

行政院國家科學委員會專題研究計畫 成果報告

子計畫二:無線通訊環境下國語中文之聲學及語言處理基礎  
技術之研究(3/3)

計畫類別：整合型計畫

計畫編號：NSC91-2219-E-002-036-

執行期間：91年08月01日至92年07月31日

執行單位：國立臺灣大學資訊工程學系暨研究所

計畫主持人：李琳山

報告類型：完整報告

處理方式：本計畫可公開查詢

中華民國 92 年 8 月 28 日

行政院國家科學委員會專題研究計畫成果報告  
無線通訊環境下以語音擷取網路資訊相關技術之研究 (3/3)

> 子計畫二：無線通訊環境下國語中文之聲學及語言  
處理基礎技術之研究

計畫編號：NSC 91-2219-E-002-036

執行期限：91年8月1日至92年7月31日

主持人：李琳山 國立台灣大學資訊工程學系

E-mail: lslee@gate.sinica.edu.tw

## ABSTRACT

It was previously proposed to use the Principal Component Analysis (PCA) to derive the data-driven temporal filters for obtaining robust features in speech recognition, in which the first principal components are taken as the filter coefficients [1,2]. In this report, a multi-eigenvector approach is proposed instead, in which the first M eigenvectors obtained in PCA are weighted by their corresponding eigenvalues and summed to be used as the filter coefficients. Experimental results showed that the multi-eigenvector filters offer significant recognition performance as compared to the previously proposed PCA-derived filters under all different conditions tested with the AURORA2 database, especially when the training and testing environments are highly mismatched.

### 摘要

前人提出以主成份分析產生與語料相關的時域濾波器，以獲得較具強健性的語音辨識特徵參數；這種方法即是把第一個特性向量作為濾波器係數。在本報告中，我們提出了多特性向量法，也就是將由主成份分析所得到的前M個特性向量乘上對應的特性值後相加來作為濾波器係數。在AURORA 2語料的各種條件下，實驗結果皆顯示多特性向量濾波器比起前人提出的主成份分析濾波器有著顯著的進步，尤其是在訓練及測試環境嚴重不匹配時。

## 1. INTRODUCTION

Real applications of speech recognition strongly demand the recognition performance to be robust with respect to environmental changes. However, the recognition accuracy of almost all existing recognition systems drops dramatically when there is a mismatch between the training and testing

conditions. Substantial researches have been made in this area. One category of such approaches, among many others, is focused on finding a set of robust feature representation for signals, so that it is less sensitive to various environmental distortions. Cepstral Mean Subtraction (CMS), Cepstral Normalization (CN) [4], and Relative Spectral (RASTA) [5] techniques are typical examples of this category. Some of them can be considered as pre-filtering on the time trajectories of speech features in order to alleviate the harmful effects of various distortions and corruptions. Such approaches have been proved to be able to improve the robustness of the recognition performance significantly, and the subject of this report is also along this direction.

The filters used by CMS, CN and others are independent of the recognition environment. Although they are very effective, there is no guarantee that these solutions are optimal for a specific recognition application environment. The filtering coefficients optimized for a specific recognition task or application environment via data-driven approaches based on some optimization criteria can become highly desirable. The criterion of Linear Discriminant Analysis (LDA) has been widely applied [6,7] in such approaches in the optimization process. In recent works, the criteria of Linear Discriminant Analysis (LDA), Principal Component Analysis (PCA) and Minimum Classification Error (MCE) were also used as the optimization criteria to derive the data-driven temporal filters. It was shown [2] that all these data-driven temporal fi-

although derived from different criteria, can offer reasonable recognition performance improvements for recognition tasks with mismatched conditions.

In this report, PCA is again used as the optimization criterion for deriving temporal filters as well. However, different from the previous works [1, 2], here we take into consideration multi-eigenvectors rather than the first principal component only as used previously [1, 2], and the resulted temporal filters are the properly weighted linear combination of these multi-eigenvectors. It is shown that these new temporal filters can offer significant improvements over the previously obtained PCA-derived filters based on the first principal components only.

The remainder of this report has 4 sections. The approach to obtain the multi-eigenvector temporal filters is first described in section 2. The experimental setup and the experimental results are then presented and discussed in sections 3 and 4. Section 5 finally gives the concluding remarks.

## 2. TEMPORAL FILTER DESIGN USING PRINCIPAL COMPONENT ANALYSIS (PCA)

Given an ordered sequence of  $K$ -dimensional feature vectors  $\mathbf{x}(n)$  as shown in Figure 1(a), with time index  $n=1, \dots, N$ , and feature index  $k=1, 2, \dots, K$ ,  $\mathbf{x}(n)=[x(n, 1) \ x(n, 2), \dots, x(n, k), \dots, x(n, K)]^T$ ,  $n=1, \dots, N$  (1) then the  $k$ 'th time trajectory of  $\mathbf{x}(n)$  is the sequence  $[x(1, k) \ x(2, k) \ \dots \ x(N, k)]$ , denoted as  $y_k(n)$ , where  $y_k(n)=x(n, k)$ . Now

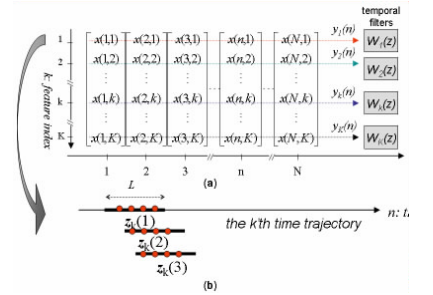
we'd like to design an  $L$ -sample FIR  $W_k(z)$  to be performed on the  $k$ 'th trajectory  $y_k(n)$ . First, an  $L$ -rectangular window is shifted along  $k$ 'th time trajectory to obtain sequences of  $L$ -dimensional vectors  $n=1 \dots N-L+1$ ,

$$\mathbf{z}_k(n)=[y_k(n) \ y_k(n+1) \ y_k(n+2) \ \dots \ y_k(n+L-1)]^T. \quad (2)$$

So  $\mathbf{z}_k(n)$  is the windowed vector of the time trajectory started at the time  $n$ , on which the  $L$ -sample FIR filter is applied, as depicted in Figure 1(b).

### 2.1 Previously proposed PCA-derived temporal filters [1]

In the previous work [1], the  $L$ -dimensional vectors,  $\mathbf{z}_k(n)$ , are as the samples of a random vector hence the mean vector and the covariance matrix of  $\mathbf{z}_k$  can be calculated,



**Figure 1.** The representation of trajectories of feature sequence the time index  $n$  and the feature  $i$  and (b) the windowed vector  $\mathbf{z}_k(n)$  along a time trajectory

$$\hat{\mu}_{z_i} = \frac{1}{N-L+1} \sum_{n=1}^{N-L+1} z_i(n)$$

$$\hat{\sigma}_{z_i}^2 = \frac{1}{N-L+1} \sum_{n=1}^{N-L+1} (z_i(n) - \hat{\mu}_{z_i})^2$$

Following the procedure of PCA, the components of the first eigenvector  $\tilde{\boldsymbol{o}}_k$  corresponding to the largest eigenvalue of the covariance matrix  $\Sigma_{z_k}$  is then taken as the coefficients of the  $L$ -sample filter, which maps the  $L$ -dimensional random vector  $\boldsymbol{z}_k$  into a one-dimensional random variable with maximum variance. Such process is carried out for each time trajectory, thus yielding a separate FIR filter for each time trajectory.

## 2.2 The multi-eigenvector temporal filtering approach

Based on the PCA theory [8], assume  $\tilde{\boldsymbol{o}}_{i,k}$ ,  $i=1, 2, \dots, L$  are the  $L$  distinct normalized eigenvectors of the covariance matrix  $\Sigma_{z_k}$  with corresponding decreasing eigenvalues  $\tilde{\epsilon}_{i,k}$ ,  $i=1, 2, \dots, L$ , i.e.,  $\tilde{\epsilon}_{1,k} \geq \tilde{\epsilon}_{2,k} \geq \dots \geq \tilde{\epsilon}_{L,k}$ , and  $y_{i,k}$ ,  $i=1, 2, \dots, L$  are the random variables representing the projections of the random vector  $\boldsymbol{z}_k$  on  $\tilde{\boldsymbol{o}}_{i,k}$ . It can then be shown that the variance of  $y_{i,k}$  is equal to  $\tilde{\epsilon}_{i,k}$ . As a result, with the previously proposed PCA-derived temporal filters, the filter coefficients used are the components of the eigenvector  $\tilde{\boldsymbol{o}}_{i,k}$ , therefore the filter output is the random variable  $y_{i,k}$  with variance being the largest eigenvalue  $\tilde{\epsilon}_{i,k}$ . Therefore  $y_{i,k}$  can be viewed as the most “expressive” 1-dimensional representation of  $\boldsymbol{z}_k$ , and this is apparently why the recognition can be improved.

However, from the above discussion it is clear that there is still some part of information carried by  $\boldsymbol{z}_k$  which was not used at all in the previously proposed PCA-derived filters, i.e., the other  $y_{i,k}$ 's,  $i=2, \dots, L$ , which also carry some

information of  $\boldsymbol{z}_k$  that may be help improving the recognition perfor. With these observations, the multi-eigenvector temporal filtering approach proposed in this report are

$$\bar{\boldsymbol{w}}_k = \sum_{i=1}^M \lambda_{i,k} \tilde{\boldsymbol{o}}_{i,k}, \quad \boldsymbol{w}_k = \frac{\bar{\boldsymbol{w}}_k}{\|\bar{\boldsymbol{w}}_k\|} = \frac{1}{\sqrt{\sum_{i=1}^M \lambda_{i,k}^2}} \bar{\boldsymbol{w}}_k,$$

where  $\boldsymbol{w}_k$  is the new coefficients  $L$ -sample filter for time trajectory summation is over the first  $M$  eigenvectors with larger corresponding eigenvalues and  $1 < M \leq L$ . Note that the length of the temporal filter vector  $\boldsymbol{w}_k$  is normalized to unity, so as to be consistent with the eigenvector used in the previously proposed PCA-derived temporal filters. Also, the first  $M$  eigenvectors are weighted by the corresponding eigenvalues in Eq. (5), therefore the outputs of the new temporal filter can be viewed as samples of a new random variable

$$v_k = \boldsymbol{w}_k^T \boldsymbol{z}_k = \frac{1}{\sqrt{\sum_{i=1}^M \lambda_{i,k}^2}} \sum_{i=1}^M \lambda_{i,k} y_{i,k}.$$

In other words, in addition to  $y_{1,k}$  used previously, with the proposed multi-eigenvector filters the other components of  $\boldsymbol{z}_k$  are also included here, while weighted by their corresponding variance values. The parameter  $M$ , or the number of eigenvectors used here, can be determined empirically. In the experiments presented below,  $M$  is chosen to be 3. Such a choice can be verified by the fact that  $\tilde{\epsilon}_{1,k}$ ,  $\tilde{\epsilon}_{2,k}$ , and  $\tilde{\epsilon}_{3,k}$  are always significantly larger than the other eigenvalues.

In the experiments below, it will

be shown that the new multi-eigenvector filters obtained with the approach here are low-pass filters whose characteristics in speech signals is able to enhance the syllabic-rate information (about 4Hz) in speech signals. However, with the low-pass characteristics the slowly-varying channel bias components may also be emphasized. In order to handle this problem, in the experiments below the original Mel-Frequency Cepstral Coefficients (MFCC) are first processed by *Cepstral Normalization* (CN) [4] in order to properly reduce the low-frequency components. Therefore multi-eigenvector temporal filters discussed here are derived from, and performed on, the normalized cepstral coefficients. In other words, the CN process was first performed on the MFCC's of the training speech database to obtain the multi-eigenvector filters. These filters are then applied to the time trajectories of the MFCC's of both the training and testing database to obtain the new feature parameters. These new parameters are finally used in model training as well as the testing experiments.

### 3. EXPERIMENTAL SETUP

The AURORA2 database distributed by ETSI committee is used for the experiments here. It contains several sets of noisy speech with additive noise of different characteristics and levels plus some channel effects, as representatives for real-world environments. The clean speech data in the AURORA2 database consisting of 8440 utterances of English connected digits were used to obtain the data-driven

temporal filters in the experiments. Each utterance was first converted to 13-dimensional Mel-frequency cepstral coefficients (12 MFCCs + log cepstral coefficient) using the AURORA Front-end, as defined in the documentation. The resulted 8440 sets of MFCC's were then processed by *Cepstral Normalization* (CN) algorithm such that the mean and variance parameters were normalized to 0 and 1 respectively. Then for each time trajectory of feature vectors, two versions of temporal filters were constructed, one with the previously proposed PCA-based approach with a single eigenvector [1,2], and another with the multi-eigenvector approach proposed here as described in section 2. The length  $L$  of the temporal filters was empirically set to be 15.

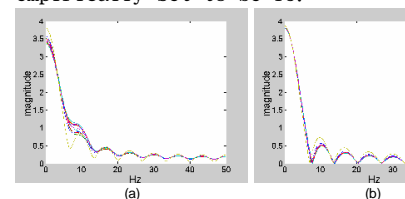


Figure 2. The frequency responses of (a) the 13 multi-eigenvector temporal filters and (b) the 13 previously proposed PCA-derived filters

Figure 2 shows the frequency responses of the obtained temporal filters: (a) the new multi-eigenvector approach proposed here in this report, and (b) the previously proposed PCA approach with a single eigenvector for the 13 time trajectories, trained on the clean speech database of AURORA2. From these figures we can have some observations as follows.

1. Both sets of the temporal filters are low-pass. They don't attenuate the

modulation frequency components, although CMS or RASTA does. This is why the *Cepstral Normalization* (CN) process is useful here.

2. The 13 previously proposed PCA-derived filters for the 13 temporal trajectories are very close, which the 13 multi-eigenvector temporal filters are obviously different around 0~15Hz.

3. The previously proposed PCA-derived filters have the shape of a main-lobe and several side-lobes; but the multi-eigenvector filters proposed here have a wider pass-band without apparent zeros.

4. For the new multi-eigenvectors filters proposed here, the modulation-frequency components around 0~4.5Hz are specially emphasized, therefore they may somehow enhance the syllabic-rate information (about 4Hz [5]) in the speech signals, which is a possible reason for the recognition performance improvements as found in the data below.

#### 4. EXPERIMENTAL RESULTS

In the recognition experiments, there are two training modes in AURORA2: *clean speech training* and *multi-condition training*. In the multi-condition training, the acoustic models were trained with speech data under different noisy conditions, added with different types of noise at different levels and so on. For each training mode, three sets (Sets A, B and C) of utterances artificially contaminated by different types of noise (subway, babble, car, etc.) at different SNR levels (ranging from -5dB to 20dB) were tested. Since the proposed approach

here only has to do with the feature extraction, all the following procedures for training and recognition are exactly identical to the related experiments stated in the documentation.

In the training process 13-dimensional normalized MFCC features i.e., the MFCC features but processed by CN, were used to construct the multi-eigenvector and the previously proposed PCA-derived temporal filters and these two sets of temporal filters were then applied on these 13-dimensional normalized MFCC features. The resulting 13-dimensional new features plus delta and delta-delta features were components in the finally 39-dimensional feature vectors. These new feature vectors of the HMM for each digit were trained. Similarly, testing phase the clean and corrupted testing speech data were first converted to MFCC's, processed and then individually processed above two sets of temporal filters optimized with the training data, to various sets of feature vectors for testing.

##### 4.1 Recognition results

Table 1 lists the recognition results respectively for the baseline experiment (Baseline), i.e., with the original features without any further processing and the experiments with the features processed by CN only (CN), first by the multi-eigenvector filters (CN+MEV) and then by the previously proposed PCA-derived filters (CN+PCA), first by the multi-eigenvector filters (CN+MEV+PCA) and then by the multi-eigenvector filters (CN+MEV+PCA).

proposed here (CN+M-eigen). The results include those for two training modes, clean speech training and multi-condition training, and three testing sets, sets A, B, and C. The length of the filters,  $L$ , for both of the two latter cases are set to be 15, and the number of eigenvectors,  $M$ , for the multi-eigenvector filters proposed here in eq. (6) is set to be 3. The word accuracy listed in Table 1 is the average of the recognition rates between 0–20dB. The overall word error rate (WER) improvements in Table 1 were calculated with respect to the baseline results. Note that with the proposed multi-eigenvector filters proposed in this report applied after CN, very significant improvements over CN only or CN plus the previously proposed PCA-derived filters were obtained in both clean speech and multi-condition training modes, and the improvements in the clean training mode are specially high. This verified that the new multi-eigenvector temporal filters are particularly effective when the

training and testing environments mismatched. Retained comparison data in Table 1 indicates that the proposed multi-eigenvector filters performed better than the previously proposed PCA-derived filters in all cases, for all the testing sets in all training modes.

From the right part of Table 1, we observe that in the multi-condition training mode the new multi-eigenvector filters can successfully improve the recognition performance of CN-processed speech features, but this is not always true for the previously proposed PCA-derived filters.

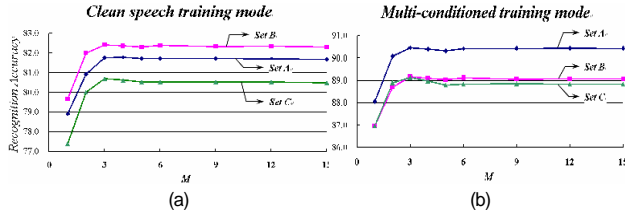
#### 4.2 Choice of the parameter $M$

We mentioned in section 2.2 that the parameter  $M$ , or the number of eigenvectors to be included, is chosen to be 3 because the eigenvalues  $\tilde{\epsilon}_{i,k}$  becomes much smaller for  $i > 3$ . Here some

Approaches	Clean speech training				Multi-condition training		
	Set A	Set B	Set C	Overall WER improvements	Set A	Set B	Set C
Baseline	61.34	55.76	66.14	-	87.82	86.22	83.78
CN	70.18	70.78	66.37	20.62	89.67	88.07	86.16
CN+PCA	78.99	79.69	77.47	45.31	87.91	86.90	86.76
CN+M-eigen	81.90	82.78	80.83	53.33	90.43	89.22	89.11

**Table 1.** Word accuracy for the three testing sets A, B, C under clean - and multi-condition training modes. The overall word error rate (WER) improvements were calculated with respect to the baseline experiments.

experimental results for different values of  $M$  are presented. The results for the same recognition experiments, as in Table 1, i.e., for the two training modes and three test sets, but using different values of  $M$  ( $M=1, 2, 3, 4, 5, 6, 9, 12, 15$ ) are depicted in Figure 3(a)(b). Of course here the case of  $M=1$  is exactly the previously proposed PCA-derived filters (CN+PCA in Table 1). It can be found clearly from this figure that very sharp improvements in performance were obtainable as  $M$  was increased from 1 to 3 in all cases, but the improvements turned out to be saturated if  $M$  was further increased from 3. As mentioned above, the eigenvalues  $\check{e}_{4,k} \sim \check{e}_{15,k}$  are very small relative to  $\check{e}_{1,k} \sim \check{e}_{3,k}$ ; therefore the information projected on the first three eigenvectors  $\check{o}_{1,k} \sim \check{o}_{3,k}$  are much more important, and actually dominate the recognition processes.



**Figure 3.** Recognition accuracy for multi-eigenvector filters with different  $M$  values: (a) clean speech training and (b) multi-condition training

### 4.3 Choice of the filter

#### length $L$

When designing the temporal filters including the multi-eigenvector filters as proposed here, it is well known that if the length of the temporal filters  $L$  is smaller, the width of the pass-band of the filters will be larger. This is clear by comparing Figure 2(a) for  $L=15$  with Figure 4(a)(b) for multi-eigenvector filters with  $M=3$  but  $L=10$  and 20 respectively. From these figures we can observe that the 3 dB width of the pass-band are around 6~7Hz and 4Hz for the cases of  $L=10$  and 20 respectively, while that for  $L=15$  is about 4.5Hz. Since the syllabic-rate of human speech is roughly around 4Hz, so the characteristics of the filters for the case  $L=15$  may exactly emphasize the syllabic-rate information and thus the performance roughly saturates at  $L=15$ , as can be found by the recognition accuracy shown in Table 2.

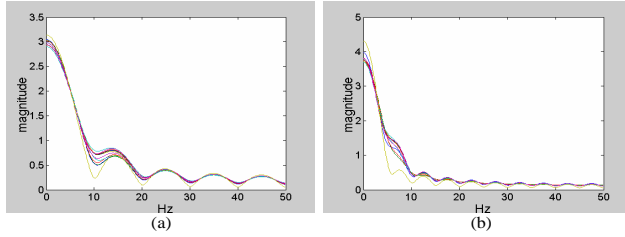


Figure 4. The frequency responses of the 13 multi-eigenvector filters as  $L$  is set to (a) 10 (b) 20

	Clean speech training		
$L$	Set A	Set B	Set C
10	79.95	80.22	77.23
15	81.90	82.78	80.83
20	81.75	82.64	80.79

Table 2. Recognition accuracy for multi-eigenvector filters with  $M=3$  and different values of  $L$  with clean speech training mode

#### 4.4 Further comparison between the two versions of temporal filters with some other metric

In this subsection, we'd like to compare the multi-eigenvector filters proposed here with the previously proposed PCA-derived filters using a metric different from recognition accuracy. The metric used is the average of the normalized distance between the corrupted features,  $\hat{x}$ , and the corresponding clean speech features,  $x$ ,

$$d = E \left[ \frac{\|\hat{x} - x\|}{\|x\|} \right], \quad (7)$$

where the average is taken over all the testing speech. This metric is to provide an estimated measure of the robustness of the temporal-filtered MFCC features with respect to the corruption. Smaller values of  $d$  imply that the features are less influenced by the corruption. Table 3 compares this distance measure  $d$  for the three testing sets A, B, C under different SNR values. We see that the multi-eigenvector filters proposed here gives smaller averaged normalized distance in all cases. This offers another explanation why the proposed multi-eigenvector filters give better recognition accuracy.

	Approaches	20db	15db	10db	5db	0db	-5db
Set A	CN+PCA	0.6826	0.7554	0.8281	0.9081	0.9932	1.0855
	CN+M-eigen	0.6211	0.6946	0.7688	0.8521	0.9473	1.0587
Set B	CN+PCA	0.6693	0.7394	0.8150	0.9053	1.0016	1.1086
	CN+M-eigen	0.6071	0.6778	0.7556	0.8493	0.9555	1.0782
Set C	CN+PCA	0.6836	0.7561	0.8305	0.9123	1.0057	1.1029
	CN+M-eigen	0.6238	0.7002	0.7808	0.8763	0.9902	1.1098

Table 3. The average normalized distance between clean and corrupted speech features

under various SNRs and two different temporal filter techniques

## 5. CONCLUSION

In this report, we proposed a multi-eigenvector approach of designing data-driven temporal filters, in which more than one eigenvectors are weighted by their corresponding eigenvalues and summed to form the filter coefficients. Very encouraging experimental results have been obtained and this approach is also shown to be particularly effective when the training and testing environments are highly mismatched.

## 6. REFERENCE

- [1] J-W. Hung. et al “Comparative Analysis for Data-Driven Temporal Filters Obtained Via Principal Component Analysis (PCA) and Linear Discriminant Analysis (LDA) In Speech Recognition”, Eurospeech 2001
- [2] J-W. Hung and L-S. Lee “Data-Driven Temporal Filters Obtained Via Different Optimization Criteria Evaluated On AURORA2 Database”, ICSLP 2002
- [3] S. Furui, "Cepstral analysis technique for automatic speaker verification". IEEE Trans. Acoust. Speech Signal Process. 1981
- [4] O. Viikki and K. Laurila, “Noise robust HMM-based speech recognition using segmental cepstral feature vector normalization,” in ESCA NATO Workshop Robust Speech Recognition Unknown Communication Channels, Pont-a-Mousson, France, 1997, pp. 107–110.
- [5] H. Hermansky and N. Morgan, “RASTA processing of speech”. IEEE Trans. Speech Audio Process. 2, 578-589, 1994
- [6] C. Avendano, S. van Vuuren and H. Hermansky, "Data Based Filter Design for RASTA-like Channel Normalization in ASR" ICSLP 96
- [7] S. van Vuuren and H. Hermansky, "Data-driven Design of RASTA-like Filters", Eurospeech 97
- [8] K. Fukunaga, “Introduction to statistical pattern recognition”, E.2<sup>nd</sup>, Academic Press, 1990