テンソル分解に基づく言語情報表現を用いた言語識別に関する検討* ☆鈴木 颯、齋藤 大輔、峯松 信明 (東大)

1 はじめに

言語識別は、入力音声の言語を判定する技術であ り、多言語対応の議事録作成支援アプリケーションや 国際コールセンター等で用いられている。入力音声か ら言語を特定するために必要な特徴を抽出すること で適切な識別が可能となることが期待されるが、音 声は話者やチャネル等の条件によって多様に変化し、 これら非言語的特徴の変動による識別性能低下が課 題の一つとなっている。近年では、これに着目して提 案された i-vector が標準的な言語情報表現の手法に なっているが [1]、これは各発話を Gaussian Mixture Model (GMM) でモデル化し、GMM の各分布の平均 ベクトルを連結した GMM supervector (GMM-SV) を次元圧縮することによって得られる言語特徴量で ある。

GMM-SV や i-vector は元来話者識別の分野で提 案された特徴量であり [2, 3]、この分野では近年、テ ンソル分解に基づく話者情報表現が提案されている [4]。これは、行および列がそれぞれ GMM の要素と 平均ベクトルに対応するような行列によって一発話 を表現し、多数話者分の行列をテンソルとして扱い、 テンソル解析を導入することで複数要因からの音響 的変動の分離を行なうという手法であり、言語識別に も適用可能であると考えられる。本稿では、テンソル 分解に基づく言語表象の効果を実験的に検討するこ とを目的とする。

2 先行研究

2.1 GMM-SV

GMM-SV は、一発話を GMM でモデル化し、各 分布の平均ベクトルを一列に連結した特徴量である [2]。GMM の各分布は音素や単音のような音響的な 要素を捉えていることが期待できるため、GMM-SV はそれらの音響的な特徴を発話単位で表現した特徴 量と考えることができる。

GMM-SV では一発話の GMM を推定することにな るが、これはあらかじめ様々な音声から統計的に話者・ 言語非依存の Universal Background Model (UBM) を構築しておき、UBM を事前分布として入力発話に 対する最大事後確率基準による推定 (MAP 推定) を行 なうことで得られる (Fig. 1 参照)。これにより、発話



GMM の各分布のインデックスと UBM の各分布の インデックスの対応付けが保たれるため、GMM-SV が発話間で比較可能になる。また、一発話という少量 のサンプルから安定に GMM を推定できるという利 点もある。

2.2 i-vector

ー発話から抽出された GMM-SV M は、UBM の GMM-SV m と、話者やチャネルの変化による音声 のばらつきをモデル化した低次元空間への射影行列 T (Total variability matrix) を用いて

$$\boldsymbol{M} = \boldsymbol{m} + T\boldsymbol{w} \tag{1}$$

と分解される。この w が、i-vector と呼ばれる [3]。 GMM-SV には言語識別においてノイズとなりうる話 者性・チャネルの成分が含まれたままであるが、言語 識別分野における i-vector では低次元空間への射影 によってこれらの成分の除去を行なっている。

Total variability matrix による射影は Principal Component Analysis (PCA) に相当する。話者識別 分野における i-vector は、GMM-SV に対する PCA という観点から考えると、次に述べる声質変換分野 における固有声 (EV) に基づく話者情報表現と基本 的に同一視できる。

2.3 固有声に基づく話者情報表現

固有声は音声認識におけるモデル適応法として提案 された技術であるが [5]、声質変換分野ではこれを適 用した固有声変換法 (Eigenvoice Conversion; EVC) が提案されている [6]。GMM に基づく声質変換では、 入力話者の特徴量 X_t と出力話者の特徴量 Y_t の結 合 GMM を用いて入出力の変換のモデル化を行なう が、EVC ではこの結合 GMM の出力話者側の平均べ

^{*}A study of language identification using tensor-based language identity representation. by S. Suzuki, D. Saito, and N. Minematsu (The University of Tokyo)

クトルを事前学習用のS人の話者の特徴を用いて表した EV-GMM $\lambda^{(EV)}$ に基づいて声質変換を行なう。

まず、事前学習用話者毎に GMM を学習し、GMM-SV を抽出する。S 個の GMM-SV に対して特異値分 解 (SVD) を行なうことでバイアスベクトルと $K(\leq S)$ 個の基底ベクトルを求め、これによって話者空間 を構築する (Fig. 2 参照)。すなわち、出力話者の GMM-SV $M^{(tar)}$ を以下のようにバイアスベクトル と基底ベクトルの線型結合で表す。

$$\boldsymbol{M}^{(tar)} = \boldsymbol{B}\boldsymbol{w} + \boldsymbol{b} \tag{2}$$

但し B は K 個の基底ベクトルからなる行列である。 このように出力話者の GMM-SV が K 次元の重みベ クトル w によって制御されるため、w が出力話者を 表現していると考えることができる。

重みベクトル w は、出力話者の音声データを用い て最尤基準に基づいて推定することができる。出力 話者の特徴量系列を $Y^{(tar)}$ とすると、w は以下のよ うに推定できる。

$$\hat{\boldsymbol{w}} = \operatorname*{argmax}_{\boldsymbol{w}} \int p(\boldsymbol{X}, \boldsymbol{Y}^{(tar)} | \boldsymbol{\lambda}^{(EV)}, \boldsymbol{w}) \mathrm{d}\boldsymbol{X} \qquad (3)$$

$$= \operatorname*{argmax}_{\boldsymbol{w}} p(\boldsymbol{Y}^{(tar)} | \boldsymbol{\lambda}^{(EV)}, \boldsymbol{w})$$
(4)

EM アルゴリズムを用いて、以下の更新式を得る。

$$\hat{\boldsymbol{w}} = \left\{ \sum_{m=1}^{M} \overline{\gamma}_{m}^{(tar)} \boldsymbol{B}_{m}^{\top} \boldsymbol{\Sigma}_{m}^{(YY)^{-1}} \boldsymbol{B}_{m} \right\}^{-1} \sum_{m=1}^{M} \boldsymbol{B}_{m}^{\top} \boldsymbol{\Sigma}_{m}^{(YY)^{-1}} \overline{\boldsymbol{Y}}_{m}^{(tar)}$$
(5)

$$\overline{\gamma}_m^{(tar)} = \sum_{t=1}^T \gamma_{m,t}, \ \overline{\boldsymbol{Y}}_m^{(tar)} = \sum_{t=1}^T \gamma_{m,t} (\boldsymbol{Y}_t^{(tar)} - \boldsymbol{b}_m^{(0)})(6)$$

$$\gamma_{m,t} = P(m|\boldsymbol{Y}_t^{(tar)}, \boldsymbol{\lambda}^{(EV)}, \boldsymbol{w})$$
(7)

EVC では事前学習用話者の情報を活用しているた め、データが少量でも効率的に出力話者を表現するこ とが可能になっている。しかし、GMM-SV は GMM の各分布の平均ベクトルが単純に一列に並んだ構成 になっているため、分布同士の関係性までは考慮され ていない。GMM の各分布には相関の高いものも低 いものも存在すると考えられるため、それらの関係性 を適切に考慮することでより精密な話者情報表現あ るいは言語情報表現が可能になることが期待できる。

3 テンソル分解に基づく言語情報表現

3.1 概要

GMM の各分布の関係性に着目した手法として、 テンソル分解に基づく話者情報表現を用いた声質変 換[7]と話者識別[4]が検討されている。[4,7]では、 GMM の各分布の平均ベクトル群を行列で表し、複





ソルに対して、PCA の射影行列の一つの計算法であ る SVD を拡張した Tucker 分解と呼ばれるテンソル 分解を用いて話者情報を表現している。本稿では、こ の手法を適用して言語情報表現を行なう。

3.2 多重線形解析

テンソルは、行列表現を一般化した多次元配列表 現である。テンソルにおける個々のインデックスは モードと呼ばれ、特定のモードに沿ってスライスする 平坦化操作によってテンソルを行列の形で表現でき る。 $A \in \mathcal{R}^{I_1 \times I_2 \times I_3}$ を3階のテンソルとすると、こ れをそれぞれモード 1, 2, 3に沿って平坦化した行列 $A_{(1)}, A_{(2)}, A_{(3)}$ は Fig. 3のようになる。このような 平坦化行列を用いて、テンソルと行列間の積が定義で きる。テンソル $\mathcal{G} \in \mathcal{R}^{I_1 \times \cdots \times I_N}$ と行列 $B \in \mathcal{R}^{J_n \times I_n}$ のモード n積 $A = \mathcal{G} \times_n B$ はモード nの平坦化行 列による演算 $A_{(n)} = B \cdot G_{(n)}$ によって定義される。

3.3 SVD と Tucker 分解

行列 $A \in \mathcal{R}^{M \times N}$ は、以下のように分解することができる。

$$\boldsymbol{A} = \boldsymbol{U}\boldsymbol{S}\boldsymbol{V}^{\top} \tag{8}$$

これを SVD と呼ぶ。但し、 $U \in \mathcal{R}^{M \times M}, S \in \mathcal{R}^{M \times N}, V \in \mathcal{R}^{N \times N}$ で U, V は直交行列、S は K

個の特異値からなる対角行列である。ここで、K は A のランクを表す ($K \leq \min(M, N)$)。行列 A の SVD を 2 階のテンソルの分解として書き直すと以下 のようになる。

$$\boldsymbol{A} = \boldsymbol{S} \times_1 \boldsymbol{U} \times_2 \boldsymbol{V} \tag{9}$$

これを以下のように高階のテンソルに拡張したもの を Tucker 分解と呼ぶ [8]。

$$\mathcal{A} = \mathcal{S} \times_1 \mathcal{U}_1 \times_2 \mathcal{U}_2 \times_3 \mathcal{U}_3 \tag{10}$$

但し、 U_1, U_2, U_3 が直交行列の場合はコアテンソル Sは密なテンソルとなる。

3.4 Tucker 分解による言語情報表現

ー発話の GMM の各分布の平均ベクトルを並べて $M \times D$ の行列で表現する。ここで M は混合数、D は 平均ベクトルの次元である。始めに全発話の平均を求 め、これをバイアス行列 b として各発話の行列から減 算しておく。学習データの発話数を N としたとき、N個の $M \times D$ 行列を 3 階のテンソル $M \in \mathcal{R}^{M \times D \times N}$ として考えて Tucker 分解を行なうと以下のように なる。

$$\mathcal{M} = \mathcal{G}^{M \times D \times N} \times_1 \boldsymbol{U}^{(M)} \times_2 \boldsymbol{U}^{(D)} \times_3 \boldsymbol{U}^{(N)} \quad (11)$$

但し $U^{(M)} \in \mathcal{R}^{M \times M}, U^{(D)} \in \mathcal{R}^{D \times D}, U^{(N)} \in \mathcal{R}^{N \times N}$ であり、それぞれ GMM 混合数、平均ベク トルの次元、発話インデックスの効果を捉えている。 ここで、以下のように M の第 3 モードを固定する ことで、特定の発話を表す行列が得られると考えら れる。

$$\hat{\boldsymbol{\mu}}^{(n)} = \boldsymbol{\mathcal{G}} \times_1 \boldsymbol{U}^{(M)} \times_2 \boldsymbol{U}^{(D)} \times_3 \boldsymbol{U}^{(N)}(n,:) \quad (12)$$

 $U^{(N)}(n,:) \in \mathcal{R}^{1 \times N}$ を重みパラメータ、その他のモード積の項を言語空間の基底と捉えるとSVDと等価になる。一方、本稿ではGMMの各分布の関係性を考慮するため以下のようなグルーピングを考える。

$$\hat{\boldsymbol{\mu}}^{(n)} = \boldsymbol{U}^{(M)} \left\{ \boldsymbol{\mathcal{G}} \times_2 \boldsymbol{U}^{(D)} \times_3 \boldsymbol{U}^{(N)}(n,:) \right\}$$
(13)

$$= \boldsymbol{U}^{(M)} \boldsymbol{W}_{n}^{\dagger} \tag{14}$$

 $U^{(M)}$ を基底、 $W_n \in \mathcal{R}^{D \times M}$ を重み行列とする。次 元圧縮の観点から、基底を縮約することで任意の発 話は以下のような行列で表される。

$$\boldsymbol{\mu}^{(new)} = \boldsymbol{U}^{(M)} \boldsymbol{W}_{(new)}^{\top} + \boldsymbol{b}$$
(15)

 $U^{(M)} \in \mathcal{R}^{M \times K} (K \leq N)$ が表現行列、 $W_{(new)} \in \mathcal{R}^{D \times K}$ が重み行列となり、 $D \times K$ の重み行列によって発話の言語情報を表現する。

Τa	ble 1 音響分析条件
サンプリング	8 bit / 8 kHz
窓	25 ms length / 10 ms shift
音響的特徴	MFCC (7 次元) + SDC (49 次元)

[4] では式 (15) の最小二乗解として重み行列 W を 算出しているが、これも EVC の重みベクトルと同様 に、以下のような最尤基準に基づいて W の推定を 行なうことができる [7]。

$$\hat{\boldsymbol{W}} = \operatorname*{argmax}_{\boldsymbol{W}} \int p(\boldsymbol{X}, \boldsymbol{Y}^{(tar)} | \boldsymbol{\lambda}^{(EV)}, \boldsymbol{W}) d\boldsymbol{X} \quad (16)$$
$$= \operatorname*{argmax}_{\boldsymbol{W}} p(\boldsymbol{Y}^{(tar)} | \boldsymbol{\lambda}^{(EV)}, \boldsymbol{W}) \quad (17)$$

EVC と同様に EM アルゴリズムを用いて、以下の更 新式を得る。

$$\operatorname{vec}(\boldsymbol{W}) = \left[\sum_{m=1}^{M} \overline{\gamma}_{m}^{(tar)} \boldsymbol{U}_{m}^{\top} \boldsymbol{U}_{m} \otimes \boldsymbol{\Sigma}_{m}^{(YY)^{-1}}\right]^{-1} \operatorname{vec}(\boldsymbol{C}) \quad (18)$$

$$\boldsymbol{C} = \sum_{t=1}^{T} \sum_{m=1}^{M} \gamma_{m,t} \boldsymbol{\Sigma}_{m}^{(YY)^{-1}} (\boldsymbol{Y}_{t}^{(tar)} - \boldsymbol{b}_{m}^{(0)}) \boldsymbol{U}_{m} \quad (19)$$

$$\boldsymbol{U}_m = \boldsymbol{U}^{(M)}(m, :) \in \mathcal{R}^{1 \times K}$$
(20)

$$\boldsymbol{b}_m^{(0)} = \boldsymbol{b}(m,:)^\top \in \mathcal{R}^{D \times 1}$$
(21)

ここで、vec() は行列を列ベクトルに展開する演算子 である。

4 言語識別実験

テンソル分解に基づく手法 (Tensor-based) の有効 性を検証するため、i-vector と固有声に基づく手法 (EV-based) と合わせて3つの手法で言語識別実験を 行なった。

4.1 実験条件

The National Institute of Standards and Technology Language Recognition Evaluation (NIST LRE) の 2003・2005・2007 年版を用いて、言語識別実験を 行なった。このコーパスには電話での会話が継続長 3s, 10s, 30s の 3 つのカテゴリーに分かれて収録され ている。NIST LRE 2007 Evaluation Plan に従い、 アラビア語・ベンガル語・ペルシア語・ドイツ語・日 本語・韓国語・ロシア語・タミル語・タイ語・ベトナ ム語・中国語・英語・ヒンドゥスタン語の14 言語を 識別対象の言語とした。

各音声データについて Table 1 の条件で音響分析を 行ない、Voice Activity Detection (VAD) と Cepstral Mean and Variance Normalization (CMVN), Vocal Tract Length Normalization (VTLN) を施した音響 特徴量系列を計算した。UBM は 2,048 混合、対角 共分散行列の条件で 24,577 発話から学習した。EVbased, Tensor-based における基底の個数 K はそれぞ

Table 2 14 言語 Cの認識率 [%]								
	WA			UA				
	3s	10s	30s	3s	10s	30s		
i-vector	44.40	66.36	79.19	36.90	60.95	77.00		
EV-based (MMSE)	35.87	53.20	65.20	28.05	43.92	56.41		
EV-based (ML)	38.83	58.80	70.11	32.10	52.05	63.42		
Tensor-based (MMSE)	38.23	60.56	75.67	31.42	53.56	70.64		
Tensor-based (ML)	31.23	56.81	70.44	26.08	51.75	66.27		

れ 600, 256 とした。i-vector の次元数は 600 とした。 識別には線形サポートベクターマシン (SVM) を用い、 23,665 発話で学習、NIST LRE 2007 Evaluation Plan で指定された 6,474 発話(3s, 10s, 30s 各 2,158 発話) でテストを行なった。音響特徴量系列と i-vector の計 算には KALDI Toolkit¹、GMM の学習には Hidden Markov Model Toolkit (HTK)²、SVM の実装には LIBLINEAR³を用いた。評価として、全テストデー タに対する認識率である Weighted Accuracy (WA) と各言語の認識率の平均である Unweighted Accuracy (UA) を計算した。

4.2実験結果

Table 2 に実験結果を示す。EV-based, Tensorbased で重みをそれぞれ式 (2), (15) の最小二乗解 として求めた場合を MMSE、重みをそれぞれ式 (4), (17)の最尤基準で求めた場合を ML と表記している。 MMSE に関しては WA, UA ともに Tensor-based の 認識率がEV-based よりも高くなっており、GMM の 各分布の関係性を考慮した Tensor-based の有効性を 示せた。一方、MMSE から ML になることで EVbased での認識率は向上したが、Tensor-based での認 識率は低下した。そのため ML に関しては、Tensorbased が EV-based の認識率を上回ったのは WA, UA ともに 30s のみとなった。Tensor-based における重 み行列は 56 × 256、EV-based における重みベクト ルは 600 次元となっており、Tensor-based の方が表 現力が高いために過学習を生じやすいことが考えら れる。

また、3 手法のうち認識率が最高となったのは ivector であった。i-vector と EV-based は、GMM-SV に対しての PCA という点では共通しているが、 i-vector では SVD による決定的な PCA ではなく Probablistic PCA (PPCA) に基づいて確率的なアプ ローチで言語空間の基底を求めているという違いが ある。両者の性能差はこれに起因すると考えられ、 Tensor-based にも PPCA を導入することで性能向 上が期待できる。

おわりに 5

本稿では、学習データセットをテンソルで表現して 分解することによって言語空間を構築し、GMM の 各分布の関係性に着目した言語情報表現を用いた言 語識別の検討を行なった。同様の検討を話者認識タス クとして行なった [4] と同様に、固有声に基づく言語 情報表現と比較して識別性能向上が確認できた。

今後は、テンソル分解に基づく重み行列を推定する 際の過学習を防ぐため、最尤推定と事前分布 (UBM) との内挿による MAP 的な重み行列の推定を検討し たい。また、PPCA を導入して確率的に言語空間の 基底を求めることについても検討したい。

本稿では識別に線形 SVM を用いたが、i-vector の 識別でよく用いられている Probabilistic Linear Discriminant Analysis (PLDA) を拡張し、重み行列を そのままの形状で入力できるような行列変量 PLDA を将来的には導入したい。

謝辞 本研究の一部は科研費・若手研究 (B) (25730105)、及び基盤研究 (A) (26240022) の助成を 受けたものである。

参考文献

- [1] N. Dehak et al., "Language Recognition via Ivectors and Dimensionality Reduction," Proc. IN-TERSPEECH, pp. 857–860, 2011.
- [2] W. M. Campbell et al., "Support Vector Machines using GMM Supervectors for Speaker Verification," IEEE Signal Processing Letters, vol. 13, pp. 308–311, 2006.
- [3] N. Dehak *et al.*, "Front-End Factor Analysis for Speaker Verication," IEEE Transactions on Audio, Speech, and Language Processing, vol. 19, no. 4, pp. 788–798, 2011.
- [4] チン・トゥアン・トゥー他、"テンソル分解に基づ く話者情報表現を用いた話者識別の検討."日本音 響学会春季講演論文集, pp. 217-220, 2015.
- [5] R. Kuhn et al., "Rapid Speaker Adaptation in Eigenvoice Space," IEEE Transactions on Speech and Audio Processing, vol. 8, no. 6, pp. 695-707, 2000.
- [6] T. Toda et al., "Eigenvoice Conversion Based on Gaussian Mixture Model," Proc. INTER-SPEECH, pp. 2446–2449, 2006.
- [7] D. Saito et al., "One-to-Many Voice Conversion Based on Tensor Representation of Speaker Space," Proc. INTERSPEECH, pp. 653–656, 2011.
- [8] L. R. Tucker, "Some mathematical notes on three-mode factor analysis," Psychometria, vol. 31, No. 3, pp. 279–311, 1966.

¹http://kaldi.sourceforge.net/

²http://htk.eng.cam.ac.uk/

³http://www.csie.ntu.edu.tw/~cjlin/liblinear/