

ギブス分布に基づく拡張ベイズ推定によるテキスト指定形話者照合

九州工業大学 ○植木琢也 溝部祐太 西田健 黒木秀一

Text-Prompted Multistep Speaker Verification Using Gibbs-Distribution-Based Extended Bayesian Inference

Takuya Ueki, Yuta Mizobe, Takeshi Nishida, and Shuichi Kurogi
Kyushu Institute of technology

Abstract: This paper presents a method of text-prompted multistep speaker verification for reducing verification errors. The method is developed for our speech processing system which utilizes competitive associative nets (CAN2s) for learning piecewise linear approximation of nonlinear speech signal to extract feature vectors of pole distribution from piecewise linear coefficients reflecting nonlinear and time-varying vocal tract of the speaker. This paper focuses on reducing verification errors by means of multistep verification using Gibbs-distribution-based extended Bayesian inference (GEBI) in text-prompted speaker verification. The effectiveness of GEBI and the comparison to BI (Bayesian inference) is shown and analyzed by means of experiments using real speech signals.

1. まえがき

本論文ではテキスト指定形話者照合の方法を提示する。テキスト指定形話者照合は詐称者や数字列のなりすましに対応するために開発された[1]。本論文では照合誤差の低減に焦点を当て、多段ベイズ(BI)推定を用いた方法とギブス分布に基づく拡張ベイズ(GEBI)推定を用いた方法の比較実験を行い、GEBI推定の有効性を示す。

さらに、近年音声信号が非線形であるという多くの報告があり、非線形を扱う手法により高い認識性能を実現できるのではないかと考え、本論文では非線形関数を区分的に線形近似する能力を持つニューラルネットの1つとして提案されている競合連想ネット(CAN2:Competitive Associative Net 2)を使用する。これまでにCAN2はアンサンブル学習の1つであるバギング法をCAN2に適用したbagging CAN2を用いて、音声時系列を性能よく再現し、また、再現する際に抽出される音声時系列の極の分布の配置パターンを特徴パターンとして用いた話者認識は、高い認識率を得ている[4,5]。

2. 多段ベイズ(BI)推定とギブス分布に基づく拡張ベイズ(GEBI)推定

多段ベイズ(BI)推定の問題を簡単に提示し(詳細については[5]を参照)、ギブス分布に基づく拡張ベイズ(GEBI)推定を導入する。まず、話者の集合を $S = \{s_i | i \in I_S\}$ 、数字の集合を $D = \{d_i | i \in I_D\}$ とする。さらに、学習機械として実数関数を学習する回帰学習機械(RLM: Regression Learning Machine; 詳細は後述)を用いることとし、その集合を $\text{RLM}^{[M]} = \{\text{RLM}^{[m]} | m \in M\}$ と表す。話者照合の場合、 $\text{RLM}^{[m]}$ は話者 $m \in M = S$ の音声信号を学

習し、数字照合の場合、 $\text{RLM}^{[m]}$ は数字 $m \in M = D$ の音声信号を学習するものとする。 $\text{RLM}^{[m_i]} (m_i \in M)$ の出力を $v^{[m_i]}$ ($= 1$ or -1) とし、 $\text{RLM}^{[M]}$ の出力を $\mathbf{v}^{[M]} = (v^{[m_1]}, \dots, v^{[m_M]})$ と表す。さらに、第 t 段(ステップ)での $\text{RLM}^{[M]}$ の出力を $\mathbf{v}_t^{[M]}$ とし、 $t = 1, 2, \dots, T$ に対する出力時系列を $\mathbf{v}_{1:T}^{[M]} = \mathbf{v}_1^{[M]} \mathbf{v}_2^{[M]} \dots \mathbf{v}_T^{[M]}$ と表す。すると多段BIによる $m (\in M = S \text{ or } D)$ に対する事後確率は次式で与えられる。

$$p_B(m | \mathbf{v}_{1:T}^{[M]}) = \frac{1}{Z_t} p_B(m | \mathbf{v}_{1:t-1}^{[M]}) p(\mathbf{v}_t^{[M]} | m) \quad (1)$$

ここで Z_t は正規化定数であり、 $\sum_{m \in M} p_B(m | \mathbf{v}_{1:t}^{[M]}) = 1$ をみたす。さらに $t = 1, 2, \dots$ とし次式を得る。

$$p_B(m | \mathbf{v}_{1:t}^{[M]}) = \frac{1}{Z_t} \exp\left(-t \left(\tilde{L}_{1:t}^{[m]} - \frac{1}{t} \log p_0(m)\right)\right) \quad (2)$$

ここで $p_0(m) = p_B(m | \mathbf{v}_{1:0}^{[M]})$ は事前確率を表し、 $\tilde{L}_{1:t}^{[m]} \equiv -\frac{1}{t} \left(\sum_{k=1}^t \log p(\mathbf{v}_k^{[M]} | m)\right)$ は正規化対数尤度である。また t 、 $\sum_{m \in M} p_B(m | \mathbf{v}_{1:t}^{[M]}) = 1$ より $m \in M$ と $m_v = \arg \max_{m_i \in M} p_B(m_i | \mathbf{v}_{1:t}^{[M]})$ の確率の比率は $t \rightarrow \infty$ のとき

$$\begin{aligned} r_{B,i,v} &\equiv \frac{p_B(m_i | \mathbf{v}_{1:t}^{[M]})}{p_B(m_v | \mathbf{v}_{1:t}^{[M]})} \\ &= \frac{p_B(m_i)}{p_B(m_v)} \exp\left(-t(\tilde{L}_{1:t}^{[m_i]} - \tilde{L}_{1:t}^{[m_v]})\right) \\ &\rightarrow \begin{cases} 1, m_i = m_v \\ 0, m_i \neq m_v \end{cases} \end{aligned} \quad (3)$$

となる。したがって、入力音声が無登録の場合でも、

登録話者または登録数字 $m = m_v$ に対する学習機械の事後確率 $p_B(m | \mathbf{v}_{1:T}^{[M]})$ が非常に大きくなる. この問題を回避するために以下に示すギブス分布を導入する.

$$p_G(m | \mathbf{v}_{1:T}^{[M]}) \equiv \frac{1}{Z_t} \exp\left(-\beta\left(\tilde{L}_{1:t}^{[m]} - \frac{1}{t} \log p_0(m)\right)\right) \quad (4)$$

ここで $p_0(m) = p_G(m | \mathbf{v}_{1:0}^{[M]})$ は事前確率を表し, β は逆温度である. 以下に示すように $m \in M$ と $m_v = \arg \max_{m \in M} p_G(m | \mathbf{v}_{1:T}^{[M]})$ の確率の比率は t の増加に伴い, $\frac{1}{t}$ より小さい値に収束する.

$$r_{G,i,v} \equiv \frac{p_G(m_i | \mathbf{v}_{1:t}^{[m]})}{p_G(m_v | \mathbf{v}_{1:t}^{[m]})} \rightarrow \exp\left(-\beta(\tilde{L}_{1:t}^{[m_i]} - \tilde{L}_{1:t}^{[m_v]})\right) \rightarrow c_i^\beta < 1 \quad (5)$$

このことは上記のBI推定の問題を回避できる可能性を示唆する. ここで式(4)より, GEBI確率の漸化式を以下のように導出する.

$$p_G(m | \mathbf{v}_{1:t}^{[M]}) \equiv \frac{1}{Z_t} p_G(m | \mathbf{v}_{1:t-1}^{[M]})^{\beta_t / \beta_{t-1}} p(\mathbf{v}_t^{[M]} | m)^{\beta_t} \quad (6)$$

式(6)において, $\beta_t = \beta/t$, ($t \geq 1$), $\beta_0 = 1$ とする. なお, 従来のBI推定法では $\beta_t = 1$, ($t \geq 0$) で与えられていた.

3. テキスト指定形話者認識システム

3.1 話者と数字の一段認識

CAN2を用いたテキスト指定形話者認識システムのブロック図を図1に示す. 一般に話者認識は音声をアナログ信号からデジタル信号へ変換する前処理, 特徴量抽出, パターンマッチング, 話者の決定の4工程で構成される. 話者認識は話者照合と話者識別に分類され, 本論文ではある者が本人の主張している通りの人物であるかを判断する話者照合を行う. また, テキストとして数字を使用しテキスト指定形話者照合を実現する. 話者 s または数字 d に対する回帰学習機械 $\text{RLM}^{[m]}$ ($m = s \text{ or } d$) の目的関数は次式で表される.

$$y^{[m]}(\mathbf{q}) = \begin{cases} 1, & \text{if } \mathbf{q} \in \mathbf{Q}^{[m]} \\ -1, & \text{otherwise} \end{cases} \quad (7)$$

ここで $\mathbf{Q}^{[m]}$ は $m (= s \text{ or } d)$ に対する特徴ベクトル \mathbf{q} の集合を表す. この目的関数の学習後, $\text{RLM}^{[m]}$ は連続関数 $\hat{y}^{[m]} = \hat{f}^{[m]}(\mathbf{q}^{[m]})$ を出力するが, これを次式より2値化して一段照合あるいは2クラス分類を行う.

$$v^{[m]} = \begin{cases} 1, & \text{if } \hat{y}^{[m]} \geq y_\theta^{[m]} \\ -1, & \text{otherwise} \end{cases} \quad (8)$$

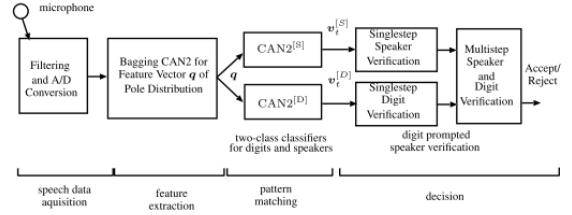


Fig.1. Diagram of text-prompted speaker verification system using CAN2s

すなわち, $v^{[m]} = 1$ の場合 m を受理し, それ以外の場合は棄却する. ここで閾値 y_θ は2.3に示すように安定した照合結果が得られるように調整する.

3.2 数字と話者の多段照合

テキスト指定形話者照合を実行するため, 数字列と話者の多段照合を並列に行い, その結果を組み合わせ, 数字列が受理された場合のみ話者照合を行う. 話者または数字の多段照合の学習機械 $m \in M (= S \text{ or } D)$ の特徴ベクトル $\mathbf{q} \in \mathbf{Q}$ から得た $\text{RLM}^{[M]}$ の二値化出力ベクトル $\mathbf{v}^{[M]} = (v^{[m_1]}, \dots, v^{[m_M]})$ の条件付き確率を式(9)に示す.

$$p(\mathbf{v}^{[M]} | s) = \prod_{m_i \in M} p(v^{[m_i]} | m) \quad (9)$$

さらに, 話者または数字の入力時系列 $m_{1:T} = m_1 m_2 \dots m_T$ に対する出力時系列 $\mathbf{v}_{1:T}^{[M]}$ を $\mathbf{v}_{1:T}^{[M]} = \mathbf{v}_1^{[M]} \mathbf{v}_2^{[M]} \dots \mathbf{v}_T^{[M]}$ とする. 出力 $\mathbf{v}_{1:T}^{[M]}$ が参照ト入力(正解とすべき入力) $m_{1:T}^{[r]}$ に対する応答であるかを判定するために, GEBI推定を使用し, 式(10)(11)に示すように $t = 1, 2, \dots, T$ のときの2つの再帰的な事後確率を算出する.

$$p_G(m_{1:t}^{[r]} | \mathbf{v}_{1:t}^{[M]}) = \frac{1}{Z_t} p_G(m_{1:t-1}^{[r]} | \mathbf{v}_{1:t-1}^{[M]})^{\beta_t / \beta_{t-1}} p(\mathbf{v}_t^{[M]} | m_t^{[M]})^{\beta_t} \quad (10)$$

$$p_G(\overline{m}_{1:t}^{[r]} | \mathbf{v}_{1:t}^{[M]}) = \frac{1}{Z_t} p_G(\overline{m}_{1:t-1}^{[r]} | \mathbf{v}_{1:t-1}^{[M]})^{\beta_t / \beta_{t-1}} p(\mathbf{v}_t^{[M]} | \overline{m}_t^{[M]})^{\beta_t} \quad (11)$$

ここで, Z_t は $p_G(m_{1:t}^{[r]} | \mathbf{v}_{1:t}^{[M]}) + p_G(\overline{m}_{1:t}^{[r]} | \mathbf{v}_{1:t}^{[M]}) = 1$ を満たす正規化定数である. また,

$$p(\mathbf{v}_t^{[M]} | \overline{m}_t^{[r]}) = \sum_{m \in M \setminus \{m_t^{[r]}\}} p(\mathbf{v}_t^{[M]} | m) / (|M| - 1) \quad (12)$$

とする.

$t = T$ の時の数字の照合を式(13)により行う.

$$V_{1:T}^{[D]} = \begin{cases} 1, & \text{if } p_G(d_{1:T}^{[r]} | \mathbf{v}_{1:T}^{[D]}) \geq p_\theta^{[D]} \\ -1, & \text{otherwise} \end{cases} \quad (13)$$

さらに T ステップにおける話者と数字の照合, すな

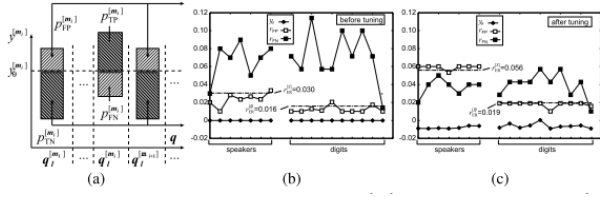


Fig.2. (a) Schematic relationship between the threshold $y_\theta^{[m_i]}$ and the performance ratios $r_{FP}^{[m_i]}$, $r_{FN}^{[m_i]}$, $r_{TP}^{[m_i]}$, $r_{TN}^{[m_i]}$. The horizontal axis indicates all training feature vectors $\mathbf{q} = \mathbf{q}_l^{[m_i]} (l=1,2,\dots)$ obtained from speakers $m=s \in S$ or digits $m=d \in D$. The vertical axis indicates the corresponding output $y^{[m_i]}$ of $\text{RLM}^{[m_i]}$. Experimental results (see 3) of y_θ, r_{FP}, r_{FN} for the original $y_\theta = 0$ and tuned y_θ are shown in (b) and (c), respectively.

わちテキスト指定形話者照合の照合は

$$V_{i,T}^{[SD]} = \begin{cases} 1, & \text{if } (V_{i,T}^{[D]} = 1) \wedge (p_G(s_{i,T}^{[r]} | \mathbf{v}_{i,T}^{[S]}) \geq p_\theta^{[S]}) \\ -1, & \text{otherwise} \end{cases} \quad (14)$$

より行う。 $p_\theta^{[D]}$ と $p_\theta^{[S]}$ は閾値である。以上より、 $V_{i,T}^{[SD]} = 1$ のとき入力を受理し、それ以外の際は棄却する。

3.3 回帰学習機械の閾値 y_θ の調整

全ての話者と全ての数字に同じ閾値を使用するために、以下の手順により話者 $m=s \in S$ と数字 $m=d \in D$ それぞれについての式(8)における閾値 $y_\theta = y_\theta^{[m]}$ を調整する。まず、入力 m の特徴ベクトル \mathbf{q} は $\text{RLM}^{[m_i]} (m_i \in S)$ によってFP(偽陽性)とFN(偽陰性)に分類される。ランダムに選択されたデータに対するそれらの割合は次のように推定される。

$$r_{FP}^{[m_i]} = \frac{1}{|M|-1} \sum_{m \in M \setminus m_i} p(v^{[m_i]} = 1 | m) \quad (15)$$

$$r_{FN}^{[m_i]} = p(v^{[m_i]} = -1 | m)$$

また、TP (真陽性) とTN (真陰性) の比率は $r_{TP}^{[m_i]} = 1 - r_{FN}^{[m_i]}$ と $r_{TN}^{[m_i]} = 1 - r_{FP}^{[m_i]}$ である。ただし、各 $\text{RLM}^{[m_i]}$ は誤差率の平均二乗誤差 $\langle (v^{[m_i]} - y^{[m_i]})^2 \rangle$ を最小化するように学習してある。

$$r_{ER}^{[m_i]} = \frac{1}{|M|} = (r_{FN}^{[m_i]} + (|M|-1)r_{FP}^{[m_i]}) \quad (16)$$

図2(a)より、 $y_\theta^{[m_i]}$ の値を大きくすると $r_{FP}^{[m_i]}$ が減少し、 $r_{FN}^{[m_i]}$ が増加していることがわかる。この関係より、三つの要素 $(y_{\theta,n}^{[m_i]}, r_{FP,n}^{[m_i]}, r_{FN,n}^{[m_i]})$ からなる $\delta^{[m_i]}$ を得る。

試用閾値は

$$y_{\theta,n}^{[m_i]} = (n/n_y) (y_P^{[m_i]} - y_N^{[m_i]}) + y_N^{[m_i]} \quad (17)$$

である。ただし、 $n \in I_y = \{0,1,2,\dots,n_y\}$ とする。また $y_P^{[m_i]}, y_N^{[m_i]}$ は全ての学習データの特徴ベクトル \mathbf{q} の出力 $y^{[m_i]}$ の真と偽の平均であり、 n_y は分割数である。 $\delta^{[m]} (m \in M)$ のデータを式(18)により分割する。

$$\delta_l^{[m]} = \left\{ (y_\theta, r_{FP}, r_{FN}) \in \delta^{[m]} \mid l = \arg \min_{l \in I_y} |r_{FP} - r_l| \right\} \quad (18)$$

ただし、

$$r_l = (l/n_y) (r_1 - r_0) + r_0, (l \in I_y)$$

$$r_0 = \min \{ r_{FP,n}^{[m]} \mid m \in M, n \in I_y \} \quad (19)$$

$$r_1 = \max \{ r_{FP,n}^{[m]} \mid m \in M, n \in I_y \}$$

とする。その後 $(y_\theta, r_{FP}, r_{FN}) \in \delta^{[m]}$ の平均と差異を計算する。ここで $(E_{\delta_l^{[m]}}(y_\theta), E_{\delta_l^{[m]}}(r_{FP}), E_{\delta_l^{[m]}}(r_{FN}))$ と $(V_{\delta_l^{[m]}}(y_\theta), V_{\delta_l^{[m]}}(r_{FP}), V_{\delta_l^{[m]}}(r_{FN}))$ はそれぞれ $m \in M, l \in I_y$ とする。 $l \in I_y$ を r_{FP} および r_{FN} の全ての分散の和を最小化するものとし、式(20)のように示す。

$$\tilde{l} = \arg \min_{l \in I_y} \sum_{m \in M} (V_{\delta_l^{[m]}}(r_{FP}) + V_{\delta_l^{[m]}}(r_{FN})) \quad (20)$$

これより、全ての $m \in M$ における r_{FP} と r_{FN} の分散を小さくすることができる。

4. 実験

4.1 実験方法

音声データは研究室の静かな部屋の中でサンプリング周期8[kHz]、分解能16[bit]の下でサンプリングしたを使用した。7人の話者の集合を $S = \{\text{fHS}, \text{fMS}, \text{mKK}, \text{mKO}, \text{mMT}, \text{mNH}, \text{mYM}\}$ 、10個の数字の集合を $D = \{\text{/zero/}, \text{/ichi/}, \text{/ni/}, \text{/san/}, \text{/yon/}, \text{/go/}, \text{/roku/}, \text{/nana/}, \text{/hachi/}, \text{/kyu/}\}$ とする。話者の各数字は2か月間で異なる日時で10サンプルとったものである。数字は $x = x_{s,d,l} (s \in S, d \in D, w \in W)$ によって発声されたものとし、データセットを $X = \{x = x_{s,d,l} \mid s \in S, d \in D, l \in L\}$ とする。これらの音声時系列について、out-of-bag(OOB)によって極分布による評価実験を行った。OOBはLOOCV(Leave-one-set-out-cross-validation)より小さい傾きと分散を有すると期待されている[7]。なお、実験結果はテスト話者 s の数字列 $d_{i,T} = d_1 d_2 \dots d_T$ とそれに対応する参照話者 $s^{[r]}$ の数字列 $d_{i,T}^{[r]} = d_1^{[r]} d_2^{[r]} \dots d_T^{[r]}$ について、 $r_{CD} = n_{CD}/T$ の比率で正しい数字 $d_i = d_i^{[r]}$ を含むという条件の下ランダムに $d_i, s, d_i^{[r]}, s^{[r]}$ を選択した。

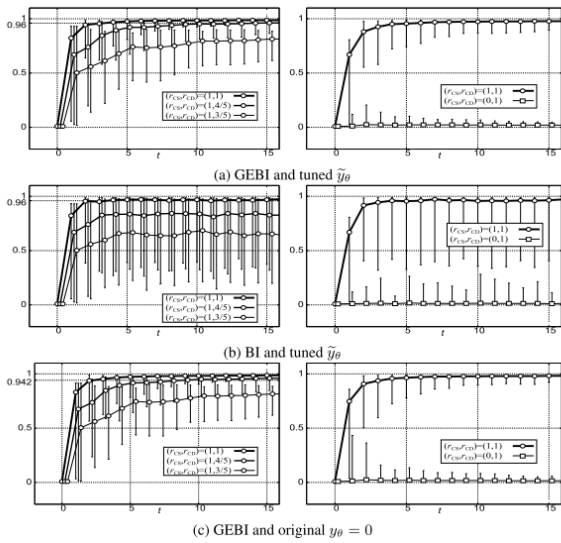


Fig. 3. Experimental result of multistep probability for digits (left) and speakers (right). The plus and minus error bars indicate RMS (root mean square) of positive and negative errors from the mean, respectively. The curves for different datasets are shifted slightly and horizontally to avoid crossovers.

4.2 実験結果と考察

本実験ではテスト数字列は5桁とし、また $T = 15 = 5 \times 3$ である。表1に実験結果を示す。表1の $(r_{CS}, r_{CD}) = (0,1)$ の行は不正解話者 ($r_{CS} = 0$) が5桁の数字を全て正しく ($r_{CD} = 1$) 発声した結果を示し、その下にはそれぞれ正解話者 ($r_{CS} = 1$) が5桁の数字を $r_{CD} = 1, 0, 4/5, 3/5$ と発声したときの結果である。誤棄却率FRRは $FRR = 100 - r_{acc}^{[SD]} [\%]$ で与えられる。方法(a)では $(r_{CS}, r_{CD}) = (1,1)$ において、ステップ数 t が増加すると誤棄却率は減少する。さらに、 $(r_{CS}, r_{CD}) = (1, 1/3)$ のにおいて、誤受理率FARも減少する。この単調減少の性質は以下の式(4)に示すGEBIの確率 $p_G(m | \mathbf{v}_{lr}^{[M]})$ がステップ数 t の増加によって単調に収束する性質から得られる。式(2)で与えられるBIの確率 $p_B(m | \mathbf{v}_{lr}^{[M]})$ もまたステップ数 t の増加に伴い単調に収束するが、GEBI推定に比べて収束が遅い。すなわち、 $\tilde{L}_{lr}^{[m]} \gg |\log p_0(m)|/t$ を満たす t において $p_G(m | \mathbf{v}_{lr}^{[M]})$ は収束に達するが $p_B(m | \mathbf{v}_{lr}^{[M]})$ は収束に達しないことがわかる。図3(b)では確率曲線の誤差範囲が大きく、しかも変動している。従って、テキスト指定形多段話者照合の照合誤差について、GEBI推定は有効であると考えられる。

Table 1. Experimental result of acceptance rates $r_{acc}^{[D]}$ and $r_{acc}^{[SD]}$ achieved by the methods using (a) GEBI and tuned \tilde{y}_θ (c) GEBI and original $y_\theta = 0$, and four datasets with $(r_{CS}, r_{CD}) = (0,1), (1,1), (1,1/4)$ and $(1,1/3)$. The values of $r_{acc}^{[D]}$ and $r_{acc}^{[SD]}$ are expressed by the rate[%] to the total 1000 test sequences for each case. The thresholds are $p_\theta^{[D]} = 0.96$ for (a) and (b), 0.942 for (c), $p_\theta^{[S]} = 0.5$ and $T = 15$.

r_{CS}	r_{CD}	(a) GEBI & \tilde{y}_θ		(b) BI & \tilde{y}_θ		(c) GEBI & $y_\theta = 0$	
		$r_{acc}^{[D]}$	$r_{acc}^{[SD]}$	$r_{acc}^{[D]}$	$r_{acc}^{[SD]}$	$r_{acc}^{[D]}$	$r_{acc}^{[SD]}$
0	1	97.5	0.0	94.9	0.0	98.4	0.0
1	1	98.1	98.1	96.3	94.4	98.0	98.0
1	4/5	76.3	76.3	72.5	70.5	84.8	84.8
1	3/5	0.2	0.2	44.5	43.7	0.8	0.8

5. 結論

照合誤差を低減するためにGEBIを用いたテキスト指定形多段話者照合の方法を提案した。実音声信号を用いた確率および実験結果の分析により、GEBIの確率はBIの確率より安定ではるかに誤差率を低減することを示した。本論文ではテスト話者は登録話者のみであったが、今後の追加研究として、テスト話者が未登録話者の場合の方法を検討する。

参考文献

- [1] Beigi, H.: Fundamentals of speaker recognition. Springer-Verlag New York Inc (C) (2011)
- [2] Melin, H., Lindberg, J.: Prompting of Passwords in Speaker Verification Systems, Fonetik-97, Phonum 4, Umea University, Sweden, May 28-30. (1997)
- [3] Kurogi, S., Ueno, T. and Sawa, M.: A batch learning method for competitive associative net and its application to function approximation. Proc. SCI2004, Vol. V, pp.24-28 (2004)
- [4] Kurogi, S., Mineishi, S. and Sato, S.: An analysis of speaker recognition using bagging CAN2 and pole distribution of speech signals. Proc. ICONIP2010, Part I, LNCS 6443, pp.363-370 (2010)
- [5] Mizobe, Y., Kurogi, S., Tsukazaki, T., and Nishida, T.: Multistep speaker identification using Gibbs-distribution-based extended Bayesian inference for rejecting unregistered speaker. Proc. ICONIP2012, Part V, LNCS 7667, pp.247-255 (2013)
- [6] Campbell, J.P.: Speaker Recognition: A Tutorial. Proc. the IEEE, Vol. 85, No.9, pp.1437-1462 (1997)
- [7] Kurogi, S.: Improving generalization performance via out-of-bag estimate using variable size of bags, J. Japanese Neural Network Society, Vol. 16, no. 2, pp.81-92 (2009)