

音響特性補正の導入による肉伝導音声変換の収録環境適応

宮本 大輔[†], 中村 圭吾[†], 戸田 智基[†], 猿渡 洋[†], 鹿野 清宏[†]

[†] 奈良先端科学技術大学院大学 情報科学研究科

肉伝導音声変換は Non-Audible Murmur (NAM) マイクロフォンで収録される肉伝導音声の音質向上に効果的である。この手法では、肉伝導音声から空気伝導音声へ変換するための確率モデルが事前に学習される。肉伝導音声の音響特性は、NAM マイクロフォンの圧着位置などの収録環境に敏感であるため、実際の使用においては学習時と変換時の音響特性の不一致により、しばしば変換音質が大きく劣化する。この問題に対して、我々は肉伝導音声変換のための Cepstrum Mean Subtraction (CMS) と Constrained Structural Maximum A Posteriori Linear Regression (CSMAPLR), または Signal Bias Removal (SBR) と CSMAPLR の組み合わせに基づく教師無しの音響特性補正法を提案する。実験結果から、提案手法により音響特性の不一致に起因する変換音質の劣化が大幅に低減されることを示す。

Adaptive Approach to Varying Recording Conditions in Body Transmitted Voice Conversion Based on Acoustic Compensation

Daisuke MIYAMOTO[†] Keigo NAKAMURA[†] Tomoki TODA[†] Hiroshi SARUWATARI[†]
Kiyohiro SHIKANO[†]

[†] Graduate school of Information Science, Nara Institute of Science and Technology

Body transmitted voice conversion is very effective for enhancing body transmitted speech recorded with Non-Audible Murmur (NAM) microphone. In this method, a probabilistic model to convert body transmitted speech into natural speech is trained previously. Because acoustic characteristics of body transmitted speech is sensitive to recording conditions such as a location of NAM microphone, significant degradation of the conversion performance is often caused in practical situations by acoustic mismatches between the training and the conversion processes. To alleviate this problem, we propose unsupervised acoustic compensation methods based on combination of Cepstrum Mean Subtraction (CMS) and Constrained Structural Maximum A Posteriori Linear Regression (CSMAPLR), or combination of Signal Bias Removal (SBR) and CSMAPLR for body transmitted voice conversion. Experimental results demonstrate that the proposed methods significantly reduce the quality degradation of the converted speech caused by the acoustic mismatches.

1 はじめに

近年、携帯電話の普及により様々な場所で音声コミュニケーションをとることが可能になった。しかし、人ごみや道路の側などの騒音環境下では相手に声を伝えることは困難である。また、図書館などの静穏環境下では、声を出すことで周囲に迷惑を掛けてしまうため小声で話さなければならず、同様に相手に声を伝えるのは困難である。このような状況下においても快適な音声コミュニケーションを実現するため、外部雑音に頑健であり、小さな声を収録可能な体内伝導音収録用マイクroフォンが開発されている^{1, 2)}。

我々は、体内伝導音収録用マイクroフォンの一つとして Non-Audible Murmur (NAM) マイクroフォン²⁾に着目している。NAM マイクroフォンは耳介後方に圧着して、声道内の空気振動を皮下の筋肉組織を通じて収録する。このマイクroフォンの特徴は、通常の空気伝導マイクroフォンと比べてダイナミックレンジが大きく、通

常の音量の音声から周囲に聞こえないほど小さな咳き声である NAM まで収録できることである。しかし、NAM マイクroフォンで収録される肉伝導音声は、口唇の放射特性の欠落や体内伝導によるローパス特性などの影響により、通常の空気伝導音声と比べて音質が大きく劣化する問題がある。

肉伝導音声の音質改善のために、統計的声質変換³⁾に基づく肉伝導音声変換が提案されている⁴⁾。肉伝導音声変換では、同一内容を発声した空気伝導音声と肉伝導音声の平行データを 50 文程度用いて、事前に変換モデルを学習しておく。この変換モデルを用いることで、肉伝導音声から空気伝導音声への変換が可能となる。学習時と変換時でデータの収録環境が同一であれば、肉伝導音声変換は極めて有効であり、音質が大きく改善される。しかし、肉伝導音声はマイクroフォンを体表に直接圧着して収録するため、NAM マイクroフォンの圧着位置などの収録環境が変化することで、収録される肉伝

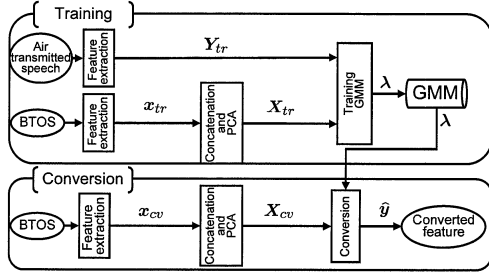


Fig. 1 Framework of body transmitted voice conversion.

導音声の音響特性が大きく変化する。収録環境の変化に起因する音響特性の変化は空気伝導音声と比べて大きく、その結果、肉伝導音声変換による変換音質が大きく劣化してしまう。

この問題を解決するため、我々は肉伝導音声変換のための音響特性補正技術の研究に取り組んでいる⁵⁾。これまでに、Cepstrum Mean Subtraction (CMS)⁶⁾、Constrained Structural Maximum A Posteriori Linear Regression (CSMAPLR)⁷⁾に基づく音響特性補正の導入を行い、その有効性を確認している⁵⁾。CMSに基づく補正法では、1発話ごとに入力データをバイアス補正することで変換音質の劣化を大きく抑制できる。また、CSMAPLRに基づく補正法では、適応データを用いて新しい収録環境に変換モデルを適応させることで、数文以上の適応発声でCMSを上回る性能が得られる。本稿では、1文ごとのバイアス補正とモデル適応を組み合わせさせた手法として、CMSとCSMAPLR、Signal Bias Removal (SBR)⁸⁾とCSMAPLRを組み合わせさせた手法をそれぞれ導入する。客観および主観評価実験結果から、各音響特性補正法の性能を示す。

2 肉伝導音声変換

本稿では通常発声をNAMマイクロフォンで収録した肉伝導通常音声 (Body Transmitted Ordinary Speech: BTOS) から、通常の空気伝導マイクロフォンで収録された通常音声への変換を扱う。Fig. 1に肉伝導音声変換の枠組みを示す。

2.1 肉伝導通常音声 (BTOS)

空気伝導音声と比較すると、BTOSは特に高域の周波数成分が大きく減衰している。このため、BTOSは酷くこもった音声となる。さらに、無声摩擦音などの高域に大きなパワーをもつ音素は、音韻特徴が失われやすい。これらの要因から、BTOSは空気伝導音声に比べて音質が大きく劣化しており、明瞭性も低い。

2.2 音響特徴量

本稿で用いる肉伝導音声変換では、スペクトル特徴量としてメルケプストラムを使用する。肉伝導によって欠落した情報を補うため、入力特徴量としてスペクトルセグメントを用いる。時刻 t におけるメルケプストラムベクトル x_t の前後 n フレームを連結させたベクトル $c_t = [x_{t-n}^T, \dots, x_t^T, \dots, x_{t+n}^T]^T$ (T は転置) を作成し、主成分分析 (PCA) により次元圧縮することで時刻 t におけるスペクトルセグメント X_t を得る。

$$X_t = Dc_t - d \quad (1)$$

ここで D は PCA 変換行列、 $d = D\bar{c}$ であり、 \bar{c} は学習データ全体における c_t の平均ベクトルを表す。また、時刻 t における出力特徴量として、静的特徴量 y_t と動的特徴量 Δy_t を連結させた静的・動的特徴量 $Y_t = [y_t^T, \Delta y_t^T]^T$ を用いる。

2.3 最尤基準に基づく特徴量変換³⁾

学習処理では、入力と出力の時間方向のアライメントをとった特徴量の結合確率密度を以下のように Gaussian Mixture Model (GMM) でモデル化する。

$$P(Z_t | \lambda) = \sum_{m=1}^M w_m \mathcal{N}(Z_t; \mu_m^{(Z)}, \Sigma_m^{(ZZ)}) \quad (2)$$

ここで Z_t は入力特徴量と出力特徴量の結合ベクトルであり $Z_t = [X_t^T, Y_t^T]^T$ で構成される。記号 $\mathcal{N}()$ は正規分布を表す。混合数 M の GMM のモデルパラメータセット λ は各分布 m の重み w_m 、平均ベクトル $\mu_m^{(Z)}$ 、共分散行列 $\Sigma_m^{(ZZ)}$ で構成される。ここで $\mu_m^{(Z)}$ と $\Sigma_m^{(ZZ)}$ はそれぞれ以下のように構成される。

$$\mu_m^{(Z)} = \begin{bmatrix} \mu_m^{(X)} \\ \mu_m^{(Y)} \end{bmatrix} \quad (3)$$

$$\Sigma_m^{(ZZ)} = \begin{bmatrix} \Sigma_m^{(XX)} & \Sigma_m^{(XY)} \\ \Sigma_m^{(YX)} & \Sigma_m^{(YY)} \end{bmatrix} \quad (4)$$

$\mu_m^{(X)}$ と $\mu_m^{(Y)}$ はそれぞれ分布 m の入力特徴量の平均ベクトルと出力特徴量の平均ベクトルを示す。 $\Sigma_m^{(XX)}$ と $\Sigma_m^{(YY)}$ はそれぞれ分布 m の入力特徴量、出力特徴量の共分散行列であり、 $\Sigma_m^{(XY)}$ 、 $\Sigma_m^{(YX)}$ は相互共分散行列を示す。

特徴量変換を行う際は、入力および出力の特徴量系列ベクトルをそれぞれ $X = [X_1^T, \dots, X_T^T]^T$ 、 $Y = [Y_1^T, \dots, Y_T^T]^T$ とし、以下の尤度関数の最大化に基づき特徴量変換を行う。

$$P(Y|X, \lambda) \simeq P(m|X, \lambda) P(Y|X, m, \lambda) \quad (5)$$

ここで、 $m = [m_1, \dots, m_T]^T$ は分布系列を表す。まず、準最適な分布系列 \hat{m} を次式で決定する。

$$\hat{m} = \arg \max_m P(m|X, \lambda) \quad (6)$$

この分布系列 \hat{m} を用いて変換特徴量の静的特徴量系列 \hat{y} を次式で得る。

$$\hat{y} = \arg \max_y P(Y|X, \hat{m}, \lambda) \quad (7)$$

このとき $Y = E\hat{y}$ であり、 E は静的特徴量系列から静的・動的特徴量系列へと変換する行列を表す。さらに、系列内変動を考慮することで、大幅に変換音質が改善される³⁾。

本変換法においては、確率密度 $P(Z_t | \lambda)$ が学習時と変換時で変化しないという前提が必要となる。

3 音響特性補正法

NAMマイクロフォンは体表に直接圧着して使用するため、収録される肉伝導音声の音響特性は圧着位置や接触圧などの収録環境に敏感である。このような収録環境の変化に対応するために音響特性の補正を試みる。本稿ではより簡便で実用的な方法として、出力音声である空気伝導音声や言語情報を用いず、入力音声である肉伝導音声のみを用いた補正法を提案する。

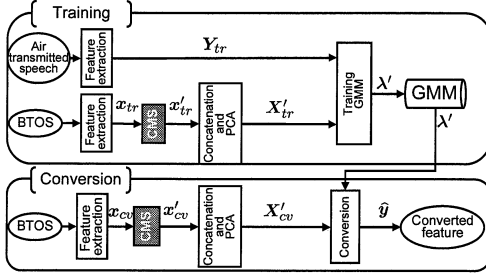


Fig. 2 Framework of body transmitted voice conversion with acoustic compensation based on CMS.

3.1 Cepstrum Mean Subtraction (CMS)に基づく補正法⁵⁾

CMS⁶⁾は静的な乗法性歪みの影響を除去する手法である。学習時と変換時それぞれの入力特徴量に対してCMSを行い、収録環境の違いによって変化する静的な乗法性歪みの影響を除去し、変換音質の劣化を抑制する。

Fig. 2にCMSに基づく音響特性補正を導入した肉伝導音声変換の枠組みを示す。Section 2.2で示した入力メルケプストラムに対してCMSを導入すると、正規化された特徴量が次式で与えられる。

$$x'_t = x_t - \bar{x} \quad (8)$$

\bar{x} は x_t の長時間平均であり、本稿では発話ごとに計算された平均ベクトルを用いる。この正規化されたメルケプストラムを用いてスペクトルセグメント X'_t を求める。

$$X'_t = D'c'_t - d' \quad (9)$$

このとき、 $c'_t = [x'_{t-n}, \dots, x'_t, \dots, x'_{t+n}]^T$ 、 $d' = D'c'$ であり、 D' と c' はそれぞれPCA変換行列と学習データ全体における c'_t の平均ベクトルを示す。

3.2 Constrained Structural Maximum A Posteriori Linear Regression (CSMAPLR)に基づく補正法⁵⁾

CSMAPLR⁷⁾はモデル適応手法の一つであり、CMSと比べてより複雑な変換処理を行う。Fig. 3にCSMAPLRに基づく音響特性補正を導入した肉伝導音声変換の枠組みを示す。

CSMAPLRはConstrained Maximum Likelihood Linear Regression (CMLLR)⁹⁾と同様、入力特徴量に対してアフィン変換を行うことでモデルとデータの不一致を補正する。入力特徴量 X_t の変換特徴量は次式で与えられる。

$$\hat{X}_t = AX_t + b = W\xi_t \quad (10)$$

このとき、 $W = [b, A]$ は変換行列であり、 ξ_t は特徴量の拡張ベクトル $\xi_t = [1, X_t^T]^T$ である。ここで、Fig. 4に示すような回帰木を用いて、適応データ量に従って複数の変換行列を推定し、変換を行うことで複雑な変換を実現する。各分布に使用する変換行列は、閾値以上のデータ量を持つ上位ノード(クラス) r で推定される変換行列 $W_r = [b_r, A_r]$ を使用する。CSMAPLRでは、変換行列を推定する際に木構造を利用したMAP推定を行

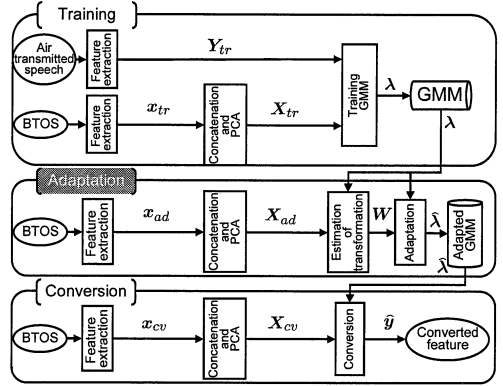


Fig. 3 Framework of body transmitted voice conversion with acoustic compensation based on CSMAPLR.

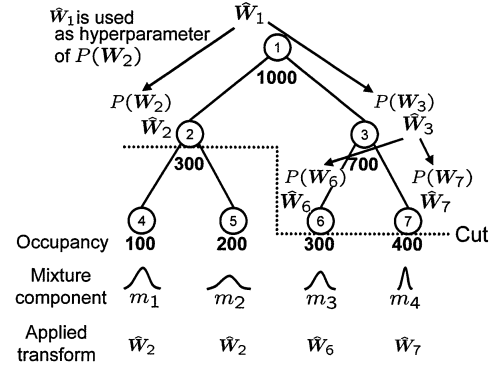


Fig. 4 An example of usage of regression tree with CSMAPLR when setting occupancy threshold to 300.

う。Fig. 4に示すように、ハイパーパラメータとして親ノードの変換行列の推定値を使用することで、過学習を抑制する。

出力特徴量を用いない教師無し適応を行うため、適応データ特徴量系列 X の周辺分布の尤度を最大化する変換行列 \hat{W}_r を推定する。

$$\hat{W}_r = \arg \max_{W_r} \int P(X, Y | W_r, \lambda) P(W_r) dY \quad (11)$$

ここで事前分布 $P(W_r)$ は行列変量正規分布で表され、次式で定義される。

$$P(W_r) \propto |\Psi|^{-(d+1)/2} |\Phi|^{-d/2} \cdot \exp \left[-\frac{1}{2} \text{tr} \{ (W_r - H)^T \Psi^{-1} (W_r - H) \Phi^{-1} \} \right] \quad (12)$$

ここで、 d は入力特徴量の次元数であり、 Ψ 、 Φ 、 H はそれぞれハイパーパラメータである。 H は回帰木の親ノードで推定された変換行列を用い、 Ψ 、 Φ はそれぞれ $\Psi^{-1} = C \cdot I$ 、 $\Phi^{-1} = I$ とする。 I は単位行列を示

し、 C は $P(W_r)$ の大きさを決定する係数である。 \hat{W}_r はEMアルゴリズムで求められ

$$\hat{w}_{ri} = (\alpha c_i + \mathbf{n}^{(ri)}) \mathbf{V}^{(rii)^{-1}} \quad (13)$$

で与えられる。ここで、 \hat{w}_{ri} は \hat{W}_r の i 行成分、 c_i は A_r の i 行目の余因子ベクトルの拡張ベクトルである。 α は二次方程式を解くことで得られる値である⁹⁾。 $\mathbf{n}^{(ri)}$ と $\mathbf{V}^{(rii)}$ はそれぞれ以下の式で与えられる。

$$\mathbf{V}^{(rii)} = \mathbf{G}^{(rii)} + C \cdot \mathbf{I} \quad (14)$$

$$\mathbf{n}^{(ri)} = \mathbf{k}^{(ri)} + C \cdot \mathbf{h}(i) \quad (15)$$

$$\mathbf{G}^{(rij)} = \sum_{m=1}^{M_r} p_m(i, j) \sum_{t=1}^T \gamma_m(t) \xi_t \xi_t^\top \quad (16)$$

$$\mathbf{k}^{(ri)} = \sum_{m=1}^{M_r} p_m(i) \mu_m \sum_{t=1}^T \gamma_m(t) \xi_t^\top - \sum_{j=1, j \neq i}^d w_{rj} \mathbf{G}^{(rij)} \quad (17)$$

ここで、 $p_m(i)$ と $p_m(i, j)$ はそれぞれ共分散行列 $\Sigma_m^{(XX)}$ の i 行成分と (i, j) 要素を示す。また、 M_r と $\gamma_m(t)$ はそれぞれクラス r に属する分布数と分布 m の時間 t における生起確率 $\gamma_m(t) = P(m | \mathbf{X}_t, W_r, \lambda)$ を示す。

本稿ではCSMAPLRで推定された変換行列をモデル空間に適用し、適応された変換モデルを用いて変換処理を行う。適応されたモデルパラメータは以下の形で与えられる。

$$\hat{\mu}_m^{(Z)} = \begin{bmatrix} A_r' \mu_m^{(X)} - b_r' \\ \mu_m^{(Y)} \end{bmatrix} \quad (18)$$

$$\hat{\Sigma}_m^{(Z)} = \begin{bmatrix} A_r' \Sigma_m^{(XX)} A_r'^\top & A_r' \Sigma_m^{(XY)} \\ \Sigma_m^{(YX)} A_r'^\top & \Sigma_m^{(YY)} \end{bmatrix} \quad (19)$$

このとき、 $\hat{\mu}_m^{(Z)}$ と $\hat{\Sigma}_m^{(ZZ)}$ はそれぞれ適応された分布 m の平均ベクトルと共分散行列を表す。また、 $A_r' = A_r^{-1}$ 、 $b_r' = A_r^{-1} b_r$ である。

3.3 バイアス補正とCSMAPLRを組み合わせた補正法

1文ごとにバイアス補正を行うCMSは、適応データ数で変換を推定するCSMAPLRと比べると、発話ごとの変動に対して頑健であると考えられる。そこで、Fig. 5に示すようにCMSとCSMAPLRを組み合わせることで、更なる性能改善を試みる。学習時は、CMSのみを用いた補正と同様に変換モデルを学習する。適応時は、適応データから得られるCMSによって補正された特徴量を用いて、CSMAPLRによって回帰行列を推定し、変換モデルを適応する。変換時は、適応された変換モデルを用い、CMSのみを用いた補正と同様に変換処理を行う。

CMSによって得られるバイアスは1文の平均として計算される。この場合、後段のCSMAPLRで用いられる評価基準を一切考慮しておらず、CSMAPLRにおける回帰行列の推定に悪影響を及ぼす可能性がある。そこで、バイアスを尤度基準で推定するSignal Bias Removal (SBR)⁸⁾を導入する。また、本稿ではSBRによってバイアスを推定する際にCMSの結果を初期値として与える。Fig. 6にSBRとCSMAPLRを組み合わせた補正を導入した肉伝導音声変換の枠組みを示す。

学習時は、まずCMSのみを用いて補正を行いGMMを学習する。得られたGMMを用いて、1文ごとにSBR

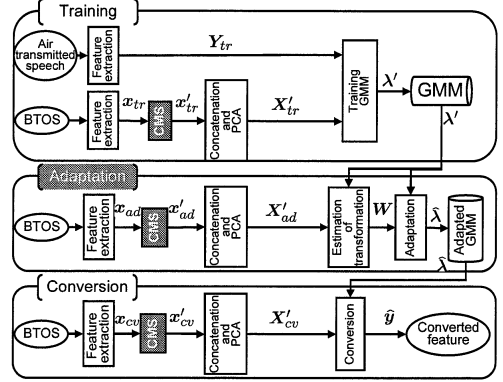


Fig. 5 Framework of body transmitted voice conversion with acoustic compensation based on the method combining CMS and CSMAPLR.

により補正を行う。SBRによって得られる特徴量は、バイアス g を用いて

$$\tilde{X}_t = X'_t - g \quad (20)$$

で与えられる。学習時は出力特徴量の情報が得られるため、尤度を最大化する入力特徴量のバイアス \hat{g} は

$$\hat{g} = \arg \max_g P(X', Y | g, \lambda) \quad (21)$$

として、各発話に対して推定される。 \hat{g} はEMアルゴリズムで求められ、次式で与えられる。

$$\hat{g} = B^{(Z)^{-1}} (l^{(Z)} - q^{(Z)}) \quad (22)$$

このとき $q^{(Z)}$ 、 $l^{(Z)}$ 、 $B^{(Z)}$ は以下で与えられる。

$$q^{(Z)} = \sum_m^M [P_m^{(XX)} P_m^{(XY)}] \mu_m^{(Z)} \sum_t^T \tilde{\gamma}_m^{(Z)}(t) \quad (23)$$

$$l^{(Z)} = \sum_m^M [P_m^{(XX)} P_m^{(XY)}] \sum_t^T \tilde{\gamma}_m^{(Z)}(t) Z'_t \quad (24)$$

$$B^{(Z)} = \sum_m^M P_m^{(XX)} \sum_t^T \tilde{\gamma}_m^{(Z)}(t) \quad (25)$$

ここで、 $\tilde{\gamma}_m^{(Z)}(t) = P(m | Z'_t, g, \lambda)$ であり、 $Z'_t = [X'_t{}^\top, Y_t{}^\top]^\top$ である。また、 $P_m^{(XX)}$ 、 $P_m^{(XY)}$ 、 $P_m^{(YX)}$ 、 $P_m^{(YY)}$ はそれぞれ $\Sigma_m^{(XX)}$ 、 $\Sigma_m^{(XY)}$ 、 $\Sigma_m^{(YX)}$ 、 $\Sigma_m^{(YY)}$ と同じサイズであり、 $\Sigma_m^{(ZZ)^{-1}}$ を構成する行列である。

$$\Sigma_m^{(ZZ)^{-1}} = \begin{bmatrix} P_m^{(XX)} & P_m^{(XY)} \\ P_m^{(YX)} & P_m^{(YY)} \end{bmatrix} \quad (26)$$

得られたバイアスの最尤推定値を用いて特徴量の補正を行い、GMMのパラメータを再学習する。さらに、得られたGMMを用いて上記処理を繰り返す。本学習処理は、発話毎の適応パラメータとGMMパラメータを同一の尤度関数に基づいて更新するものであり、発話適応学習とみなせる。

適応時は学習時に得られたGMMを用いてSBRによるバイアス推定を行う。また、変換時はCSMAPLRで

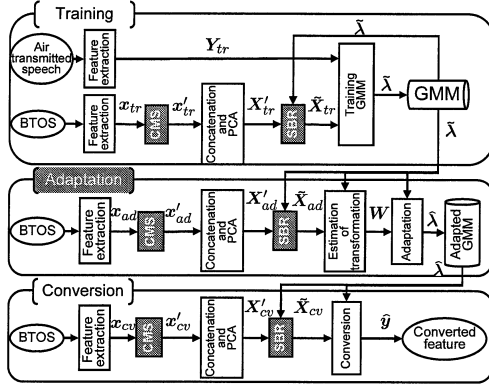


Fig. 6 Framework of body transmitted voice conversion with acoustic compensation based on the method combining SBR and CSMAPLR.

適応された GMM を用いて SBR によるバイアス推定を行う。適応および変換時は、入力特徴量のみを用いてバイアスを推定するため、 X' の周辺分布を最大化する \hat{g} を各文ごとに推定する。

$$\hat{g} = \arg \max_g \int P(X', Y | g, \lambda) dY \quad (27)$$

\hat{g} は EM アルゴリズムで求められ次式で与えられる。

$$\hat{g} = B^{(X)^{-1}}(l^{(X)} - q^{(X)}) \quad (28)$$

このとき $q^{(X)}$, $l^{(X)}$, $B^{(X)}$ は、生起確率 $\tilde{\gamma}_m^{(X)}(t) = P(m | X'_t, g, \lambda)$ を用いて以下で与えられる。

$$q^{(X)} = \sum_m \Sigma_m^{(X,X)^{-1}} \mu_m^{(X)} \sum_t \tilde{\gamma}_m^{(X)}(t) \quad (29)$$

$$l^{(X)} = \sum_m \Sigma_m^{(X,X)^{-1}} \sum_t \tilde{\gamma}_m^{(X)}(t) X'_t \quad (30)$$

$$B^{(X)} = \sum_m \Sigma_m^{(X,X)^{-1}} \sum_t \tilde{\gamma}_m^{(X)}(t) \quad (31)$$

4 評価実験

4.1 実験条件

実験用音声として NAM マイクロフォンとヘッドセットマイクロフォンを同時に装着し、BTOS と空気伝導音声を同時に収録した。2名の男性話者から、NAM マイクロフォンの圧着位置のみ変化させ（首の左右を入れ替え）、他の収録環境は可能な限り保持した2つの環境で同一内容の発話を収録した。まず、NAM マイクロフォンの圧着位置を固定したまま音素バランス文を100文収録した。次に NAM マイクロフォンの圧着位置を変更し、再度同一内容の音素バランス文を100文収録した。1話者は1組（100 × 2文）、もう1話者は2組の音声を収録し、合計3組の音声セットを実験に用いた。音声のサンプリング周波数は8 kHz とした。

収録環境が異なる条件での変換音質の評価を行うために、Fig. 7 に示すようなデータセットを用いた。一方の

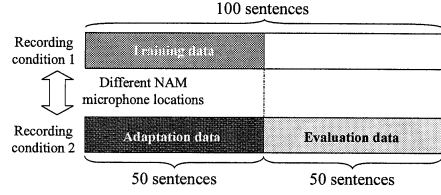


Fig. 7 Data sets for experiments.

環境の50文を変換モデル学習に用い、もう一方の環境の、モデル学習に使用しなかった内容の50文を用いて評価を行った。CSMAPLRの変換行列推定に用いる適応文は評価文と同一環境の、評価に用いない50文から必要量を選択し使用した。この適応データ50文は、同一環境の変換モデル学習時にも使用した。CMS及びSBRによる補正では、各文ごとにバイアスを算出した。交差検定として、学習環境と評価環境を入れ替えて同一の評価を行った。

音響特性補正法として、収録環境が変化した場合の、1) 補正無し、2) CMS、3) CSMAPLR、4) CMS と CSMAPLR の組み合わせ、5) SBR と CSMAPLR の組み合わせを導入した肉伝導音声変換をそれぞれ評価した。また、6) 同一の収録環境で変換を行った場合の評価も行った。これらの変換手法に対して客観的・主観的に評価を行った。客観評価実験の評価指標は目標音声と変換音声のパワー項を含むメルケプストラム歪みを用いた。主観評価実験は、一対比較法で音質を評価した。刺激音声として、上述の1), 2), 3), 5), 6) を行った場合の変換音声を提示した。CSMAPLRに使用する適応文は1文あるいは10文を用いた場合の変換音声をそれぞれ評価した。刺激音声は、異なる補正法を用いて変換された2種類の音声をランダムに提示した。被験者は補正法のすべての組み合わせを評価し、1文あたり20通りの組み合わせで補正法を評価した。被験者は9名とした。

スペクトル特徴量として、0から16次までのメルケプストラム係数を用いた。また、入力特徴量として前後4フレーム連結したベクトルを34次元に圧縮したスペクトルセグメントを用いた。GMMの混合数は64とした。CSMAPLRにおいて、クラス決定のためのデータ量の閾値を1000とし、ハイパーパラメータ C の値は 10^5 として実験を行った。

4.2 実験結果

Fig. 8 に客観評価実験結果を示す。この結果より、適応文数が少ない場合 CMS の方が性能が高く、適応文数が多くなると CSMAPLR は CMS を上回る性能が得られることが分かる。CMS と CSMAPLR を組み合わせた補正では、CSMAPLR のみの補正と比べると、適応文数が極端に少ない場合は性能が高いが、適応文数が多くなると性能が悪化することが分かる。これは、CMS によるバイアス推定が不安定であるため、CSMAPLR の回帰行列推定精度が低下していることが原因であると考えられる。一方で、SBR と CSMAPLR を組み合わせることで、適応文数が大きくなってでも CSMAPLR のみの補正と同程度の性能が得られる。適応文数が十分に大きくなると CSMAPLR のみの補正を導入した枠組みで変換した結果に漸近することから、発話ごとの変動は GMM

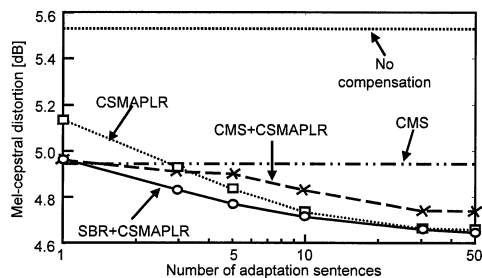


Fig. 8 Mel-cepstral distortion as a function of the number of adaptation sentences. The mel-cepstral distortion in the matched condition is 3.93 dB.

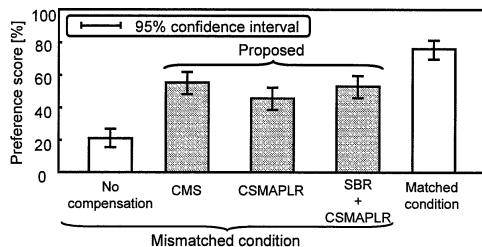


Fig. 9 Results of preference test on speech quality. One sentence was used for adaptation data in the CSMAPLR compensation.

で表現できる程度であり、発話適応処理を行う必要が無いことが分かる。また、適応文数が極端に少ない場合には CSMAPLR のみの補正と比べて大きく性能が向上していることが分かる。このことから、適応分数が少ない場合に CSMAPLR で生じる過学習の影響は、SBR を用いて CSMAPLR の推定精度を改善させることにより、効果的に緩和できることが分かる。

Fig. 9 に適応文を 1 文用いた場合、Fig. 10 に適応文を 10 文用いた場合の主観評価実験結果を示す。適応文が 1 文の場合、CMS に基づく補正と SBR と CSMAPLR の組み合わせに基づく補正が同程度の性能を示しており、CSMAPLR のみの補正を用いる場合よりも音質が改善される傾向がある。また、適応文が 10 文程度得られた場合、CSMAPLR 単体で用いる場合と SBR と組み合わせで用いた場合で同程度の音質であり、これらは CMS のみを用いた場合よりも音質が良いことが分かる。これらの結果は客観評価実験結果と同様の傾向を示しており、妥当であると考えられる。

5 おわりに

学習時と変換時で音響特性が異なる場合における肉伝導音声変換の音質改善のため、CMS と CSMAPLR の組み合わせ、および SBR と CSMAPLR の組み合わせに基づく音響特性補正の提案を行った。客観および主観評価実験結果から、バイアス補正による発話ごとの変動を補正する必要性は低いが、SBR が CSMAPLR の推定精度を向上させ、極端に適応文が少ない場合において

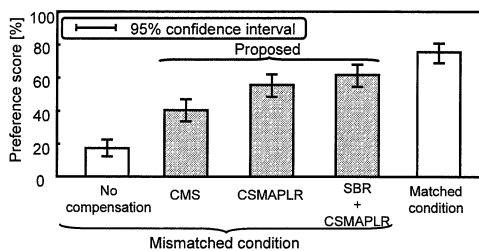


Fig. 10 Results of preference test on speech quality. Ten sentences were used for adaptation data in the CSMAPLR compensation.

特にその効果が顕著に見られることが分かった。

謝辞

本研究の一部は、総務省の SCOPE および文部科学省の科研費若手 (A) により実施したものである。

参考文献

- 1) Y. Zheng, Z. Liu, M. Sinclair, J. Droppo, L. Deng, A. Acero, and X. Huang, "Air- and Bone-Conductive Integrated Microphones for Robust Speech Detection and Enhancement", *Proc. ASRU*, pp. 249–254, 2003
- 2) Y. Nakajima, H. Kashioka, N. Campbell, and K. Shikano, "Non-Audible Murmur Recognition", *IE-ICE Trans. Information and Systems*, Vol. E89–D, No. 1, pp. 1–8, 2006
- 3) T. Toda, A. W. Black, and K. Tokuda, "Voice Conversion Based on Maximum Likelihood Estimation of Spectral Parameter Trajectory", *Trans. ASLP*, Vol. 15, No. 8, pp. 2222–2235, 2007
- 4) T. Toda, and K. Shikano, "NAM-to-Speech Conversion with Gaussian Mixture Models", *Proc. INTER-SPEECH*, pp. 1957–1960, 2005
- 5) 宮本 大輔, 中村 圭吾, 戸田 智基, 猿渡 洋, 鹿野 清宏, "肉伝導音声変換のための音響特性補正法", 信学技報, Jan. 2009
- 6) B. S. Atal, "Effectiveness of Linear Prediction Characteristics of the Speech Wave for Automatic Speaker Identification and Verification", *J. Acoust. Soc. America*, Vol. 55, No. 6, pp. 1304–1312, 1974
- 7) Y. Nakano, M. Tachibana, J. Yamagishi, and T. Kobayashi, "Constrained Structural Maximum A Posteriori Linear Regression for Average-Voice-Based Speech Synthesis", *Proc. INTERSPEECH*, pp. 2286–2289, 2006
- 8) M. Rahim, and B.H. Juang, "Signal bias removal by maximum likelihood estimation for robust telephone speech recognition", *IEEE Trans. Speech Audio Processing*, Vol. 4, pp. 19–30, 1996
- 9) M. J. F. Gales, "Maximum Likelihood Linear Transformation for HMM-based Speech recognition", *Computer Speech and Language*, Vol. 12, No. 2, pp. 75–98, 1998