

# Useful Derivations for i-Vector Based Approach to Data Clustering in Speech Recognition

Yu Zhang

December 16, 2011

## Contents

<b>1</b>	<b>Background</b>	<b>2</b>
1.1	Traditional maximum likelihood eigen-decomposition . . . . .	2
1.2	What's new in Kenny's paper? . . . . .	2
<b>2</b>	<b>The i-vector based approach to data clustering</b>	<b>3</b>
2.1	Data model . . . . .	3
2.2	Frontend i-vector extraction . . . . .	3
2.2.1	The i-vector solution . . . . .	3
2.2.2	Smoothed data model . . . . .	4
2.2.3	Solving i-vector . . . . .	5
2.3	Hyperparameter estimation . . . . .	6
2.3.1	Updating $\mathbf{T}$ . . . . .	7
2.3.2	Updating $\mathbf{R}$ . . . . .	7
2.3.3	Summarizing the EM algorithm . . . . .	8
<b>3</b>	<b>Implementation</b>	<b>9</b>
3.1	Training $\mathbf{T}$ and $\mathbf{R}$ . . . . .	9
3.2	Extracting $\hat{\mathbf{w}}(i)$ . . . . .	9
3.3	Computational complexity . . . . .	9
<b>4</b>	<b>Clustering of i-vectors using LBG algorithm</b>	<b>10</b>
<b>A</b>	<b>The Matrix Codebook for this memo</b>	<b>11</b>
A.1	Trace . . . . .	11
A.2	Variance of quadratic form . . . . .	11
A.3	Others . . . . .	12
<b>B</b>	<b>From a graphic model view</b>	<b>12</b>
B.1	Probabilistic PCA . . . . .	12
B.1.1	Maximum likelihood PCA . . . . .	13

# 1 Background

## 1.1 Traditional maximum likelihood eigen-decomposition

Let  $\mathbf{M}(s)$  denote the mean supervector for a speaker  $s$  formulating as:

$$\mathbf{M}(s) = \mathbf{M}_0 + \mathbf{T}\mathbf{w}, \quad (1)$$

and the likelihood function written as:

$$\prod_s \max_{\mathbf{w}} P_{\text{HMM}}(\mathbf{Y}_s | \mathbf{M}_0 + \mathbf{T}\mathbf{w}, \mathbf{R}). \quad (2)$$

The optimization proceeds by iterating the following two steps:

1. For each training speaker  $s$ , use the current estimates of  $\mathbf{T}$  and  $\mathbf{R}$  to find the  $\mathbf{w}$  which maximizes the HMM likelihood given the speaker’s training data  $\mathbf{Y}_s$ :

$$\mathbf{w}(s) = \arg \max_{\mathbf{w}} P_{\text{HMM}}(\mathbf{Y}_s | \mathbf{M}_0 + \mathbf{T}\mathbf{w}, \mathbf{R}) \quad (3)$$

2. Update  $\mathbf{T}$  and  $\mathbf{R}$  by maximizing:

$$\prod_s P_{\text{HMM}}(\mathbf{Y}_s | \mathbf{M}_0 + \mathbf{T}\mathbf{w}(s), \mathbf{R}) \quad (4)$$

## 1.2 What’s new in Kenny’s paper?

In Kenny’s paper [3], they treat  $\mathbf{w}(i)$  as a random vector with a standard normal distribution. So the estimation of  $\mathbf{T}$  and  $\mathbf{R}$  becomes to maximize the marginal likelihood function

$$\prod_s \int_{\mathbf{w}} P(\mathbf{Y}_s | \mathbf{M}_0 + \mathbf{T}\mathbf{w}, \mathbf{R}) \mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{I}) d\mathbf{w}, \quad (5)$$

where the product extends over all speakers in the training set.

Calculating the posterior distribution of  $\mathbf{w}(i)$  in the E-step rather than the maximum likelihood estimate is the key to avoiding the degeneracy problem that arises when Eq. (2) is used to estimate the eigenvoices in situations where speaker-dependent training is not feasible (e.g., when training data is sparse) or the number of eigenvoices is large compared with the number of training speakers.

The algorithm can be briefly break-down into the following steps:

1. For each training speaker  $s$ , use the current alignment of the speaker’s training data and the current estimates  $\mathbf{T}$  and  $\mathbf{R}$  to carry out MAP speaker adaptation. Use the speaker adapted model to realign the speaker’s training data.
2. The **E-step**: For each speaker  $s$ , calculate the posterior distribution of  $\mathbf{w}(i)$  using the current alignment of the speaker’s training data, the current estimates of  $\mathbf{T}$  and  $\mathbf{R}$  and the prior  $\mathcal{N}(\mathbf{w} | \mathbf{0}, \mathbf{I})$ .
3. The **M-step**: Update  $\mathbf{T}$  and  $\mathbf{R}$  by a linear regression in which the  $\mathbf{w}(i)$ ’s play the role of the explanatory variables.

## 2 The i-vector based approach to data clustering

We are considering in this memo to use i-vector based approach for more generic data clustering problems in speech recognition. Possible applications include training data clustering for multiple acoustic model training [7] and acoustic sniffing in IVN-based training [5]. So hereafter, we will change the subscript  $s$  representing the speaker to subscript  $i$  representing each individual acoustic unit that to be clustered. The granularity of the acoustic units can be one single utterance, several utterances (e.g., one conversational side in Switchboard), or even more utterances from a pre-defined source (speaker, environment, channel, etc.)

### 2.1 Data model

Suppose we are given a set of training data denoted as  $\mathcal{Y} = \{\mathbf{Y}^i | i = 1, 2, \dots, I\}$ , where  $\mathbf{Y}^i = (\mathbf{y}_1^i, \mathbf{y}_2^i, \dots, \mathbf{y}_{T_i}^i)$  is a sequence of  $D$ -dimensional feature vectors extracted from the  $i$ -th acoustic unit (e.g., utterance). From  $\mathcal{Y}$ , a Gaussian Mixture Model (GMM) can be trained using a Maximum Likelihood (ML) criterion to serve as a Universal Background Model (UBM). Let's use  $\Omega = \{c_k, \mathbf{m}_k, \mathbf{R}_k | k = 1, \dots, K\}$  to denote the set of UBM-GMM parameters where  $c_k$ 's are mixture component weights,  $\mathbf{m}_k$  and  $\mathbf{R}_k$  are  $D$ -dimensional mean and  $D \times D$  diagonal covariance matrix for the  $k^{\text{th}}$  mixture component.

Given the data model  $\Omega$ , the probability of each acoustic unit  $\mathbf{Y}^i$  can be written as:

$$p(\mathbf{Y}^i | \Omega) = \prod_{t=1}^{T_i} \sum_{k=1}^K c_k \mathcal{N}(\mathbf{y}_t^i; \mathbf{m}_k, \mathbf{R}_k). \quad (6)$$

And we denote  $\Omega^{(0)}$  hereafter the initial model parameters of the UBM.

### 2.2 Frontend i-vector extraction

#### 2.2.1 The i-vector solution

Let  $\mathbf{M}_0$  denote the  $(D \cdot K)$ -dimensional supervector by concatenating the  $\mathbf{m}_k$ 's. Given an utterance  $\mathbf{Y}^i$ , let's use another  $(D \cdot K)$ -dimensional random supervector  $\mathbf{M}(i)$  to characterize it independent of its linguistic content. In i-vector based approach,  $\mathbf{M}(i)$  is correlated with  $\mathbf{M}_0$  as:

$$\mathbf{M}(i) = \mathbf{M}_0 + \mathbf{T}\mathbf{w}(i), \quad (7)$$

where  $\mathbf{T}$  is a fixed but unknown  $(D \cdot K) \times F$  rectangular matrix of low rank (i.e.,  $F \ll (D \cdot K)$ ), and  $\mathbf{w}(i)$  is an  $F$ -dimensional random vector having a prior distribution of  $\mathcal{N}(\cdot; \mathbf{0}, \mathbf{I})$ . A graphical model representation is shown in Fig. 1. In [2],  $\mathbf{T}$  is called the *total variability matrix* and  $\mathbf{w}(i)$  the *i-vector*. They are also comparable with the *loading matrix* and *factors* in factor analysis.

Figure 1: A graphical model representation of *i*-vector approach.

Given  $\mathbf{Y}^i$ ,  $\mathbf{M}(i)$ , and  $\mathbf{R}_k$ 's, the *i*-vector is the MAP solution of the following problem:

$$\begin{aligned}\hat{\mathbf{w}}(i) &= \operatorname{argmax}_{\mathbf{w}(i)} p(\mathbf{w}(i)|\mathbf{Y}^i) \\ &= \operatorname{argmax}_{\mathbf{w}(i)} p(\mathbf{Y}^i|\mathbf{w}(i))p(\mathbf{w}(i)),\end{aligned}\quad (8)$$

where

$$p(\mathbf{Y}^i|\mathbf{w}(i)) = \prod_{t=1}^{T_i} \sum_{k=1}^K c_k \mathcal{N}(\mathbf{y}_t^i; \mathbf{M}_k(i), \mathbf{R}_k), \quad (9)$$

in which  $\mathbf{M}_k(i)$  is the  $k$ -th  $D$ -dimensional subvector of  $\mathbf{M}(i)$ .

### 2.2.2 Smoothed data model

Eq. (9) is intractable because of the summation term. So Viterbi approximation was imposed in [3] so that

$$p(\mathbf{Y}^i|\mathbf{w}(i)) \simeq \prod_{t=1}^{T_i} \max_k \mathcal{N}(\mathbf{y}_t^i; \mathbf{M}_k(i), \mathbf{R}_k). \quad (10)$$

One can also define a ‘‘smoothed’’ version of  $p(\mathbf{Y}^i|\mathbf{w}(i))$ :

$$p(\mathbf{Y}^i|\mathbf{w}(i)) \simeq \prod_{t=1}^{T_i} \prod_{k=1}^K \mathcal{N}(\mathbf{y}_t^i; \mathbf{M}_k(i), \mathbf{R}_k)^{p(k|\mathbf{y}_t^i, \boldsymbol{\Omega}^{(0)})}, \quad (11)$$

where

$$p(k|\mathbf{y}_t^i, \boldsymbol{\Omega}^{(0)}) = \frac{c_k \mathcal{N}(\mathbf{y}_t^i; \mathbf{m}_k, \mathbf{R}_k^{(0)})}{\sum_{l=1}^K c_l \mathcal{N}(\mathbf{y}_t^i; \mathbf{m}_l, \mathbf{R}_l^{(0)})} \quad (12)$$

is the Baum-Welch occupancy probability of the  $k^{\text{th}}$  component given  $\mathbf{y}_t^i$ . Note that the production terms in Eq. (11) are not strictly probability densities because they do not integrate to 1. Using Eq. (11) will lead to the same *i*-vector solution as in [2], although in that paper, the form of  $p(\mathbf{Y}^i|\mathbf{w}(i))$  was never explicitly given.

### 2.2.3 Solving *i*-vector

Given the smoothed data model in Eq. (11), the closed-form solution of the problem in Eq. (8) gives the *i*-vector extraction formula as follows:

$$\hat{\mathbf{w}}(i) = \boldsymbol{l}^{-1}(i) \mathbf{T}^\top \mathbf{R}^{-1} \boldsymbol{\Gamma} \mathbf{y}(i), \quad (13)$$

where

$$\boldsymbol{l}(i) = \mathbf{I} + \mathbf{T}^\top \boldsymbol{\Gamma}(i) \mathbf{R}^{-1} \mathbf{T}, \quad (14)$$

$\mathbf{R}$  is the  $(D \cdot K) \times (D \cdot K)$  block-diagonal matrix with  $\mathbf{R}_k$  as its  $k^{\text{th}}$  block component,  $\mathbf{\Gamma}(i)$  is a  $(D \cdot K) \times (D \cdot K)$  block-diagonal matrix with  $\gamma_k(i)\mathbf{I}_{D \times D}$  as its  $k^{\text{th}}$  block component;  $\mathbf{\Gamma}_{\mathbf{y}}(i)$  is a  $(D \cdot K)$ -dimensional supervector with  $\mathbf{\Gamma}_{\mathbf{y},k}(i)$  as its  $k^{\text{th}}$   $D$ -dimensional subvector. The statistics  $\gamma_k(i)$  and  $\mathbf{\Gamma}_{\mathbf{y},k}(i)$  are calculated as follows:

$$\begin{aligned}\gamma_k(i) &= \sum_{t=1}^{T_i} p(k|\mathbf{y}_t^i, \mathbf{\Omega}^{(0)}), \\ \mathbf{\Gamma}_{\mathbf{y},k}(i) &= \sum_{t=1}^{T_i} p(k|\mathbf{y}_t^i, \mathbf{\Omega}^{(0)})(\mathbf{y}_t^i - \mathbf{m}_k).\end{aligned}\tag{15}$$

PROOF: First,

$$\begin{aligned}\log p(\mathbf{Y}^i|\mathbf{w}(i)) &= \sum_{t=1}^{T_i} \sum_{k=1}^K p(k|\mathbf{y}_t^i, \mathbf{\Omega}^{(0)}) \log \mathcal{N}(\mathbf{y}_t^i; \mathbf{M}_k(i), \mathbf{R}_k) \\ &= \sum_{t=1}^{T_i} \sum_{k=1}^K p(k|\mathbf{y}_t^i, \mathbf{\Omega}^{(0)}) \left[ \log \frac{1}{(2\pi)^{D/2} |\mathbf{R}_k|^{1/2}} \right. \\ &\quad \left. - \frac{1}{2} (\mathbf{y}_t^i - \mathbf{m}_k - \mathbf{T}_k \mathbf{w}(i))^\top \mathbf{R}_k^{-1} (\mathbf{y}_t^i - \mathbf{m}_k - \mathbf{T}_k \mathbf{w}(i)) \right] \\ &= \sum_{t=1}^{T_i} \sum_{k=1}^K p(k|\mathbf{y}_t^i, \mathbf{\Omega}^{(0)}) \left[ \log \frac{1}{(2\pi)^{D/2} |\mathbf{R}_k|^{1/2}} - \frac{1}{2} (\mathbf{y}_t^i - \mathbf{m}_k)^\top \mathbf{R}_k^{-1} (\mathbf{y}_t^i - \mathbf{m}_k) \right. \\ &\quad \left. + \mathbf{w}^\top(i) \mathbf{T}_k^\top \mathbf{R}_k^{-1} (\mathbf{y}_t^i - \mathbf{m}_k) - \frac{1}{2} \mathbf{w}^\top(i) \mathbf{T}_k^\top \mathbf{R}_k^{-1} \mathbf{T}_k \mathbf{w}(i) \right].\end{aligned}\tag{16}$$

Only the last two terms are related with  $\mathbf{w}(i)$ , and we can define:

$$\begin{aligned}\mathcal{H}(i) &= \sum_{t=1}^{T_i} \sum_{k=1}^K p(k|\mathbf{y}_t^i, \mathbf{\Omega}^{(0)}) \left[ \mathbf{w}^\top(i) \mathbf{T}_k^\top \mathbf{R}_k^{-1} (\mathbf{y}_t^i - \mathbf{m}_k) - \frac{1}{2} \mathbf{w}^\top(i) \mathbf{T}_k^\top \mathbf{R}_k^{-1} \mathbf{T}_k \mathbf{w}(i) \right] \\ &= \mathbf{w}^\top(i) \mathbf{T} \mathbf{R}^{-1} \mathbf{\Gamma}_{\mathbf{y}}(i) - \frac{1}{2} \mathbf{w}^\top(i) \mathbf{T}^\top \mathbf{\Gamma}(i) \mathbf{R}^{-1} \mathbf{T} \mathbf{w}(i).\end{aligned}\tag{17}$$

So

$$\begin{aligned}p(\mathbf{w}(i)|\mathbf{Y}^i) &\propto p(\mathbf{Y}^i|\mathbf{w}(i))p(\mathbf{w}(i)) \\ &\propto \exp\left(\mathbf{w}^\top(i) \mathbf{T} \mathbf{R}^{-1} \mathbf{\Gamma}_{\mathbf{y}}(i) - \frac{1}{2} \mathbf{w}^\top(i) \mathbf{T}^\top \mathbf{\Gamma}(i) \mathbf{R}^{-1} \mathbf{T} \mathbf{w}(i)\right) \cdot \exp\left(-\frac{1}{2} \mathbf{w}^\top(i) \mathbf{w}(i)\right) \\ &= \exp\left(\mathbf{w}^\top(i) \mathbf{T} \mathbf{R}^{-1} \mathbf{\Gamma}_{\mathbf{y}}(i) - \frac{1}{2} \mathbf{w}^\top(i) [\mathbf{T}^\top \mathbf{\Gamma}(i) \mathbf{R}^{-1} \mathbf{T} + \mathbf{I}] \mathbf{w}(i)\right) \\ &= \exp\left(\mathbf{w}^\top(i) \mathbf{T} \mathbf{R}^{-1} \mathbf{\Gamma}_{\mathbf{y}}(i) - \frac{1}{2} \mathbf{w}^\top(i) \mathbf{l}(i) \mathbf{w}(i)\right) \\ &\propto \exp\left(-\frac{1}{2} (\mathbf{w}(i) - \mathbf{l}^{-1}(i) \mathbf{T}^\top \mathbf{R}^{-1} \mathbf{\Gamma}_{\mathbf{y}}(i))^\top \mathbf{l}(i) (\mathbf{w}(i) - \mathbf{l}^{-1}(i) \mathbf{T}^\top \mathbf{R}^{-1} \mathbf{\Gamma}_{\mathbf{y}}(i))\right).\end{aligned}\tag{18}$$

Therefore, the ‘‘posterior’’ distribution of  $\mathbf{w}(i)$  is Gaussian of mean  $\mathbf{l}^{-1}(i) \mathbf{T}^\top \mathbf{R}^{-1} \mathbf{\Gamma}_{\mathbf{y}}(i)$  and covariance  $\mathbf{l}^{-1}(i)$ . ■

### 2.3 Hyperparameter estimation

Given the training data  $\mathcal{Y}$ , the hyperparameters  $\mathbf{T}$  and  $\mathbf{R}$  can be estimated by maximizing the following log-likelihood function:

$$\begin{aligned}
\mathcal{F}(\mathbf{T}, \mathbf{R}) &= \log \prod_{i=1}^I \int p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i, \mathbf{w}(i)) d\mathbf{w}(i) \\
&= \sum_{i=1}^I \log \int p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i, \mathbf{w}(i)) d\mathbf{w}(i) \\
&= \sum_{i=1}^I \log \int \frac{p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i, \mathbf{w}(i))}{p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{Y}^i, \mathbf{w}(i))} p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{Y}^i, \mathbf{w}(i)) d\mathbf{w}(i) \\
&= \sum_{i=1}^I \log \int \frac{p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i, \mathbf{w}(i))}{p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{Y}^i, \mathbf{w}(i))} p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{w}(i) | \mathbf{Y}^i) d\mathbf{w}(i) + C_1,
\end{aligned} \tag{19}$$

where  $\mathbf{T}^{(0)}$  and  $\mathbf{R}^{(0)}$  are hyperparameters from last iteration,  $C_1$  is a constant independent of  $\mathbf{T}$  and  $\mathbf{R}$  (we will use  $C_j$ 's hereafter for those constants). By Jensen's inequality, we have:

$$\begin{aligned}
\mathcal{F}(\mathbf{T}, \mathbf{R}) &\geq \sum_{i=1}^I \int \log \left( \frac{p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i, \mathbf{w}(i))}{p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{Y}^i, \mathbf{w}(i))} \right) p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{w}(i) | \mathbf{Y}^i) d\mathbf{w}(i) \\
&= \sum_{i=1}^I \int \log p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i, \mathbf{w}(i)) p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{w}(i) | \mathbf{Y}^i) d\mathbf{w}(i) + C_2 \\
&= \sum_{i=1}^I \int \log p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i | \mathbf{w}(i)) p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{w}(i) | \mathbf{Y}^i) d\mathbf{w}(i) + C_3.
\end{aligned} \tag{20}$$

So an auxiliary function  $\mathcal{A}$  can be built as:

$$\begin{aligned}
\mathcal{A}(\mathbf{T}, \mathbf{R}) &= \sum_{i=1}^I \int \log p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i | \mathbf{w}(i)) p_{\mathbf{T}^{(0)}, \mathbf{R}^{(0)}}(\mathbf{w}(i) | \mathbf{Y}^i) d\mathbf{w}(i) \\
&= \sum_{i=1}^I E(\log p_{\mathbf{T}, \mathbf{R}}(\mathbf{Y}^i | \mathbf{w}(i)))
\end{aligned} \tag{21}$$

and the log likelihood function is

$$\log \mathcal{F}(\mathbf{T}, \mathbf{R}) = \sum_i \mathcal{G}(i) - \frac{1}{2} |l(i)| + \frac{1}{2} E[\mathbf{w}(i)^\top] \mathbf{T}^\top \mathbf{R}^{-1} \Gamma_{\mathbf{y}}(i) \tag{22}$$

### 2.3.1 Updating $\mathbf{T}$

Referring to Eqs (16) and (17), for updating  $\mathbf{T}$ , only the  $\mathcal{H}$  function is involved. So given the properties in Eq. (48), we have

$$\begin{aligned}
\mathcal{A} &= \sum_i E[\mathbf{w}^\top(i) \mathbf{T}^\top \mathbf{R}^{-1} \mathbf{\Gamma}_y(i) - \frac{1}{2} \mathbf{w}^\top(i) \mathbf{T}^\top \mathbf{\Gamma}(i) \mathbf{R}^{-1} \mathbf{T} \mathbf{w}(i)] + C_4 \\
&= \sum_i E[\mathbf{w}^\top(i) \mathbf{T}^\top \mathbf{R}^{-1} \mathbf{\Gamma}_y(i) - \frac{1}{2} \text{tr}(\mathbf{T}^\top \mathbf{\Gamma}(i) \mathbf{R}^{-1} \mathbf{T} E[\mathbf{w}(i) \mathbf{w}^\top(i)])] + C_4 \\
&= \sum_i \text{tr} \left( \mathbf{R}^{-1} \left( \mathbf{\Gamma}_y(i) E[\mathbf{w}^\top(i)] - \frac{1}{2} \mathbf{\Gamma}(i) \mathbf{T} E[\mathbf{w}(i) \mathbf{w}^\top(i)] \right) \mathbf{T}^\top \right) + C_4,
\end{aligned} \tag{23}$$

and

$$\frac{\partial \mathcal{A}}{\partial \mathbf{T}} = \sum_i \mathbf{R}^{-1} \left( \mathbf{\Gamma}_y(i) E[\mathbf{w}^\top(i)] - \mathbf{\Gamma}(i) \mathbf{T} E[\mathbf{w}(i) \mathbf{w}^\top(i)] \right). \tag{24}$$

Setting Eq. (24) to 0, and taking use of the diagonal structure of  $\mathbf{\Gamma}(i)$ ,  $\mathbf{T}$  can be solved line-by-line according to:

$$\mathbf{T}^m \sum_i \mathbf{\Gamma}^m(i) E[\mathbf{w}(i) \mathbf{w}^\top(i)] = \sum_i \mathbf{\Gamma}_y^m(i) E[\mathbf{w}^\top(i)], \tag{25}$$

in which  $\mathbf{T}^m$ ,  $\mathbf{\Gamma}^m$  and  $\mathbf{\Gamma}_y^m$  are the  $m^{\text{th}}$  row of  $\mathbf{T}$ ,  $\mathbf{\Gamma}$  and  $\mathbf{\Gamma}_y$ , respectively.

### 2.3.2 Updating $\mathbf{R}$

Go back to Eq. (16) and set

$$\begin{aligned}
\mathcal{G}(i) &= \sum_{t=1}^{T_i} \sum_{k=1}^K p(k | \mathbf{y}_t^i, \mathbf{\Omega}^{(0)}) \left[ \log \frac{1}{(2\pi)^{D/2} |\mathbf{R}_k|^{1/2}} - \frac{1}{2} (\mathbf{y}_t^i - \mathbf{m}_k)^\top \mathbf{R}_k^{-1} (\mathbf{y}_t^i - \mathbf{m}_k) \right] \\
&= \frac{1}{2} \sum_{k=1}^K \gamma_k(i) \log |\mathbf{R}_k^{-1}| - \text{tr}(\mathbf{R}_k^{-1} \mathbf{\Gamma}_{\mathbf{y}\mathbf{y}^\top, k}(i)) + C_5,
\end{aligned} \tag{26}$$

in which

$$\mathbf{\Gamma}_{\mathbf{y}\mathbf{y}^\top, k}(i) = \sum_{t=1}^{T_i} p(k | \mathbf{y}_t^i, \mathbf{\Omega}^{(0)}) (\mathbf{y}_t^i - \mathbf{m}_k) (\mathbf{y}_t^i - \mathbf{m}_k)^\top. \tag{27}$$

It is not difficult to see that:

$$\mathcal{A} = \sum_i \mathcal{G}(i) + E[\mathcal{H}(i)]. \tag{28}$$

Because

$$\frac{\partial \mathcal{G}(i)}{\partial \mathbf{R}^{-1}} = \frac{1}{2} \sum_{k=1}^K \gamma_k(i) \mathbf{R} - \mathbf{\Gamma}_{\mathbf{y}\mathbf{y}^\top, k}(i), \tag{29}$$

and

$$\frac{\partial E[\mathcal{H}(i)]}{\partial \mathbf{R}^{-1}} = \frac{\partial \text{tr} \left( \mathbf{R}^{-1} \left( \mathbf{\Gamma}_{\mathbf{y}}(i) E[\mathbf{w}^\top(i)] - \frac{1}{2} \mathbf{\Gamma}(i) \mathbf{T} E[\mathbf{w}(i) \mathbf{w}^\top(i)] \right) \mathbf{T}^\top \right)}{\partial \mathbf{R}^{-1}} \quad (30)$$

Note that Eq. (25)

$$\begin{aligned} \sum_i \frac{\partial E[\mathcal{H}(i)]}{\partial \mathbf{R}^{-1}} &= \frac{\partial \text{tr} \left( \sum_i \mathbf{R}^{-1} \left( \mathbf{\Gamma}_{\mathbf{y}}(i) E[\mathbf{w}^\top(i)] - \frac{1}{2} \mathbf{\Gamma}(i) \mathbf{T} E[\mathbf{w}(i) \mathbf{w}^\top(i)] \right) \mathbf{T}^\top \right)}{\partial \mathbf{R}^{-1}} \\ &= \frac{1}{2} \sum_i \left( \mathbf{\Gamma}_{\mathbf{y}}(i) E[\mathbf{w}^\top(i)] \mathbf{T}^\top + \mathbf{T} E[\mathbf{w}(i)] \mathbf{\Gamma}_{\mathbf{y}}(i)^\top \right) \end{aligned} \quad (31)$$

Setting  $\frac{\partial \mathcal{A}}{\partial \mathbf{R}^{-1}}$  to 0, we have

$$\mathbf{R}_k = \frac{1}{\sum_i \gamma_k(i)} \left( \sum_i \mathbf{\Gamma}_{\mathbf{y}\mathbf{y}^\top, k}(i) - M_c \right) \quad (32)$$

where  $M_k$  denotes the  $k$ th diagonal block of the matrix

$$\frac{1}{2} \sum_i \left( \mathbf{\Gamma}_{\mathbf{y}}(i) E[\mathbf{w}^\top(i)] \mathbf{T}^\top + \mathbf{T} E[\mathbf{w}(i)] \mathbf{\Gamma}_{\mathbf{y}}(i)^\top \right). \quad (33)$$

### 2.3.3 Summarizing the EM algorithm

- **E-Step:**

$$\begin{aligned} E[\mathbf{w}(i)] &= \mathbf{l}^{-1}(i) \mathbf{T}^\top \mathbf{R}^{-1} \mathbf{\Gamma}_{\mathbf{y}}(i) \\ E[\mathbf{w}(i) \mathbf{w}^\top(i)] &= E[\mathbf{w}(i)] E[\mathbf{w}^\top(i)] + \mathbf{l}^{-1}(i) \end{aligned} \quad (34)$$

- **M-Step:**

$$\begin{aligned} \mathbf{T}^m \sum_i \mathbf{\Gamma}^m(i) E[\mathbf{w}(i) \mathbf{w}^\top(i)] &= \sum_i \mathbf{\Gamma}_{\mathbf{y}}^m(i) E[\mathbf{w}^\top(i)] \\ \mathbf{R}_k &= \frac{1}{\sum_i \gamma_k(i)} \left( \sum_i \mathbf{\Gamma}_{\mathbf{y}\mathbf{y}^\top, k}(i) - M_k \right) \end{aligned} \quad (35)$$

## 3 Implementation

### 3.1 Training $\mathbf{T}$ and $\mathbf{R}$

1. Train a GMM-UBM using all the training data.
2. Initialize  $\mathbf{R}$  by borrowing  $\mathbf{R}_k$ 's from the UBM, and  $\mathbf{T}$  randomly as:

$$\mathbf{T}^{m,f} \in [-\alpha \mathbf{R}^{m,m}, \alpha \mathbf{R}^{m,m}], \quad \forall m = 1, \dots, DK; f = 1, \dots, F, \quad (36)$$

where  $\mathbf{T}^{m,f}$  is the  $m^{\text{th}}$  line and  $f^{\text{th}}$  column of  $\mathbf{T}$ , which is similar for  $\mathbf{R}^{m,m}$ ;  $\alpha$  is typically set to 0.1.



3. For each acoustic unit  $i$ , extract the Baum-Welch statistics  $\gamma_k(i)$ ,  $\mathbf{\Gamma}_{\mathbf{y},k}(i)$  and  $\mathbf{\Gamma}_{\mathbf{y}\mathbf{y}^\top,k}(i)$  using the UBM as in Eqs. (15) and (27).
4. The E-Step: Calculate the posterior expectation  $E[\mathbf{w}(i)]$  and  $E[\mathbf{w}(i)\mathbf{w}^\top(i)]$  using the statistics and the current estimation of  $\mathbf{T}$  and  $\mathbf{R}$  with Eq. (34).
5. Repeat Steps 3 and 4 for all  $i$ 's. Accumulate necessary statistics for solving Eq. (35).
6. The M-Step: Updating  $\mathbf{T}$  and  $\mathbf{R}$  using Eq. (35).
7. If training converges, stop; otherwise go back to Step 4.

### 3.2 Extracting $\hat{\mathbf{w}}(i)$

Extracting  $\hat{\mathbf{w}}(i)$  is straightforward by using the above Steps 3 and 4 and setting:

$$\hat{\mathbf{w}}(i) = E[\mathbf{w}(i)]. \quad (37)$$

The procedure is the same for both training and testing acoustic units.

### 3.3 Computational complexity

In i-vector extraction, the most computational cost is the calculation of  $\mathbf{l}(i)$  in Eq. (14). The right-hand-side of the equation can be implemented step-by-step as:

$$\mathbf{R}^{-1} \rightarrow \mathbf{\Gamma}(i)\mathbf{R}^{-1} \rightarrow \mathbf{T}^\top\mathbf{\Gamma}(i)\mathbf{R}^{-1} \rightarrow \mathbf{T}^\top\mathbf{\Gamma}(i)\mathbf{R}^{-1}\mathbf{T}, \quad (38)$$

so the complexity is  $O(I \times DK \times F \times F)$ . In our implementation, the property that  $\mathbf{\Gamma}(i)$  is a block diagonal matrix can be used to make the computation more efficient:

1. Compute all the  $K$  block matrices of  $\mathbf{T}^\top\mathbf{R}^{-1}\mathbf{T}$ :

$$\mathbf{H}_k = \mathbf{T}^\top\mathbf{R}_k\mathbf{T} \quad \forall k = 1, \dots, K. \quad (39)$$

2. For each acoustic unit  $i$ ,

$$\mathbf{l}(i) = \sum_k \gamma_k(i)\mathbf{H}_k. \quad (40)$$

By doing so, the complexity of evaluating Eq. (14) is dominated by  $O(I \times K \times F \times F + DK \times F \times F)$ .

## 4 Clustering of i-vectors using LBG algorithm

As described above, given the training corpus, an i-vector can be extracted from each acoustic unit. Given the training set i-vectors, we use a hierarchical divisive clustering algorithm, namely LBG algorithm [4], to cluster them into multiple clusters. To measure the similarity between two i-vectors  $\hat{\mathbf{w}}(i)$  and  $\hat{\mathbf{w}}(j)$ , the both Euclidean distance and cosine similarity measure can be used.

It is quite easy to use Euclidean distance in clustering. However when cosine similarity is used, the calculation of the clustering centroid needs to be reviewed. Given the cosine similarity defined as:

$$\text{sim}(\hat{\mathbf{w}}(i), \hat{\mathbf{w}}(j)) = \frac{\hat{\mathbf{w}}^\top(i) \hat{\mathbf{w}}(j)}{\|\hat{\mathbf{w}}(i)\| \|\hat{\mathbf{w}}(j)\|}, \quad (41)$$

it can be proved that the corresponding clustering centroid,  $\mathbf{c}^{\hat{\mathbf{w}}}$  consisting of  $n$  i-vectors  $\hat{\mathbf{w}}(1), \hat{\mathbf{w}}(2), \dots, \hat{\mathbf{w}}(n)$ , can be calculated as:

$$\mathbf{c}^{\hat{\mathbf{w}}} = \underset{\mathbf{c}}{\operatorname{argmax}} \sum_{i=1}^n \text{sim}(\hat{\mathbf{w}}(i), \mathbf{c}) = \begin{cases} \frac{\sum_{i=1}^n \hat{\mathbf{w}}(i) / \|\hat{\mathbf{w}}(i)\|}{\|\sum_{i=1}^n \hat{\mathbf{w}}(i) / \|\hat{\mathbf{w}}(i)\|\|} & \text{if } \sum_{i=1}^n \hat{\mathbf{w}}(i) / \|\hat{\mathbf{w}}(i)\| \neq 0 \\ 0 & \text{otherwise} \end{cases}. \quad (42)$$

After the convergence of the LBG clustering algorithm, we obtain  $E$  clusters of i-vectors with their centroids denoted as  $\mathbf{c}_1^{\hat{\mathbf{w}}}, \mathbf{c}_2^{\hat{\mathbf{w}}}, \dots, \mathbf{c}_E^{\hat{\mathbf{w}}}$ , respectively. We use  $\mathbf{c}_0^{\hat{\mathbf{w}}}$  to denote the centroid of all the training i-vectors.

## References

- [1] F. Beaufays, V. Vanhoucke, and B. Strope, "Unsupervised discovery and training of maximally dissimilar cluster models," *Proc. Interspeech2010*, pp. 66-69.
- [2] N. Dehak, P. J. Kenny, R. Dehak, P. Dumouchel, and P. Ouellet, "Front-end factor analysis for speaker verification," *IEEE Trans. on Audio, Speech and Language Processing*, Vol. 19, No. 4, pp. 788-798, 2011.
- [3] P. Kenny, G. Boulianne, and P. Dumouchel, "Eigenvoice modeling with sparse training data," *IEEE Trans. on Speech and Audio Processing*, Vol. 13, No. 3, pp. 345-354, 2005.
- [4] Y. Linde, A. Buzo, and R. M. Gray, "An algorithm for vector quantizer design," *IEEE Trans. on Communication*, Vol. COM-28, pp. 84-95, 1980.
- [5] J. Xu, Y. Zhang, Z.-J. Yan, and Q. Huo, "An i-Vector based Approach to Acoustic Sniffing for Irrelevant Variability Normalization based Acoustic Model Training and Speech Recognition," *Proc. Interspeech2011*.

- [6] S. Young, *et al.*, The HTK Book (for HTK version 3.4), 2006.
- [7] Y. Zhang, J. Xu, Z.-J. Yan, and Q. Huo, “An i-vector based approach to training data clustering for improved speech recognition,” *Proc. Interspeech2011*.

## A The Matrix Codebook for this memo

### A.1 Trace

$$\text{tr}(AB) = \text{tr}(BA) \quad (43)$$

$$\partial \text{tr}(X) = \text{tr}(\partial X) \quad (44)$$

$$\frac{\partial}{\partial X} \text{tr}(AX^\top) = A \quad (45)$$

$$\frac{\partial}{\partial X} \text{tr}(AXBX^\top) = AXB + A^\top XB^\top \quad (46)$$

$$\frac{\partial}{\partial X} \text{tr}(AX^{-1}B) = -X^{-\top} A^\top B^\top X^{-\top} \quad (47)$$

$$\frac{\partial}{\partial X} \ln |\det(X)| = (X^{-1})^\top \quad (48)$$

### A.2 Variance of quadratic form

Assume that  $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \mathbf{R})$  and  $A$  is positive semi-definite,

$$E(\mathbf{x}\mathbf{x}^\top) = \mathbf{R} + \boldsymbol{\mu}\boldsymbol{\mu}^\top \quad (49)$$

$$E(\mathbf{x}^\top A\mathbf{x}) = \text{tr}(A\mathbf{R}) + \boldsymbol{\mu}^\top A\boldsymbol{\mu} \quad (50)$$

PROOF:

$$\begin{aligned} \therefore & A \text{ is positive semi-definite} \\ \therefore & A = C^\top C \\ \therefore & E(C^\top \mathbf{x}^\top \mathbf{x} C) = C^\top (\mathbf{R} + \boldsymbol{\mu}\boldsymbol{\mu}^\top) C \\ \therefore & E(\mathbf{x}^\top A\mathbf{x}) = E\left(\sum_i (C\mathbf{x})_i^2\right) = \text{tr}(C^\top (\mathbf{R} + \boldsymbol{\mu}\boldsymbol{\mu}^\top) C) = \text{tr}(AE(\mathbf{x}\mathbf{x}^\top)) \end{aligned}$$

Note that

$$\begin{aligned} \text{tr}(A\boldsymbol{\mu}\boldsymbol{\mu}^\top) &= \sum_i \sum_j a_{ij} \mu_i \mu_j \\ &= \sum_i A\boldsymbol{\mu}\mu_i \\ &= \boldsymbol{\mu}^\top A\boldsymbol{\mu} \end{aligned} \quad (51)$$

and let  $\boldsymbol{\mu} = E(\mathbf{w})$ ,  $\mathbf{R} = E(\mathbf{w}\mathbf{w}^\top) - E(\mathbf{w})E(\mathbf{w})^\top$ ,  $A = \mathbf{T}^\top \Gamma(i) \mathbf{R}^{-1} \mathbf{T}$  we get Eq.[23].

### A.3 Others

#### 1 Proposition:

$$\mathbf{R}\Gamma(i) = \Gamma(i)\mathbf{R} \quad (52)$$

where

$$\mathbf{R} = \text{diag}\{\mathbf{R}_1, \dots, \mathbf{R}_C\} \quad (53)$$

$$\Gamma(i) = \text{diag}\{\gamma_1\mathbf{I}, \dots, \gamma_C\mathbf{I}\} \quad (54)$$

PROOF:

$$\begin{aligned} \mathbf{R}\Gamma(i) &= \text{diag}\{\mathbf{R}_1, \dots, \mathbf{R}_C\} \text{diag}\{\gamma_1\mathbf{I}, \dots, \gamma_C\mathbf{I}\} \\ &= \text{diag}\{\gamma_1\mathbf{I}\mathbf{R}_1, \dots, \gamma_C\mathbf{I}\mathbf{R}_C\} \\ &= \Gamma(i)\mathbf{R} \end{aligned} \quad (55)$$

In Mark J.F. Gales’s paper, it assume that  $\mathbf{w}$  is only a model parameter and find the solution through “the alternative of the variables method” (there is no closed-form maximum likelihood solution for  $\mathbf{w}, \mathbf{T}, \mathbf{R}$ ).

## B From a graphic model view

### B.1 Probabilistic PCA

Probabilistic PCA is a simple example of the linear-Gaussian framework, in which all of the marginal and conditional distributions are Gaussian. We can formulate probabilistic PCA by first introducing an explicit latent variable  $w$  corresponding to the principal-component subspace. Next we define a Gaussian prior distribution  $p(w)$  over the latent variable, together with a Gaussian conditional distribution  $p(x|w)$  for the observed variable  $x$  conditioned on the value of the latent variable. Specifically, the prior distribution over  $w$  is given by a zero-mean unit-covariance Gaussian

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \mathbf{I}) \quad (56)$$

and the conditional distribution of the observed variable  $x$

$$p(\mathbf{x}|\mathbf{w}) = \mathcal{N}(\mathbf{x}|T\mathbf{w} + \boldsymbol{\mu}, \sigma^2\mathbf{I}). \quad (57)$$

We can view the probabilistic PCA model from a generative viewpoint in which a sampled value of the observed variable is obtained by first choosing a value for latent variable and then sampling the observed variable conditioned on this latent value:

$$\mathbf{x} = T\mathbf{w} + \boldsymbol{\mu} + \boldsymbol{\epsilon} \quad (58)$$

where  $w$  is an Gaussian latent variable, and  $\boldsymbol{\epsilon}$  is a zero-mean Gaussian-distributed noise variable with covariance  $\sigma^2\mathbf{I}$ . The generative process is illustrated in Figure 2.

Figure 2: Probabilistic PCA

### B.1.1 Maximum likelihood PCA