

ARTICULATORY-FEATURE BASED SEQUENCE KERNEL FOR HIGH-LEVEL SPEAKER VERIFICATION

Shi-Xiong Zhang and Man-Wai Mak

Dept. of Electronic and Information Engineering,
The Hong Kong Polytechnic University
{ zhang . sx , enmw mak } @ polyu . edu . hk

ABSTRACT

Research has shown that articulatory feature-based phonetic-class pronunciation models (AFCPMs) can capture the pronunciation characteristics of speakers. However, the scoring method used in AFCPMs does not explicitly use the discriminative information available in the training data. To harness this information, this paper proposes converting speaker models to supervectors by stacking the discrete densities in AFCPMs. An AF-kernel is constructed from the supervectors of target speakers, background speakers, and claimants. An AF-kernel based SVM is then trained to classify the supervectors. Results show that AF-kernel scoring is complementary to likelihood-ratio scoring, leading to better performance when the two scoring methods are combined.

Index Terms— Speaker verification, kernels, articulatory features, pronunciation models, SVM

1. INTRODUCTION

Speaker verification is a binary classification problem in which a speaker is authenticated based on his/her own voice. Text-independent speaker verification systems typically extract speaker features from short-term spectra of speech signals to build speaker-dependent Gaussian mixture models (GMMs) [1]. Studies have shown that combining low-level acoustic information with high-level speaker information—such as the usage or duration of particular words, prosodic features and articulatory features (AF)—can improve speaker verification performance [2–6].

However, in most systems (e.g., GMM-UBM [1], PD-AFCPM [5] and CD-AFCPM [6]), scoring is done at the frame-level, i.e., each frame of speech is scored separately and then frame-based scores are accumulated to produce an utterance-based score for classification. This frame-based scoring scheme has two drawbacks. First, treating the frames individually may not be able to fully capture the sequence information contained in the utterance. Second, the goal of speaker verification is to minimize classification errors on test utterances

rather than on individual speech frames. These drawbacks motivate us to derive a sequence-based approach in which an utterance is considered comprising a sequence of symbols and the utterance-based score can be obtained from an SVM [7] through a kernel function of the sequence of symbols.

Support vector machines (SVMs) can produce complex decision functions without a large amount of training data. However, ordinary SVMs can only classify data of fixed dimensionality whereas speech utterances are typically parameterized as variable-length sequences. This leads to the idea of GMM supervectors in which variable-length observation sequences are mapped to fixed-dimensional vectors via stacking the parameters or scores of a GMM [8, 9].

This paper derives an articulatory-feature based sequence kernel and apply it to high-level speaker verification. The method begins with extracting the observation sequences (AF labels) from the utterances of target speakers. The AF sequences are then used to train the articulatory feature-based models (called AFCPM and CD-AFCPM) as proposed in [5, 6]. These models are then converted to fixed-dimensional AF supervectors. Redundant or irrelevant features are then removed from the AF supervectors. The reduced-dimension supervectors derived from both the target speaker and background speakers are used to train an SVM with a specially designed sequence kernel. The speaker-dependent SVM will be used to compute the verification scores. Since the kernel depends on the AF-based target speaker models, we refer it to as articulatory feature (AF) kernel. The remainder of the paper will derive the AF kernel and discuss the relationship between traditional frame-based log-likelihood (LR) scoring and AF-kernel based SVM scoring. Experimental results on the NIST2000 database are presented.

2. ARTICULATORY-FEATURE BASED SUPERVECTOR EXTRACTION

2.1. Articulatory Features

Articulatory features (AFs) are representations describing the movements or positions of different articulators during speech production. In [5, 6], the manner and place properties shown in Table 1 were used for pronunciation modeling. AFs can

This work was supported by the Research Grant Council of the Hong Kong SAR Project No. PolyU5230/05E and HKPolyU Project No. A-PA6F.

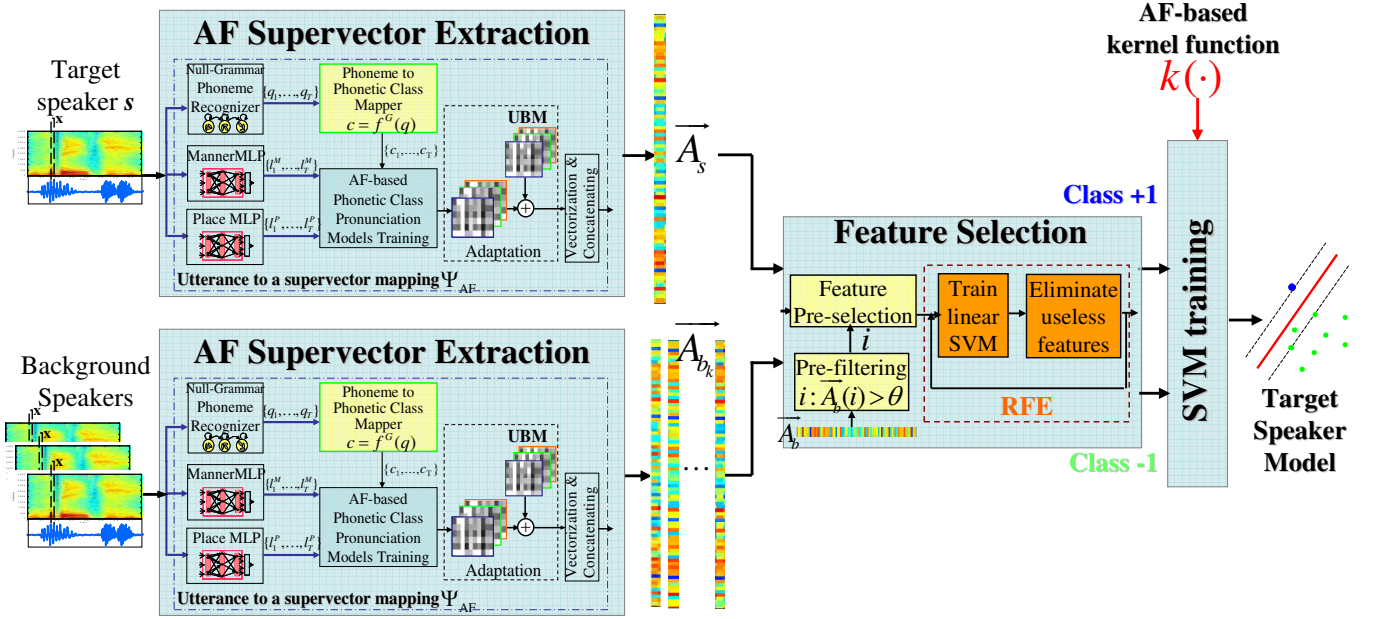


Fig. 1. The training procedure of the AF kernel-based high-level speaker verification system.

Articulatory Properties	Classes	Number of Classes
Manner(\mathcal{M})	Silence, Vowel, Stop, Fricative, Nasal, Approximant-Lateral	6
Place(\mathcal{P})	Silence, High, Middle, Low, Labial, Dental, Coronal, Palatal, Velar, Glottal	10

Table 1. The manner and place properties.

be automatically determined from speech signals using AF-based multilayer perceptrons (MLPs) shown in Fig. 1.

2.2. AF-Based Suprectors

Fig. 1 shows the training phase of a high-level speaker verification system that uses AF-kernels. The first step is to create CD-AFPCM suprectors. In [6], G phonetic-class dependent AFPCMs were trained from the AF and phoneme streams of all background speakers to represent the speaker-independent pronunciation characteristics. Each CD-AFPCM comprises the joint probabilities of the manner and place classes conditioned on a phonetic class. The training procedure begins with aligning two AF streams (l_t^m and l_t^p) obtained from the AF-MLPs and a phonetic class sequence c_t obtained from a null-grammar recognizer and a mapping function. The joint probabilities of background models $P_b^{\text{CD}}(m, p|c)$ corresponding to a particular phonetic class c is given by

$$\begin{aligned}
 P_b^{\text{CD}}(m, p|c) &= P_b^{\text{CD}}(L^{\text{M}} = m, L^{\text{P}} = p | \text{Phonetic Class} = c, \text{background}) \\
 &= \frac{\#((m, p, c) \text{ in the data of all background speakers})}{\#(*, *, c) \text{ in the data of all background speakers}}, \quad (1)
 \end{aligned}$$

where $m \in \mathcal{M}, p \in \mathcal{P}, (m, p, q)$ denotes the condition for which $L^{\text{M}} = m, L^{\text{P}} = p$, and Phonetic Class = c , * represents all possible members in that class, and $\#()$ represents the total number of frames with phoneme labels and AF labels fulfill the description inside the parentheses. We can see for each phonetic class, a total of 60 probabilities can be obtained, and a collection of G CD-AFPCM forms a AF-based universal background model (UBM).

The unadapted speaker model $P_s^{\text{CD}}(m, p|c)$ are created in the same way. And the MAP adaptation described in [10] is applied to create the adapted speaker models

$$\hat{P}_s^{\text{CD}}(m, p|c) = \beta_c P_s^{\text{CD}}(m, p|c) + (1 - \beta_c) P_b^{\text{CD}}(m, p|c), \quad (2)$$

where, $\beta_c \in [0, 1]$ is a phonetic class-dependent adaptation coefficient controlling the contribution of the speaker data and the background models (Eq. 1) on the MAP-adapted model.

For each target speaker, G speaker models (each model has 60 parameters) are stacked to form a single suprvector, called CD-AFPCM Suprvector.¹ The process maps a test utterance to a point in $60G$ -dimensional vector space. The procedure of CD-AFPCM suprvector extraction is illustrated in Fig. 1.

2.3. Feature Selection

The dimensionality of CD-AFPCM and PD-AFPCM is 720 (when $G = 12$) and 2760 (46 phonemes in English), respectively. Many of these features, however, may be redundant or having low discriminative power. Therefore, extracting the relevant features from the high-dimensional suprectors is expected to improve verification performance.

¹The procedure is also applicable to PD-AFPCM with $G = 46$. For clarity, we focus on CD-AFPCM in the sequel.

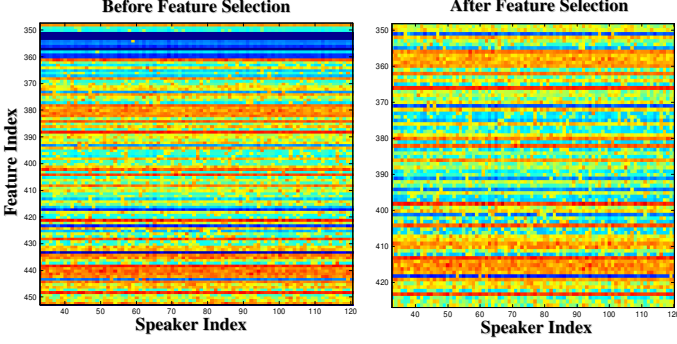


Fig. 2. Effect of feature selection on CD-AFCPM supervectors. A row with small variation (almost identical color intensity) suggests that the corresponding feature is not speaker dependent and therefore can be removed without sacrificing classification accuracy.

The feature selection process is divided into two steps. In Step 1, irrelevant features are weeded out using a pre-filtering approach. More precisely, the background CD-AFCPMs are vectorized and feature elements with value smaller than a threshold are removed. This step avoids the numerical difficulty that may encounter when the CD-AFCPMs supervectors are normalized during the evaluation of the kernel function (see Eqs. 12 and 13). Then, in Step 2, based on the remaining features in Step 1, a CD-AFCPM supervector is constructed for each target speaker, and 618 CD-AFCPM supervectors are constructed from 618 background speakers. Then, for each speaker, feature elimination is done by applying SVM-RFE [11] to the dataset formed by the speaker’s CD-AFCPM supervector (positive class) and background speakers’ CD-AFCPM supervectors (negative class). Note that Step 2 is applied to each of the target speakers, meaning that each speaker has their own feature set.

The effect of feature selection is shown in Fig. 2, where each column corresponds to a portion of a speaker’s CD-AFCPM supervector. A row with small variation (almost identical color intensity) suggests that the corresponding feature is not speaker dependent and therefore can be removed without sacrificing classification performance. We can see from Fig. 2 (right panel) that features of low discriminative power have been eliminated.

3. ARTICULATORY FEATURE-BASED KERNELS

3.1. Articulatory Feature-Based LR Scoring

Given a test utterance $X_1^T = \{X_1, \dots, X_t, \dots, X_T\}$, we can express the frame-based likelihood-ratio (LR) score as follows:

$$S_{\text{CD-AFCPM}}(X_1^T) = \frac{1}{T} \sum_{t=1}^T \left(\log \frac{\hat{P}_s^{\text{CD}}(l_t^M, l_t^P | c_t)}{\hat{P}_b^{\text{CD}}(l_t^M, l_t^P | c_t)} \right) \quad (3)$$

where $f^G(q_t)$ is one of the mapping function proposed in [6], c_t is the phonetic class determined by a null-grammar recog-

nizer, and $l_t^M \in \mathcal{M}$ and $l_t^P \in \mathcal{P}$ are the labels determined by the manner and place MLPs, respectively. Grouping all frames that are classified to the same phonetic class, we can further express the LR score as:

$$\begin{aligned} S_{\text{CD-AFCPM}}(X_1^T) &= \sum_{c=1}^G \frac{1}{T} \left(\sum_{\substack{m \in \mathcal{M} \\ p \in \mathcal{P}}} \left(\sum_{\substack{f^G(q_t)=c \\ l_t^M=m, \\ l_t^P=p}} \left(\log \frac{\hat{P}_s^{\text{CD}}(l_t^M = m, l_t^P = p | c)}{\hat{P}_b^{\text{CD}}(l_t^M = m, l_t^P = p | c)} \right) \right) \right) \\ &= \sum_{c=1}^G \frac{T_c}{T} \left\{ \frac{1}{T_c} \sum_{i=1}^{60} \left(\left(\log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_i | c)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_i | c)} \right) \sum_{t: \{f^G(q_t)=c\}} 1 \right) \right\} \\ &= \sum_{c=1}^G \frac{T_c}{T} \left\{ \frac{1}{T_c} \sum_{i=1}^{60} \left(\left(\log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_i | c)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_i | c)} \right) N_{i,c} \right) \right\} \\ &= \sum_{c=1}^G \frac{T_c}{T} \left\{ \sum_{i=1}^{60} \left(\left(\log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_i | c)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_i | c)} \right) \frac{N_{i,c}}{T_c} \right) \right\} \end{aligned} \quad (4)$$

where $\mathcal{L}_1 = \{l_t^M = \text{'Vowel'}, l_t^P = \text{'High'} \text{ for any } t\}, \dots, \mathcal{L}_{60} = \{l_t^M = \text{'Lateral'}, l_t^P = \text{'Glottal'} \text{ for any } t\}$, $N_{i,c}$ is the number of frames belonging to phonetic class c and \mathcal{L}_i , and T_c is the number of frames belonging to phonetic class c .

Assume that a test utterance is produced by a claimant cl claiming a speaker identity s . Then, we can obtain the CD-AFCPM of the claimant as follows:

$$\begin{aligned} P_{cl}^{\text{CD}}(m, p | c) &= P_{cl}^{\text{CD}}(L^M = m, L^P = p | \text{PhoneClass} = c, \text{claimant} = cl) \\ &= P_{cl}^{\text{CD}}(\mathcal{L}_i | c) \\ &= \frac{\#((m, p, c) \text{ in the utterances of the claimant})}{\#((*, *, c) \text{ in the utterances of the claimant})} \\ &= \frac{N_{i,c}}{T_c}, \end{aligned} \quad (5)$$

where index i corresponds to the i -th combination of the manner and place class (m, p) . Setting $G = 12$ and substituting Eq. 5 into Eq. 4, we obtain:

$$\begin{aligned} S_{\text{CD-AFCPM}}(X_1^T) &= \sum_{c=1}^G \frac{T_c}{T} \left(\sum_{i=1}^{60} \left(\left(\log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_i | c)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_i | c)} \right) P_{cl}^{\text{CD}}(\mathcal{L}_i | c) \right) \right) \\ &= \sum_{c=1}^G \left\langle \begin{bmatrix} \log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_1 | c)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_1 | c)} \\ \log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_2 | c)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_2 | c)} \\ \dots \\ \log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_{60} | c)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_{60} | c)} \end{bmatrix}_{60}, \begin{bmatrix} \frac{T_c}{T} P_{cl}^{\text{CD}}(\mathcal{L}_1 | c) \\ \frac{T_c}{T} P_{cl}^{\text{CD}}(\mathcal{L}_2 | c) \\ \dots \\ \frac{T_c}{T} P_{cl}^{\text{CD}}(\mathcal{L}_{60} | c) \end{bmatrix}_{60} \right\rangle \end{aligned}$$

$$= \left\langle \begin{bmatrix} \log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_1|c=1)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_1|c=1)} \\ \dots \\ \log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_{60}|c=1)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_{60}|c=1)} \\ \log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_1|c=2)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_1|c=2)} \\ \dots \\ \log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_{60}|c=2)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_{60}|c=2)} \\ \dots \\ \log \frac{\hat{P}_s^{\text{CD}}(\mathcal{L}_{60}|c=12)}{\hat{P}_b^{\text{CD}}(\mathcal{L}_{60}|c=12)} \end{bmatrix}_{720}, \begin{bmatrix} w_1 P_{cl}^{\text{CD}}(\mathcal{L}_1|c=1) \\ \dots \\ w_1 P_{cl}^{\text{CD}}(\mathcal{L}_{60}|c=1) \\ w_2 P_{cl}^{\text{CD}}(\mathcal{L}_1|c=2) \\ \dots \\ w_2 P_{cl}^{\text{CD}}(\mathcal{L}_{60}|c=2) \\ \dots \\ w_{12} P_{cl}^{\text{CD}}(\mathcal{L}_{60}|c=12) \end{bmatrix}_{720} \right\rangle,$$

where $w_i = \frac{T_i}{T}$ ($i = 1, \dots, 12$), which gives

$$S_{\text{CD-AFCPM}}(X_1^T) = \left\langle \log \frac{\vec{A}_s}{\vec{A}_b}, \vec{A}'_c \right\rangle = \left\langle \log \frac{\vec{A}_s}{\vec{A}_b}, \vec{w} * \vec{A}_c \right\rangle \quad (6)$$

where \vec{A}_s, \vec{A}_b and \vec{A}_c stand for the AF supervector of the

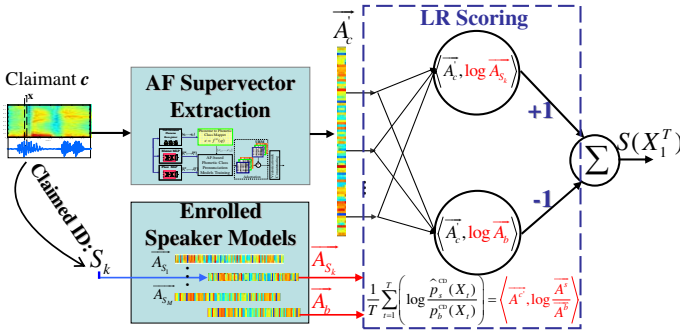


Fig. 3. An alternative implementation of the traditional log-likelihood scoring in CD-AFCPM speaker verification.

speaker, background and claimant, respectively, $\log \frac{\vec{X}}{\vec{Y}} \equiv \left[\log \frac{x_1}{y_1}, \dots, \log \frac{x_N}{y_N} \right]^T$ and $\vec{X} * \vec{Y} \equiv [x_1 y_1, \dots, x_N y_N]^T$, where x_i and y_i are elements of \vec{X} and \vec{Y} , respectively. We can see from Eq. 6 that the traditional frame-based LR scoring for discrete models can be computed in another way: compute the dot product between a speaker-dependent supervector derived from the models of the target speakers and a weighted supervector obtained from the claimant's model.

Further, denote \vec{A}'_c as the weighted AF supervector of the claimant, we have

$$\begin{aligned} S_{\text{CD-AFCPM}}(X_1^T) &= \left\langle \log \frac{\vec{A}_s}{\vec{A}_b}, \vec{A}'_c \right\rangle \\ &= \left\langle \vec{A}'_c, \log \frac{\vec{A}_s}{\vec{A}_b} \right\rangle = \left\langle \vec{A}'_c, \log \vec{A}_s \right\rangle - \left\langle \vec{A}'_c, \log \vec{A}_b \right\rangle. \end{aligned} \quad (7)$$

Eq. 7 suggests an alternative approach to implementing the traditional LR scoring. This is shown in Fig. 3. It shows that if we can replace the fixed '+1' and '-1' multiplication factors in the LR scoring block by varying weights, the result may probably be improved. These weights can be optimally determined by SVM training. In order to apply SVM and to make sure that the training algorithm converges to a stable solution, the function inside the cycle in Fig. 3 should satisfy the Mercer condition [7]. Unfortunately, $f(\vec{X}, \vec{Y}) = \left\langle \vec{X}, \log \vec{Y} \right\rangle$ does not satisfy the Mercer condition because $\left\langle \vec{X}, \log \vec{Y} \right\rangle$ cannot be written as $\left\langle \Phi(\vec{X}), \Phi(\vec{Y}) \right\rangle$. Therefore, we derive an AF kernel function that satisfies the Mercer condition in the next section.

3.2. Deriving Kernels from Similarity Scores

Given AF-based supervectors \vec{A}_s and \vec{A}_b obtained from the utterances of speaker s and background speakers into a fixed-dimension input space, the similarity score between the model deriving from the test utterance X_1^T of claimant c and the model of speaker s can be computed by a similarity (discriminant) function:

$$\text{Similarity}(\vec{A}_c, \vec{A}_s) = \left\langle \vec{A}'_c, \log \frac{\vec{A}_s}{\vec{A}_b} \right\rangle \quad (8)$$

where Eq. 7 has been used in the derivation.

Our goal is to make Eq. 8 symmetric and satisfy the Mercer condition. To this end, we expand $\log(x)$ at $x = 1$ as a Taylor series:

$$\begin{aligned} \log(x) &= \sum_{n=0}^{\infty} \frac{\log^{(n)}(1)}{n!} (x-1)^n \\ &= (x-1) - \frac{1}{2}(x-1)^2 + \mathcal{O}((x-1)^3). \end{aligned} \quad (9)$$

Because the speaker models are adapted from the UBMs, $\frac{\vec{A}_s}{\vec{A}_b} \rightarrow \vec{1}$. Therefore, we can ignore the high orders of $\left(\frac{\vec{A}_s}{\vec{A}_b} - \vec{1}\right)$ and approximate Eq. 8 as:

$$\begin{aligned} \text{Similarity}(\vec{A}_c, \vec{A}_s) &= \left\langle \vec{A}'_c, \log \frac{\vec{A}_s}{\vec{A}_b} \right\rangle \approx \left\langle \vec{A}'_c, \left(\frac{\vec{A}_s}{\vec{A}_b} - \vec{1} \right) \right\rangle \\ &= \left\langle \vec{A}'_c, \frac{\vec{A}_s}{\vec{A}_b} \right\rangle - \left\langle \vec{A}'_c, \vec{1} \right\rangle = \left\langle \vec{A}'_c, \frac{\vec{A}_s}{\vec{A}_b} \right\rangle - 1, \end{aligned} \quad (10)$$

where the last term is due to

$$\sum_{i=1}^{60G} A'_{c,i} = \sum_{k=1}^G \sum_{i=60(k-1)+1}^{60k} w_i A_{c,i} = \frac{1}{T} \sum_{k=1}^G T_k = 1.$$

Because the last term can be dropped without affecting the classification, we have

$$\begin{aligned} \text{Similarity}(\vec{A}_c, \vec{A}_s) &\approx \left\langle \vec{A}_c, \frac{\vec{A}_s}{\sqrt{A_b}} \right\rangle = \left\langle \frac{\vec{A}_c}{\sqrt{A_b}}, \frac{\vec{A}_s}{\sqrt{A_b}} \right\rangle \\ &= \left\langle \frac{\vec{w} \cdot \vec{A}_c}{\sqrt{A_b}}, \frac{\vec{A}_s}{\sqrt{A_b}} \right\rangle \approx \left\langle \frac{\sqrt{w_b} \cdot \vec{A}_c}{\sqrt{A_b}}, \frac{\sqrt{w_b} \cdot \vec{A}_s}{\sqrt{A_b}} \right\rangle \end{aligned} \quad (11)$$

$$\text{where } \vec{w}_b = \left[\overbrace{\frac{T_1^b}{T}, \dots, \frac{T_1^b}{T}}^{60}, \overbrace{\frac{T_2^b}{T}, \dots, \frac{T_2^b}{T}}^{60}, \dots, \overbrace{\frac{T_G^b}{T}, \dots, \frac{T_G^b}{T}}^{60} \right]^T$$

contains the phonetic-class weights obtained from the UBM, which is used to approximate \vec{w} in Eq. 6. The approximation aims to make the similarity measure symmetric.

Finally, we write the similarity function (Eq. 11) as a kernel:

$$K_{\text{AF}}(\vec{A}_c, \vec{A}_s) = \left\langle \frac{\sqrt{w_b} \cdot \vec{A}_c}{\sqrt{A_b}}, \frac{\sqrt{w_b} \cdot \vec{A}_s}{\sqrt{A_b}} \right\rangle = \langle \varphi(\vec{A}_c), \varphi(\vec{A}_s) \rangle \quad (12)$$

where the mapping $\varphi(\cdot)$ is defined as:

$$\varphi(\vec{X}) = \frac{\sqrt{w_b} \cdot \vec{X}}{\sqrt{A_b}}. \quad (13)$$

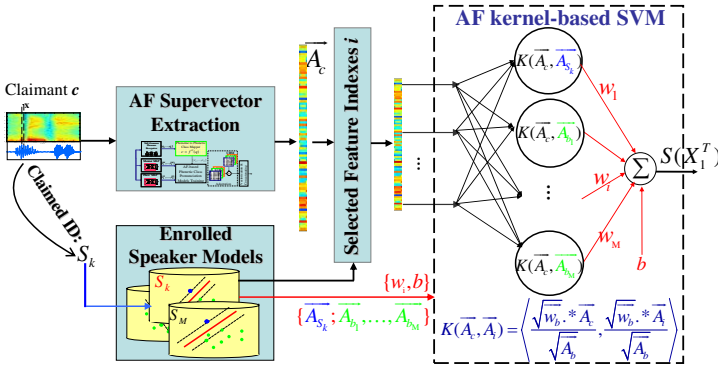


Fig. 4. The verification phase of an AF-kernel based speaker verification system.

Figures 5(a) and 5(b) show the un-normalized \vec{A}_s and the normalized \vec{A}_s (i.e., $\varphi(\vec{A}_s)$) for 150 speakers, respectively. For clarity, only feature elements with indexes between 531 and 660 are shown. Evidently, without normalization, some features have a large but almost constant value across all speakers (e.g., rows with dark-red color). These features will cause problems in SVM classification because they will dictate the decision boundary of the SVM, even though they contain little speaker-dependent information. This problem has

been largely alleviated by the normalization, as demonstrated in Fig. 5(b). In particular, the normalization has the effect of keeping all features within a comparable range, which helps prevent the large but almost constant features from dominating the classification decision. Figure 4 shows the scoring procedure during the verification phase. Comparing Fig. 4 and Fig. 3 suggests that scoring based on the AF-kernel is more general. The scores produced by the AF-kernel may also complement the ones produced by the LR-based method, which will be demonstrated in the next section.

Note that Eq. 12 depends on the speaker model, which in turn depends on the target speaker's utterances. Denote \vec{O}_s as the stacking of all acoustic vectors of speaker s . We can consider the training of speaker s 's CD-AFCPM as finding a function $\Psi_{\text{AF}}(\vec{O}_s)$ that maps the variable-dimension super-vector \vec{O}_s in \mathbb{R}^{TD} to an AF-suprvector \vec{A}_s of fixed-dimension in \mathbb{R}^{60G} , where T is the number of acoustic frames in the training utterances, D is the dimension of the acoustic vectors, and G is the number of phonetic classes. Therefore, we can express the AF kernel as:

$$\begin{aligned} \tilde{K}_{\text{AF}}(\vec{O}_c, \vec{O}_s) &= \left\langle \frac{\sqrt{w_b} \cdot \Psi_{\text{AF}}(\vec{O}_c)}{\sqrt{\Psi_{\text{AF}}(\vec{O}_b)}}, \frac{\sqrt{w_b} \cdot \Psi_{\text{AF}}(\vec{O}_s)}{\sqrt{\Psi_{\text{AF}}(\vec{O}_b)}} \right\rangle \\ &= \langle \Phi_{\text{AF}}(\vec{O}_c), \Phi_{\text{AF}}(\vec{O}_s) \rangle \end{aligned} \quad (14)$$

where the mapping $\Phi_{\text{AF}}(\cdot)$ is defined as:

$$\Phi_{\text{AF}}(\vec{O}) = \frac{\sqrt{w_b} \cdot \Psi_{\text{AF}}(\vec{O})}{\sqrt{\Psi_{\text{AF}}(\vec{O}_b)}}. \quad (15)$$

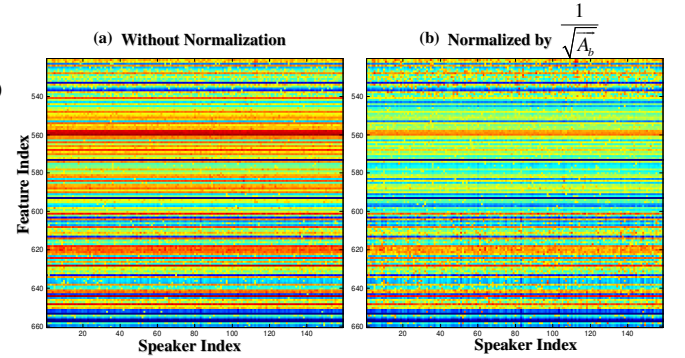


Fig. 5. The effect of the normalization term $\frac{1}{\sqrt{A_b}}$ in the mapping $\varphi(\cdot)$.

4. EXPERIMENTS AND RESULTS

4.1. Procedures

NIST99, NIST00, SPIDRE, and HTIMIT were used in the experiments. NIST99 was used for creating the background models and mapping functions, and the female part of NIST00 was used for creating speaker models and for performance evaluation. HTIMIT and SPIDRE were used for training the

AF-MLPs and the null-grammar phone recognizer, respectively. The phone recognizer uses standard 39- D vectors comprising MFCCs, energy, and their derivatives. The AF-MLPs use 38- D vectors comprising 19- D MFCCs and their first derivative computed every 10ms.

4.2. Results and Discussions

Kernel Function	CD-AFCPM Supervector	PD-AFCPM Supervector
$K(\bar{A}_c, \bar{A}_s) = \langle \bar{A}_c, \bar{A}_s \rangle$	26.12%	28.63%
$K(\bar{A}_c, \bar{A}_s) = \left\langle \frac{\sqrt{w_b} \cdot \bar{A}_c}{\sqrt{A_b}}, \frac{\sqrt{w_b} \cdot \bar{A}_s}{\sqrt{A_b}} \right\rangle$	24.14%	27.14%

Table 2. The EERs of AF kernel-based speaker verification systems using PD-AFCPM and CD-AFCPM supervectors without feature selection.

The results of using AF-kernels (without feature selection) for computing the verification scores are shown in Table 2. Evidently, normalization helps reduce the EER significantly. Similar to the results in LR-based scoring approach, CD-AFCPM is superior to PD-AFCPM under the AF-kernel framework. Comparing the result in [6] (EER = 23.46%) and the 2nd row of Table 2 (EER = 24.14%) suggests that without feature selection, the AF-kernel is inferior to the conventional LR-based scoring.

Kernel Function	No Feature Selection	Feature Selection
$K(\bar{A}_c, \bar{A}_s) = \left\langle \frac{\sqrt{w_b} \cdot \bar{A}_c}{\sqrt{A_b}}, \frac{\sqrt{w_b} \cdot \bar{A}_s}{\sqrt{A_b}} \right\rangle$	24.14%	23.90%

Table 3. EER achieved by the AF kernel-based speaker verification system using CD-AFCPM supervectors with and without feature selection.

We conjecture that the inferiority is caused by the irrelevant features in the supervectors. To verify this conjecture, we performed the same experiment but with the irrelevant features removed by SVM-RFE [11]. The results are shown in Table 3. Although feature selection can reduce the EER from 24.14% to 23.90%, it is still higher than that achieved by the LR-based method, which is 23.46%. From the experiment we can see that the EER can be reduced by reducing the dimensionality of the supervectors. To further reduce the EER, we fused the LR-based method and the AF-kernel scoring method. The EER is reduced to 23.02%, as shown in the DET plots in Fig. 6. The p -values [12] between the EERs of the fusion and non-fusion cases are all smaller than 0.00001, suggesting that the differences in EERs are statistically significant. This result illustrates that the AF-kernel based method and LR-based method are complementary to each other.

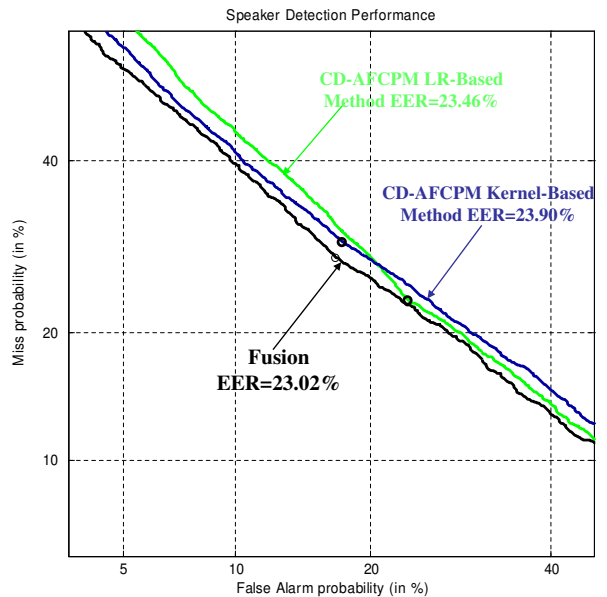


Fig. 6. DET produced by the AF-kernel scoring method with and without feature selection. CD-AFCPM supervectors were used in the AF-kernel.

5. REFERENCES

- [1] 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.
- [2] D. Reynolds, et al., "The superSID project: Exploiting high-level information for high-accuracy speaker recognition," in *Proc. International Conference on Audio, Speech, and Signal Processing*, Hong Kong, April 2003, vol. 4, pp. 784–787.
- [3] J. P. Campbell, D. A. Reynolds, and R. B. Dunn, "Fusing high- and low-level features for speaker recognition," in *Proc. Eurospeech*, 2003, pp. 2665–2668.
- [4] D. Klusacek, J. Navratil, D. A. Reynolds, and J. P. Campbell, "Conditional pronunciation modeling in speaker detection," in *Proc. ICASSP'03*, 2003, vol. 4, pp. 804–807.
- [5] K. Y. Leung, M. W. Mak, and S. Y. Kung, "Adaptive articulatory feature-based conditional pronunciation modeling for speaker verification," *Speech Communication*, vol. 48, no. 1, pp. 71–84, 2006.
- [6] S. X. Zhang, M. W. Mak, and Helen H. Meng, "Speaker verification via high-level feature based phonetic-class pronunciation modeling," *IEEE Trans. on Computers*, vol. 56, no. 9, pp. 1189–1198, 2007.
- [7] V. N. Vapnik, *The Nature of Statistical Learning Theory*, Springer-Verlag, New York, 1995.
- [8] W. M. Campbell, D. E. Sturim, and D. A. Reynolds, "Support vector machines using GMM supervectors for speaker verification," *IEEE Signal Processing Letters*, vol. 13, pp. 308–311, 2006, May.
- [9] V. Wan and S. Renals, "SVM-SVM: Support vector machine speaker verification methodology," in *Proc. ICASSP'03*, 2003, vol. II, pp. 221–224.
- [10] S.X. Zhang and M. W. Mak, "A new adaptation method for speaker-model creation in high-level speaker verification," in *Advances in Multimedia Information Processing (PCM'2007)*, Hong Kong, Springer LNCS 4810, 2007, pp. 325–335.
- [11] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene selection for cancer classification using support vector machines," *Machine Learning*, vol. 46, pp. 389–422, 2002.
- [12] L. Gillick and S. Cox, "Some statistical issues in the comparison of speech recognition algorithms," in *Proc. ICASSP'89*, 1989, pp. 532–535.