
SNR 15	1.99	0.70
SNR 10	4.65	0.77
SNR 5	16.76	0.93
SNR 0	49.96	0.95
SNR - 5	81.46	1.13
Average	14.89	0.82

^{/ork} である.

$$\boldsymbol{e}_{t} = [1, \boldsymbol{y}_{t-u}^{\top}, \dots, \boldsymbol{y}_{t-1}^{\top}, \boldsymbol{y}_{t}^{\top}, \boldsymbol{y}_{t+1}^{\top}, \dots, \boldsymbol{y}_{t+u}^{\top}]^{\top}$$
(10)

なお、A_kは重み付き最小二乗誤差基準で学習する.

$$\hat{\boldsymbol{A}}_{k} = \operatorname*{argmin}_{\boldsymbol{A}_{k}} \sum_{j} p(k|\boldsymbol{y}_{j}) ||\boldsymbol{x}_{j} - \boldsymbol{A}_{k} \boldsymbol{e}_{j}||^{2} \qquad (11)$$

これは解析解を得ることができる. A_k は,

$$\hat{\boldsymbol{A}}_k = \boldsymbol{X} \boldsymbol{P} \boldsymbol{E}^\top (\boldsymbol{E} \boldsymbol{P} \boldsymbol{E}^\top)^{-1}$$
(12)

として計算することができる.ここで, $X \in \mathcal{R}^{N \times T}$ と $E \in \mathcal{R}^{(N(2u+1)+1) \times T}$ はそれぞれ出力と入力特徴 量の拡張ベクトルを並べたデータ行列, $P \in \mathcal{R}^{T \times T}$ は $p(k|y_t)$ を対角に並べた行列である.

3 従来の特徴量強調手法との比較

本節では、提案手法と区分的線形変換手法である SPLICE とその拡張である REDIAL との比較を述べ る.また、ニューラルネットを用いた代表的な雑音抑 圧手法である DAE との比較も行う.

3.1 SPLICE

SPLICE は特徴量強調手法の1つであり、区分的線 形変換によってノイジー音声特徴量からクリーン音 声特徴量を推定する.

$$\hat{\boldsymbol{x}}_t = \sum_i p(i|\boldsymbol{y}_t) \boldsymbol{A}_i \begin{bmatrix} 1\\ \boldsymbol{y}_t \end{bmatrix}$$
(13)

SPLICE では、ノイジー音声特徴量の確率密度関数 $p(\mathbf{y})$ を GMM でモデル化する.

$$p(\boldsymbol{y}_t) = \sum_i \pi_i^{\boldsymbol{y}} \mathcal{N}(\boldsymbol{y}_t; \boldsymbol{\mu}_i^{\boldsymbol{y}}, \boldsymbol{\sigma}_i^{\boldsymbol{y}})$$
(14)

これを用いて、ノイジー音声特徴量状態に対する事後確率 $p(i|\mathbf{y}_t)$ を

$$p(i|\boldsymbol{y}_t) = \frac{p(\boldsymbol{y}_t|i)p(i)}{\sum_{i'} p(\boldsymbol{y}_t|i')p(i')}$$
(15)

として計算する.式 (14) のように,SPLICE ではノ イジー音声のみを用いて領域分割の学習を行う.しか し、これは学習データの雑音環境に過学習してしま うことが予想され、クリーン音声特徴量空間でモデ ル化を行うべきだと考えられる.

そこで、ノイジー音声特徴量領域とクリーン音声特 徴量領域のそれぞれで GMM によるモデル化と重み 付けを行った場合の比較を行った。Table 1 は Aurora 2 のテストセット A (ノイズクローズド環境)を用い た認識結果である。ここで、クリーン音声特徴量領域 で重み付けを行ったオラクルの SPLICE は、

$$\hat{\boldsymbol{x}}_t = \sum_{i^*} p(i^* | \boldsymbol{x}_t) \boldsymbol{A}_{i^*} \begin{bmatrix} 1 \\ \boldsymbol{y}_t \end{bmatrix}$$
(16)

$$p(i^*|\boldsymbol{x}_t) = \frac{p(\boldsymbol{x}_t|i^*)p(i^*)}{\sum_{i^{*'}} p(\boldsymbol{x}_t|i^{*'})p(i^{*'})}$$
(17)

としてクリーン音声状態 i^* に対する事後確率 $p(i^*|\mathbf{x}_t)$ を計算する.オラクルの結果は雑音の大きな環境でも頑健にクリーン音声を推定することができており、これは事後確率 $p(i^*|\mathbf{y}_t)$ を精確に推定することができていることに加え、クリーン音声ドメインにおける分割がより線形変換を仮定することができるためだと考えられる.なお、NMN-SPLICE は雑音の平均成分を除去することで分割をクリーン音声特徴量領域に近づくことが期待されている [2].

3.2 REDIAL

REDIAL は SPLICE を拡張したものであり, ク リーン音声特徴量に近づけた空間で分割を行う. 学習 段階では,まずクリーン音声の状態をラベルとした LDA により観測ノイジー特徴量を次元圧縮し,ノイ ジー特徴量をクリーン音声状態の識別が容易な特徴 量へ次元圧縮する.LDA の次元圧縮行列 L を対応す るクリーン音声特徴量の状態インデクスをソフトラ ベルとして学習する.次に,LDA により次元圧縮を 行ったベクトル $v_j = Ld_j$ を用いて K^* 混合の GMM を学習する.

$$p(\boldsymbol{v}_t) = \sum_{k^*=1}^{K^*} \pi_{k^*}^{\boldsymbol{v}} \mathcal{N}(\boldsymbol{v}_t; \boldsymbol{\mu}_{k^*}^{\boldsymbol{v}}, \boldsymbol{\sigma}_{k^*}^{\boldsymbol{v}})$$
(18)

入力特徴量が得られた際のクリーン音声特徴量状態 インデクスに対する事後確率 $p(k^*|\mathbf{y}_t)$ を

$$p(k^*|\boldsymbol{y}_t) \simeq p(k^*|\boldsymbol{v}_t) = \frac{p(\boldsymbol{v}_t|k^*)p(k^*)}{\sum_{k^{*'}} p(\boldsymbol{v}_t|k^{*'})p(k^{*'})} \quad (19)$$

として計算する.最終的に得られた事後確率を重み として用いた区分的線形変換によりクリーン音声特 徴量を

$$\hat{\boldsymbol{x}}_t = \sum_{k^*} p(k^* | \boldsymbol{y}_t) \boldsymbol{A}_{k^*} \boldsymbol{e}_t$$
(20)

$$\boldsymbol{e}_t = [1, \boldsymbol{y}_{t-u}^\top, \dots, \boldsymbol{y}_{t-1}^\top, \boldsymbol{y}_t^\top, \boldsymbol{y}_{t+1}^\top, \dots, \boldsymbol{y}_{t+u}^\top]^\top$$
(21)

として推定する.なお、 A_{k^*} は非常に大きな行列になるため、学習の際に正則化を導入する.

提案手法と REDIAL は、ノイジー音声特徴量から クリーン音声特徴量状態を推定する部分である。し かし、LDA は線形変換であるため、クリーン音声特 徴量状態とノイジー音声特徴量との間の複雑な関係 を適切にモデル化できているとは言えない。



Fig. 3 The performance of DDAE with different numbers of hidden layers in Aurora 2 dataset.



Fig. 4 The performance of the proposed method with different numbers of hidden layers in Aurora 2 dataset.

3.3 DAE

DAE はニューラルネットを用いてノイジー音声特 徴量からクリーン音声特徴量を直接的に推定する手 法である. DDAE は DAE を多層にしたものであり, クリーン音声特徴量を

$$\hat{\boldsymbol{x}}_t = \boldsymbol{U}\boldsymbol{h}^{(n)}(\boldsymbol{d}_t) + \boldsymbol{c}$$
(22)

$$h^{(n)}(d_t) = \sigma(W^{(n)}h^{(n-1)}(d_t) + b^{(n)})$$
 (23)

$$h^{(1)}(d_t) = \sigma(W^{(1)}d_t + b^{(1)})$$
 (24)

として推定する. ここで, U は重み行列である.本 稿では, RBM により各層の初期値を決め,最小二乗 誤差基準によるバックプロパゲーションにより学習 した DDAE との比較を行う. Fig. 3 は層の数を変化 した際の Aurora 2 における単語認識誤り率をプロッ トしたものである.セットA とセットB はそれぞれ ノイズクローズド環境とノイズオープン環境であり, 各隠れ層のノード数は 1024 で統一している.結果に よると,ノイズクローズド環境では層の数を増やすと 誤り率が減るが,ノイズオープン環境ではほぼ変化 しない.これはニューラルネットが過学習しているた めだと考えられる.なお,近年リカレント構造を持っ た複雑なニューラルネットを用いたオートエンコーダ も提案されている [11].

4 実験

提案手法の有効性を Aurora 2 による連続数字読 み上げ認識実験によって示す.データはいくつかの 環境の雑音が重畳されたノイジー音声とそれに対応

Table 2	Performance comparison among	our proposal	and conventional	methods (word error	rates $\%$).
	clean condition ()	WER %)	n	ulti condi	tion (WER	%)	

	clean condition (wER. 70)			multi condition (WER. 70)				
	set A	set B	set C	average	set A	set B	set C	average
Baseline	48.93	55.80	39.23	47.98	10.57	11.89	14.33	12.27
SPLICE	14.89	19.31	21.59	18.60	9.20	14.50	15.22	12.97
REDIAL	16.70	20.59	21.14	19.48	8.98	13.26	12.45	11.56
DDAE	6.39	20.44	17.20	14.68	5.97	18.50	14.67	13.04
PROPOSED	7.04	14.93	15.54	12.51	5.64	15.20	13.29	11.38

するクリーン音声が含まれており、パラレルデータ を利用することができる.認識に用いる音響モデル は HMM を利用し、クリーン音声のみで学習したも の (clean condition) と雑音音声込みで学習したもの (multi condition) の二通りで比較した.

評価データは3種類(A,B,C)があり,セットAは 学習時と同じノイズクローズド環境,セットBは学 習時と異なるノイズオープン環境,セットCはチャ ネルオープン環境である.特徴量としてMFCCとパ ワー,その1次,2次微分を用いた.DNNの学習は KALDI[13]を利用し,入力に用いる特徴量は全て前 後3フレームの計7フレームに統一した.また,各層 のノード数は1024に設定し,推定するクリーン音声 特徴量の状態は1024に設定した.DNNのバックプ ロパゲーションにおいては,学習データ8,440発話の うち844発話を開発セットとした.また,REDIAL と同様に線形変換の学習では正則化を導入している.

まず, DNN の層の数による認識率の違いを Fig. 4 に示す. 層の数を増やすとノイズクローズド環境では エラーレートが低下するが, ノイズオープン環境では DDAE と同様に効果が表れない. これは DDAE と同 様にニューラルネットが過学習しているためだと考え られる.

次に,従来手法との比較を行う.比較を行うのは SPLICE, REDIAL, DDAE である.SPLICE の/ イジー音声状態数は1024に設定した.REDIAL はク リーン音声状態数は1024に設定し,LDA による次 元圧縮後の特徴量ベクトルは64次元にした.また, DDAE は中間層の数を5とした.

Table 2 に結果を示す. clean condition では提案 手法は非常に良い結果が得られている.一方, multi condition では REDIAL とそれほど差が出ていない. これは LDA の線形変換という制約が逆に GMM の形 状を保持できるためだと考えられる.また, DDAE と比較した場合,特にノイズオープン環境で大きな 差が生じている.これは DDAE が過学習してしまっ ているためであり, GMM で分散を考慮したクラスラ ベルの識別を行う提案手法が過学習の影響を低下す ることができているからだと考えられる.

5 まとめ

本稿では、パラレルデータを用いた新しい特徴量 強調手法を提案した。提案手法は、クリーン音声特 徴量をGMMでクラスタリングを行い、DNNによっ てクリーン音声状態をノイジー音声特徴量から推定 する。その後、この事後確率を重みとして区分的線形 変換によりクリーン音声特徴量を推定する。Aurora 2を用いた連続数字読み上げ認識タスクによる評価で 提案手法は特にノイズオープン環境で DDAE より低 いエラーレートを示すことができた.今後,提案手法 を初期値として全体をニューラルネット形式で表現し て全体を最適化することが考えられる.

参考文献

- Gales, Mark JF, "Model-based approaches to handling uncertainty," Robust Speech Recognition of Uncertain or Missing Data-Theory and Applications. Springer, Berlin, Germany, pp. 101–125, 2011.
- [2] Droppo, Jasha and Deng, Li and Acero, Alex, "Evaluation of SPLICE on the Aurora 2 and 3 tasks," International Conference on Spoken Language Processing, pp. 29–32, 2002.
- [3] Vincent, Pascal and Larochelle, Hugo and Bengio, Yoshua and Manzagol, Pierre-Antoine, "Extracting and composing robust features with denoising autoencoders," Proceedings of the 25th international conference on Machine learning, pp. 1096–1103, 2008.
- [4] Li, Jinyu and Seltzer, Michael L and Gong, Yifan, "Improvements to VTS feature enhancement," Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on, pp. 4677–4680, 2012.
- [5] Afify, Mohamed and Cui, Xiaodong and Gao, Yuqing, "Stereo-based stochastic mapping for robust speech recognition," Audio, Speech, and Language Processing, IEEE Transactions on, vol. 17, no. 7, pp. 1325–1334, 2009.
- [6] Gemmeke, Jort F and Virtanen, Tuomas and Hurmalainen, Antti, "Exemplar-based sparse representations for noise robust automatic speech recognition," Audio, Speech, and Language Processing, IEEE Transactions on, vol. 19, no. 7, pp. 2067–2080, 2011.
- [7] Senior, Andrew and Cho, Youngmin and Weston, Jason "Learning improved linear transforms for speech recognition.," Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on, pp. 1957–1960, 2012.
- [8] Suzuki, Masayuki and Yoshioka, Takuya and Watanabe, Shinji and Minematsu, Nobuaki and Hirose, Keikichi "MFCC enhancement using joint corrupted and noise feature space for highly non-stationary noise environments," Acoustics, Speech and Signal Processing (ICASSP), 2012 IEEE International Conference on, pp. 4109–4112, 2012.
- [9] Suzuki Masayuki, Yoshioka Takuya, Watanabe Shinji, Minematsu Nobuaki, and Hirose Keikichi "Feature Enhancement With Joint Use of Consecutive Corrupted and Noise Feature Vectors With Discriminative Region Weighting," IEEE TASLP, (to appear).
- [10] Hinton, Geoffrey E and Osindero, Simon and Teh, Yee-Whye "A fast learning algorithm for deep belief nets," Neural computation, vol. 18, no. 7, pp. 1527–1554, 2006.
- [11] Maas, Andrew L and Le, Quoc V and O'Neil, Tyler M and Vinyals, Oriol and Nguyen, Patrick and Ng, Andrew Y, "Recurrent Neural Networks for Noise Reduction in Robust ASR," INTERSPEECH, 2012.
- [12] http://htk.eng.cam.ac.uk/
- [13] Povey, Daniel and Ghoshal, Arnab and Boulianne, Gilles and Burget, Lukas and Glembek, Ondrej and Goel, Nagendra and Hannemann, Mirko and Motlicek, Petr and Qian, Yanmin and Schwarz, Petr and others "The kaldi speech recognition toolkit," IEEE 2011 workshop on automatic speech recognition and understanding, 2011.