

Blind Stochastic Feature Transformation for Speaker Verification over Cellular Networks

Kwok-Kwong Yiu, Man-Wai Mak, Ming-Cheung Cheung

Sun-Yuan Kung

Center for Multimedia Signal Processing
Dept. of Electronic and Information Engineering
The Hong Kong Polytechnic University, China

Dept. of Electrical Engineering
Princeton University
USA

Abstract

Acoustic mismatch between the training and recognition conditions presents one of the serious challenges faced by speaker recognition researchers today. The goal of channel compensation is to achieve performance approaching that of a “matched condition” system while avoiding the need for a large amount of training data. It is important to ensure that the channel compensation algorithms in these systems compensate the channel variation instead of speaker variation. This paper addresses the problem of unsupervised compensation in which the features of a test utterance are transformed to fit the clean speaker model and gender-dependent background model. Specifically, a feature-based transformation is estimated based on the statistical difference between a test utterance and a composite acoustic model formed by combining the speaker and background models. By transforming the features to fit both models, the transformation is implicitly constrained. Experimental results based on the 2001 NIST evaluation set show that the proposed transformation approach achieves significant improvement in both equal error rate and minimum detection cost as compared to cepstral mean subtraction, Znorm and short-time Gaussianization.

Keywords: Speaker verification; feature transformation; Channel compensation; robustness; MAP adaptation

1. Introduction

The accuracy of speaker recognition systems that enroll client speakers under one acoustic environment (e.g. using a close-talk microphone in offices) but verify claimants under another environment (e.g. using mobile phones on the street) could be significantly lower than the ones that enroll and verify speakers under the same environment. This is mainly due to the acoustic mismatch between the training and recognition conditions, which presents one of the serious challenges faced by speaker recognition researchers today. One cause of the mismatched conditions is transducer variability. Transducer variability occurs when a system is trained with speech data obtained from one type of transducer and is subsequently tested on speech data recorded from other types of transducers. The goal of channel compensation is to achieve performance approaching that of a “matched condition” system while avoiding the need for a large amount of training data.

Channel compensation can be applied in feature space, model space or score space. Feature-based compensation [1], [2] transforms channel-distorted speech features to fit clean

speaker models, whereas model-based compensation [3], [4] adapts or transforms the parameters of clean models to fit a new acoustic environment. On the other hand, score-based compensation [5] aims to minimize environment-dependent bias by normalizing the distribution of speaker scores.

Channel compensation can also be supervised or unsupervised. Supervised compensation assumes that the channel or handset characteristics are known a priori. Therefore, channel-specific compensation can be derived before recognition takes place. If handset labels are available during recognition, the corresponding channel-specific compensation can be applied to reduce the mismatch effect. Alternatively, one can detect the handset label from the speech signal during verification [2]. However, this approach may not be practical because users may use a new handset, which is not well represented in the training set, during verification. While this problem can be partially resolved by using a handset classifier with out-of-handset rejection capability [6, 7], it is difficult to find a threshold for detecting unseen handsets. On the other hand, unsupervised compensation does not assume any knowledge of the channel characteristics. In particular, it adapts speaker models or transforms speaker features to accommodate the channel variation based on verification utterances only. Therefore, handset detectors are no longer required.

In speaker verification, it is important to ensure that channel variations are suppressed so that the interspeaker distinction can be enhanced. In particular, given a claimant’s utterance recorded in an environment different from that during enrollment, one aims to transform the features of the utterance so that they become compatible with the enrollment environment. Therefore, it is not appropriate to transform the claimant’s utterance either to fit the speaker model only or to fit the background model only because the former will result in an unacceptably high FAR (false acceptance rate) and the latter an excessive FRR (false rejection rate). This paper describes a feature-based *blind* transformation approach to solving this problem. Specifically, a feature-based transformation is estimated based on the statistical difference between a test utterance and a composite acoustic model formed by combining the speaker and background models. The transformation is then used to transform the test utterance before verification. The transformation is blind in that it compensates the handset distortion without a priori information about the handset’s characteristics. Hereafter, this transformation approach is referred to as blind stochastic feature transformation (BSFT).

This work was supported by The Hong Kong Polytechnic University Grant No. G-W076 and G-T860.

2. Blind Stochastic Feature Transformation

Figure 1 illustrates a speaker verification system with BSFT, whose operations are divided into two separate phases: enrollment and verification.

1. *Environment Phase.* The speech of all client speakers are used to create a compact universal background model (UBM) Λ_b^M with M components. Then, for each client speaker, a compact speaker model Λ_s^M is created by adapting the UBM Λ_b^M using maximum a posteriori (MAP) adaptation [8]. Because verification decisions are based on the likelihood of the speaker model and background model, both models must be considered when the transformation parameters are computed. This can be achieved by fusing Λ_b^M and Λ_s^M to form a $2M$ -component composite GMM Λ_c^{2M} . During the fusion process, the means and covariances remain unchanged but the value of each mixing coefficients is divided by 2. This step ensures that the output of the composite GMM represents a probability density function.

2. *Verification Phase.* Distorted features $Y = \{\mathbf{y}_1, \dots, \mathbf{y}_T\}$ extracted from a verification utterance are used to compute the transformation parameters $\nu = \{A, \mathbf{b}\}$. This is achieved by maximizing the likelihood of the composite GMM Λ_c^{2M} given the transformed features $\hat{X} = \{\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_T\}$:

$$\hat{\mathbf{x}}_t = f_\nu(\mathbf{y}_t) = A\mathbf{y}_t + \mathbf{b}, \quad t = 1, \dots, T \quad (1)$$

where A is a $D \times D$ identity matrix for zeroth-order transformation and $A = \text{diag}\{a_1, a_2, \dots, a_D\}$ for first-order transformation, and \mathbf{b} is a bias vector. The transformed vectors \hat{X} are then fed to a full size speaker model Λ_s^N and a full size UBM Λ_b^N for computing verification scores in terms of likelihood ratio:

$$s(\hat{X}) = \log p(\hat{X} | \Lambda_s^N) - \log p(\hat{X} | \Lambda_b^N).$$

The main idea of BSFT is to transform the distorted features to fit the composite GMM Λ_c^{2M} , which ensures that the transformation compensates the acoustic distortion.

As the computation complexity of estimating SFT parameters grows with the amount of adaptation data and the total number of mixture components in the GMMs, the BSFT will become computationally intensive when the number of components is large. To perform rapid adaptation, we propose adopting a light-weight approach to computing transformation parameters. One of the positive properties of SFT is that the transformation can be estimated using GMMs with only a few components. In the light-weight approach, we synthesize a small, composite GMM (Λ_c^{2M}) from another small speaker GMM (Λ_s^M) and background GMM (Λ_b^M), both with M components where $M \ll N$. Similarly, the testing GMM (Λ_t^{2M}) was adapted from another UBM with $2M$ components. It was found that a good trade-off between performance and computation complexity can be maintained by using a suitable value of M .

Figure 2 illustrates the idea of BSFT in a classification problem with two-dimensional input patterns. Figure 2(a) plots the clean and distorted patterns of Class 1 and Class 2. The upper right (respectively lower left) clusters represent the clean (respectively distorted) patterns. The ellipses show the corresponding equal density contours. Markers ‘◆’ and ‘■’ represent the centers of the clean models. Figure 2(b) illustrates a

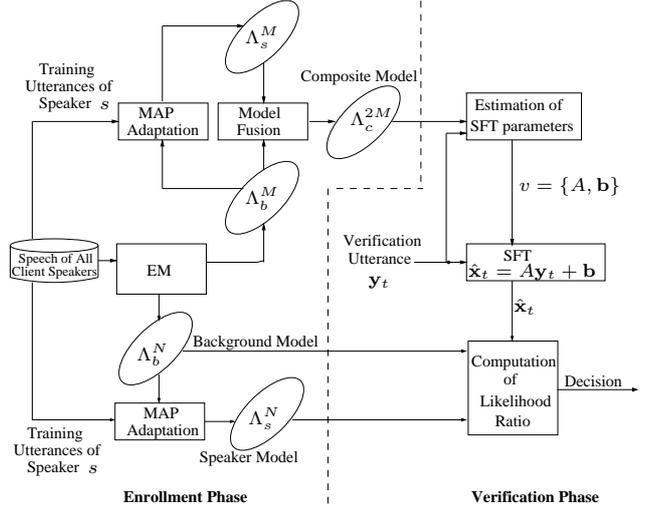


Figure 1: Estimation of BSFT parameters. The background model Λ_b^N , speaker model Λ_s^N , and composite model Λ_c^{2M} , produced during the enrollment phase, are subsequently used for verification purposes.

transformation matching the distorted data of Class 2 and the GMM of Class 1 (GMM1). Because the transformation only takes GMM1 into account, while ignoring GMM2 completely, it results in a high error rate. Similarly, the transformation in Figure 2(c) also has a high error rate. The transformation in Figure 2(d) was estimated from the distorted data of Class 1 and a composite GMM formed by fusing GMM1 and GMM2. In this case, the transformation adapts the data to a region close to both GMM1 and GMM2 because it takes both GMMs into account. Therefore, instead of transforming the distorted data to a region around GMM1 or GMM2 as in Figures 2(b) and 2(c), the transformation in Figure 2(d) attempts to compensate the distortion. The capability of BSFT is also demonstrated in a speaker verification task to be described next.

3. Experiments and Results

3.1. Experiments

Per the discussion earlier, the experiments were divided into two phases: enrollment and verification.

1. *Enrollment Phase.* A 1024-component UBM Λ_b^{1024} (i.e., $N = 1024$ in Figure 1) was trained using the training utterances of all target speakers. The same set of data was also used to train an M -component UBM (Λ_b^M in Figure 1). For each target speaker, a 1024-component speaker-dependent GMM Λ_s^{1024} was created by adapting Λ_b^{1024} using MAP adaptation [8]. Similarly, Λ_s^M was created by adapting Λ_b^M , and the two models are fused to form a composite GMM Λ_c^{2M} . The value of M was varied from 2 to 64 in the experiments.
2. *Verification Phase.* For each verification session, a feature sequence Y was extracted from the utterance of a claimant. The sequence was used to determine the BSFT parameters (A and \mathbf{b} in Eq. 1) to obtain a sequence of transformed vectors \hat{X} . The transformed vectors were then fed to Λ_s^{1024} and Λ_b^{1024} to obtain verifica-

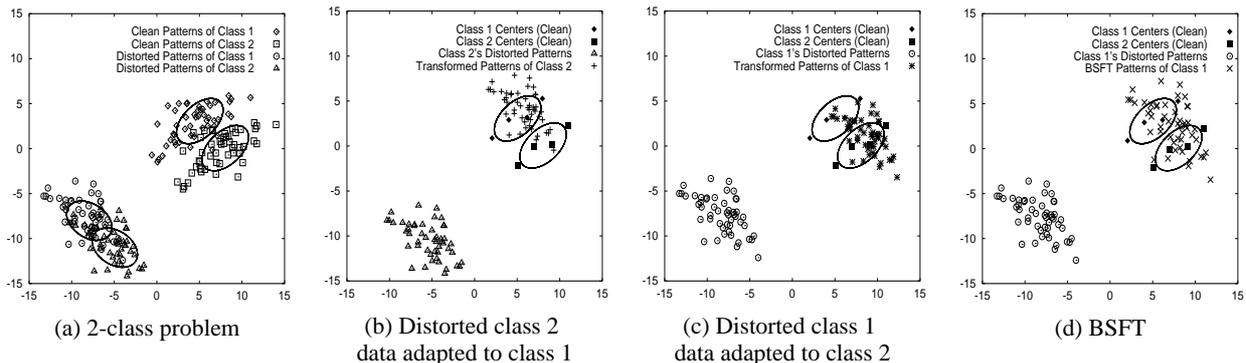


Figure 2: A 2-class problem illustrating the idea of BSFT. (a) Scatter plots of the clean and distorted patterns corresponding to Class 1 and Class 2. The thick and thin ellipses represent the equal density contours of Class 1 and Class 2, respectively. The upper right (respectively lower left) clusters contain the clean (respectively distorted) patterns. (b) Distorted patterns of Class 2 were transformed to fit Class 1’s clean model. (c) Reversely, distorted patterns of Class 1 were transformed to fit Class 2’s clean model. (d) Distorted data of Class 1 were transformed to fit the clean models of both Class 1 and Class 2 using first-order BSFT. For clarity, only the distorted patterns before and after transformation were plotted in (b)–(d).

tion scores for decision making.

The 2001 NIST speaker recognition evaluation set [9], which contains cellular phone speech of 74 male and 100 female target speakers extracted from the SwitchBoard-II Phase IV Corpus, was used in the evaluation. Each target speaker has 2 minutes of speech for training (i.e., enrollment); a total of 850 male and 1188 female utterances are available for testing (i.e., verification). Each verification utterance has length between 15 and 45 seconds and is evaluated against 11 hypothesized speakers of the same sex as the speaker of the verification utterance. Out of these 11 hypothesized speakers, one is the target speaker who produced the verification utterance. Therefore, there are one target and 10 impostor trials for each verification utterance, which amount to a total of 2,038 target trials and 20,380 impostor attempts for 2,038 verification utterances.

Mel-frequency cepstral coefficients (MFCCs) [10] and their first-order derivatives were computed every 14ms using a Hamming window of 28ms. Cepstral mean subtraction (CMS) [11] was applied to the MFCCs to remove linear channel effects. The MFCCs and delta MFCCs were concatenated to form 24-dimensional feature vectors.

Detection error tradeoff (DET) curves and equal error rates (EERs) were used as performance measures. They were obtained by pooling all scores of both sex from the speaker and impostor trials. In addition to DET curves and EERs, decision cost function (DCF) was also used as performance measure. The DCF is defined as

$$\begin{aligned} \text{DCF} &= C_{\text{Miss}} \times P_{\text{Miss}|\text{Target}} \times P_{\text{Target}} \\ &+ C_{\text{FalseAlarm}} \times P_{\text{FalseAlarm}|\text{NonTarget}} \times P_{\text{NonTarget}} \end{aligned}$$

where P_{Target} and $P_{\text{NonTarget}}$ are the prior probability of target and impostor speakers, respectively, and where C_{Miss} and $C_{\text{FalseAlarm}}$ are the costs of miss and false alarm errors, respectively. Following NIST’s recommendation [12], these parameters were set as follows: $P_{\text{Target}} = 0.01$, $P_{\text{NonTarget}} = 0.99$, $C_{\text{Miss}} = 10$ and $C_{\text{FalseAlarm}} = 1$.

3.2. Results

Because the evaluation trials in NIST01 are gender-matched, gender-dependent background models can be used for enroll-

ment and estimation of BSFT parameters. In the experiment, gender-dependent background models and MAP adaptation were used to create gender-dependent speaker models. A compact gender-dependent background model (with 64 components) was used to estimate the BSFT parameters. Figure 3 and Table 1 show the results of different approaches to channel compensation, including cepstral mean subtraction (CMS), Znrm [5] and the first-order blind stochastic feature transformation (BSFT) with different numbers of components (M). Evidently, all cases of BSFT show significant reduction in error rates when compared to CMS. In particular, first-order gender-dependent BSFT with 64 components achieves the largest error reduction. The DET curves also show that the BSFT is superior to Znrm at all operating points.

It is of interest to compare BSFT with the short-time Gaussianization approach proposed in Xiang et al. [13] because both methods transform distorted features in the feature space and their transformation parameters are estimated by the EM algorithm. In BSFT, a set of transformation parameters ν is computed by the EM algorithm in which the likelihood function of a composite GMM given the transformed test data is maximized. In short-time Gaussianization, a linear, global transformation matrix is estimated by the EM algorithm using the training data from all background speakers. The global transformation aims to decorrelate the features in the new feature space; it is applied to the distorted features before they are mapped to fit a normal distribution. The linearly transformed features are divided into a number of overlapping segments, with each segment containing a number of consecutive transformed vectors. The consecutive vectors in a segment is then sorted in ascending order. The rank of the central frame is used to find a warped feature so that its cumulative density function (CDF) matches the CDF of a standard normal distribution.

The short-time Gaussianization achieves an EER of 10.84% in the NIST 2001 evaluation set [13], whereas BSFT achieves an EER of 9.26%, which represent an error reduction of 14.58%.¹ The minimum decision cost of BSFT is also lower than that of short-time Gaussianization (0.0384 versus 0.0440).

It can be argued that the inferior performance of Gaussian-

¹Because Xiang et al. did not use Znrm in [13], their results should be compared with the one without Znrm here.

Adaptation Method	M	Equal Error Rate(%)	Minimum Detection Cost
Baseline	NA	11.44	0.0477
BSFT	2	11.29	0.0445
BSFT	4	10.27	0.0425
BSFT	8	9.77	0.0409
BSFT	16	9.48	0.0394
BSFT	32	9.38	0.0395
BSFT	64	9.26	0.0384
Znorm	NA	10.61	0.0427
BSFT+Znorm	64	8.36	0.0355

Table 1: Equal error rates (in %) and minimum decision cost achieved by cepstral mean subtraction (CMS), Znorm, and first-order blind stochastic feature transformation (BSFT) with different order and numbers of components (M) in the compact UBM Λ_b^M . Note that the number of components in the full-size speaker models and gender-dependent background models is 1024.

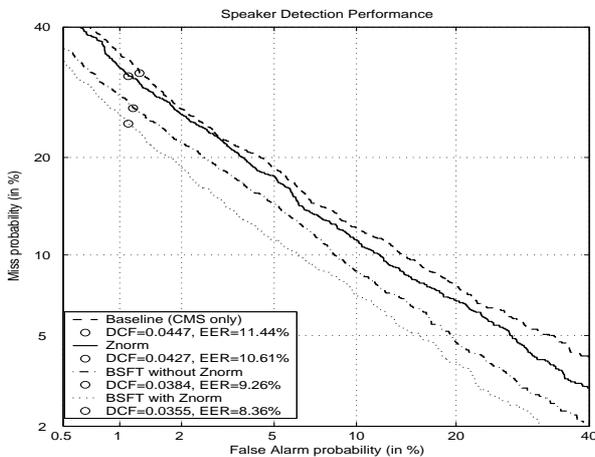


Figure 3: DET curves comparing speaker verification performance using CMS (dashed), Znorm (solid), BSFT without Znorm (dash-dot), and BSFT with Znorm (dotted). For BSFT, the number of components (M) in Λ_b^M was set to 64. The circles represent the errors at which minimum decision costs occur.

ization is due to the nonadaptive nature of its transformation parameters. However, the adaptive nature of BSFT comes with a computational price: different transformation parameters have to be computed for each speaker. Therefore, it is vital to have a cost effective computation approach for BSFT. Note that the computation complexity of estimating BSFT parameters grows with the amount of adaptation data (i.e., the value of T in Eq. 1) and the number of mixture components in the GMMs (i.e., the value of M). To reduce computation time, M should be significantly smaller than N , the number of components in the full size speaker and background models. This is particularly important for the computation of BSFT parameters during the verification phase because the computation time of this phase is a significant part of the overall verification time. The evaluations reported in this paper suggest that a good tradeoff between performance and computation complexity can be achieved by using a suitable value of M .

4. Conclusions

We have presented a new approach, namely blind stochastic feature transformation, to channel robust speaker verification and provided experimental results on the 2001 NIST evaluation set. The algorithm computes feature transformation parameters based on the statistical difference between a test utterance and a composite GMM formed by combining the speaker and background models. The transformation is then used to transform the test utterance to fit the clean speaker model and background model before verification. Experimental results show that the proposed algorithms achieves significant improvement in both equal error rate and minimum detection cost when compared to cepstral mean subtraction, Znorm and short-time Gaussianization.

5. References

- [1] A. C. Surendran, C. H. Lee, and M. Rahim, "Nonlinear compensation for stochastic matching," *IEEE Trans. on Speech and Audio Processing*, vol. 7, no. 6, pp. 643–655, 1999.
- [2] M. W. Mak and S. Y. Kung, "Combining stochastic feature transformation and handset identification for telephone-based speaker verification," in *Proc. ICASSP'02*, 2002, pp. 1701–1704.
- [3] F. Beaufays and M. Weintraub, "Model transformation for robust speaker recognition from telephone data," in *ICASSP-97*, 1997, vol. 2, pp. 1063–1066.
- [4] K. K. Yiu, M. W. Mak, and S. Y. Kung, "Environment adaptation for robust speaker verification," in *Eurospeech'03*, 2003, pp. 2973–2976.
- [5] D. A. Reynolds, "Comparison of background normalization methods for text independent speaker verification," in *Eurospeech'97*, 1997, pp. 963–966.
- [6] C. L. Tsang, M. W. Mak, and S. Y. Kung, "Divergence-based out-of-class rejection for telephone handset identification," in *Proc. Int. Conf. on Spoken Language Processing*, 2002, pp. 2329–2332.
- [7] M. W. Mak, C. L. Tsang, and S. Y. Kung, "Stochastic feature transformation with divergence-based out-of-handset rejection for robust speaker verification," *J. on Applied Signal Processing*, (to appear).
- [8] D. A. Reynolds, T. F. Quatieri, and R. B. Dunn, "Speaker verification using adapted Gaussian mixture models," *Digital Signal Processing*, vol. 10, pp. 19–41, 2000.
- [9] "The NIST year 2001 speaker recognition evaluation plan," in <http://www.nist.gov/speech/tests/spk2001/doc>.
- [10] S. B. Davis and P. Mermelstein, "Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences," *IEEE Trans. on ASSP*, vol. 28, no. 4, pp. 357–366, August 1980.
- [11] S. Furui, "Cepstral analysis technique for automatic speaker verification," *IEEE Trans. on Acoustic Speech and Signal Processing*, vol. ASSP-29, no. 2, pp. 254–272, 1981.
- [12] M. Przybocki A. Martin, "NIST's assessment of text independent speaker recognition performance 2002," in *The Advent of Biometrics on the Internet, A COST 275 Workshop*, Rome, Italy, Nov. 2002.
- [13] B. Xiang, U. Chaudhari, J. Navratil, G. Ramaswamy, and R. Gopinath, "Short-time Gaussianization for robust speaker verification," in *Proc. IEEE ICASSP02*, 2002, vol. 1, pp. 681–684.