

Analysis of the Influence of Speech Corpora in the PLDA Verification in the Task of Speaker Recognition

Lukáš Machlica and Zbyněk Zajíc

University of West Bohemia in Pilsen,
Faculty of Applied Sciences, Department of Cybernetics,
Univerzitní 22, 306 14 Pilsen
machlica@kky.zcu.cz, z Zajic@kky.zcu.cz
<http://www.kky.zcu.cz/en>

Abstract. In the paper recent methods used in the task of speaker recognition are presented. At first, the extraction of so called i-vectors from GMM based supervectors is discussed. These i-vectors are of low dimension and lie in a subspace denoted as Total Variability Space (TVS). The focus of the paper is put on Probabilistic Linear Discriminant Analysis (PLDA), which is used as a generative model in the TVS. The influence of development data is analyzed utilizing distinct speech corpora. It is shown that it is preferable to cluster available speech corpora to classes, train one PLDA model for each class and fuse the results at the end. Experiments are presented on NIST Speaker Recognition Evaluation (SRE) 2008 and NIST SRE 2010.

Key words: PLDA, latent space, fusion, supervector, FA, i-vector

1 Introduction

A major progress in the task of speaker recognition was done introducing supervector based techniques. Supervector is in fact a high dimensional feature vector obtained by the concatenation of lower dimensional vectors containing speaker dependent parameters – the most effective turned out to be the parameters related to Gaussian Mixture Models (GMMs) [1]. First attempts to incorporate supervectors into the speaker recognition task utilized Support Vector Machines (SVMs) along with distinct kernel functions [2]. Since GMMs belong to the class of generative models, whereas SVMs are based on the discrimination between classes, the techniques comprising both methods are also known as *hybrid modelling*. Subsequently, additional techniques were added to solve the problem of the change of operating conditions, namely Nuisance Attribute Projection (NAP) [3]. NAP is based on an orthogonal projection, where directions most vulnerable to environment changes are projected out. Since supervectors are of substantially high dimension (tens of thousands), which is often higher than the number of supervectors provided for the training, it is obvious that a

lot of dimensions will be correlated with each other, and that the information about the identity will be contained in a subspace of a much lower dimension. This idea was incorporated into the principle of Joint Factor Analysis (JFA) [4], where the word *joint* refers to the fact that not only the speaker, but also the channel variabilities are treated in one JFA model. However, since experiments in [5] showed that the channel space obtained by JFA does still contain some information about the speaker’s identity, JFA was slightly adjusted giving rise to *identity vectors*, or *i-vectors* [6]. The main difference between JFA and i-vectors is that i-vectors do not distinguish between speaker and channel space. They work with a *total variability space* containing simultaneously speaker and channel variabilities, whereas JFA treats both spaces individually.

Parallel to JFA a very similar approach was introduced in the image recognition called Probabilistic Linear Discriminant Analysis (PLDA) [7]. The only difference from JFA is that in PLDA ordinary feature vectors are used instead of GMM based supervectors (for details on the treatment of supervectors in JFA see [4]). Since PLDA is a generative model, it allows to compute the probability that several i-vectors originate from the same source, and thus it is well suited as a verification tool for a speaker recognition system [8]. System presented in this paper will utilize i-vectors (described in Section 3) based on GMM supervectors (see Section 2) with a PLDA model (refer to Section 4) used in the verification phase.

The crucial problem when proposing a speaker verification system composed of modules (e.g. JFA, PLDA) described above is that data from a lot of speakers are required, moreover several sessions have to be available for each speaker in order to train a reliable i-vector extractor and a PLDA model. The problem faced in this paper will address the question whether distinct speech corpora (e.g. Switchboard 1, Switchboard 2, NIST SRE 2004, NIST SRE 2006, etc.) should be pooled together and used to train one PLDA model, or if each corpus should be used individually to train a separate PLDA model. In the latter scenario the results are fused at the end. Experiments can be found in Section 5.

2 Supervector Extraction based on GMMs

At first a Universal Background Model (UBM) has to be trained. UBM is in fact a Gaussian Mixture Model (GMM), however it is trained from a set containing a lot of speakers. The speakers data should match all the conditions, in which the recognition system is going to be used. UBM consists of a set of parameters $\lambda_{\text{UBM}} = \{\omega_m, \boldsymbol{\mu}_m, \mathbf{C}_m\}_{m=1}^M$, where M is the number of Gaussians in the UBM, ω_m , $\boldsymbol{\mu}_m$, \mathbf{C}_m are the weight, mean and covariance of the m^{th} Gaussian, respectively. Let $\mathbf{O}_s = \{\mathbf{o}_{st}\}_{t=1}^{T_s}$ be the set of T_s feature vectors \mathbf{o}_{st} of dimension D belonging to the s^{th} speaker, and

$$\gamma_m(\mathbf{o}_{st}) = \frac{\omega_m \mathcal{N}(\mathbf{o}_{st}; \boldsymbol{\mu}_m, \mathbf{C}_m)}{\sum_{m=1}^M \omega_m \mathcal{N}(\mathbf{o}_{st}; \boldsymbol{\mu}_m, \mathbf{C}_m)} \quad (1)$$

be the posterior probability of m^{th} Gaussian given a feature vector \mathbf{o}_{st} . And let $\mathbf{m}_0 = [\boldsymbol{\mu}_1^T, \dots, \boldsymbol{\mu}_m^T, \dots, \boldsymbol{\mu}_M^T]^T$ be the supervector composed of UBM means. Then, for each speaker two supervectors are extracted

$$\begin{aligned} \mathbf{n}_s &= \sum_{t=1}^{T_s} \left([\gamma_1(\mathbf{o}_{st}), \dots, \gamma_m(\mathbf{o}_{st}), \dots, \gamma_M(\mathbf{o}_{st})]^T \otimes \mathbf{1}_D \right) \text{ of size } DM \times 1, \\ \mathbf{b}_s &= \sum_{t=1}^{T_s} [\gamma_1(\mathbf{o}_{st})\mathbf{o}_{st}^T, \dots, \gamma_m(\mathbf{o}_{st})\mathbf{o}_{st}^T, \dots, \gamma_M(\mathbf{o}_{st})\mathbf{o}_{st}^T]^T \text{ of size } DM \times 1, \end{aligned} \quad (2)$$

where \otimes is the Kronecker product, and $\mathbf{1}_D$ is a D dimensional vector of ones. Note that \mathbf{n}_s is the supervector containing "soft" counts of feature vectors aligned to Gaussians $1, \dots, M$, and denoting \mathbf{N}_s a diagonal matrix containing \mathbf{n}_s on its diagonal, $\mathbf{m}_s = \mathbf{N}_s^{-1}\mathbf{b}_s$ is the new Maximum Likelihood (ML) estimate of supervector \mathbf{m}_0 given the dataset \mathbf{O}_s . At last note that the Maximum Aposteriori Probability (MAP) adaptation [9] of means of the UBM according to the given data set \mathbf{O}_s expressed in the supervector notation is given as $\mathbf{m}_{\text{MAP}} = \tau\mathbf{m}_s + (1 - \tau)\mathbf{m}_0$ for some relevance factor τ .

3 i-Vector Extraction

The concept of the i-vectors extraction is based on Factor Analysis (FA) extended to handle session and speaker variabilities of supervectors to Joint Factor Analysis (JFA) [4]. Contrary to JFA, different sessions of the same speaker are considered to be produced by different speakers [5]. The generative i-vector model has the form

$$\boldsymbol{\psi}_s = \mathbf{m}_0 + \mathbf{T}\mathbf{w}_s + \boldsymbol{\epsilon}, \quad \mathbf{w}_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}) \quad (3)$$

where \mathbf{T} (of size $D \times D_w$) is called the total variability space matrix (it contains both the variabilities between speakers and the channel variabilities between distinct sessions of a speaker), \mathbf{w}_s is the s^{th} speaker's i-vector of dimension D_w having standard Gaussian distribution, \mathbf{m}_0 is the mean vector of $\boldsymbol{\psi}_s$, however often the UBM's mean supervector is taken instead as a good approximation (therefore the same notation \mathbf{m}_0 is used), and $\boldsymbol{\epsilon}$ is some residual noise with a diagonal covariance $\boldsymbol{\Sigma}$ constructed from covariance matrices $\mathbf{C}_1, \dots, \mathbf{C}_m$ of the UBM ordered on the diagonal of $\boldsymbol{\Sigma}$.

To train the matrix \mathbf{T} two steps are iterated in a sequence. Given a training set of S couples of supervectors $\mathbf{b}_s, \mathbf{n}_s$, and the diagonal matrix \mathbf{N}_s containing \mathbf{n}_s on its diagonal, these steps are:

1. use previous estimate of \mathbf{T} to extract new i-vectors for all speakers $1, \dots, S$

$$\mathbf{w}_s = (\mathbf{I} + \mathbf{T}^T \boldsymbol{\Sigma}^{-1} \mathbf{N}_s \mathbf{T})^{-1} \mathbf{T}^T \boldsymbol{\Sigma}^{-1} \bar{\mathbf{b}}_s, \quad (4)$$

2. according to the new i-vectors compute block-wise a new estimate of \mathbf{T}

$$\mathbf{T}_m = \left(\sum_{s=1}^S \bar{\mathbf{b}}_{sm} \mathbf{w}_s^T \right) \left(\sum_{s=1}^S N_{sm} \left(\mathbf{w}_s \mathbf{w}_s^T + (\mathbf{I} + \mathbf{T}^T \boldsymbol{\Sigma}^{-1} N_s \mathbf{T})^{-1} \right) \right)^{-1}, \quad (5)$$

where $\bar{\mathbf{b}}_s = \mathbf{b}_s - N_s \mathbf{m}_0$ is the centered version of \mathbf{b}_s around the mean \mathbf{m}_0 , and the index m in \mathbf{T}_m , $\bar{\mathbf{b}}_{sm}$, \mathbf{n}_{sm} (and N_{sm}) refers to blocks of \mathbf{T} , $\bar{\mathbf{b}}_s$, \mathbf{n}_s (and thus to N_{sm}) of sizes $D \times D_w$, $D \times 1$, $D \times 1$ so that $\mathbf{T}^T = [\mathbf{T}_1^T, \mathbf{T}_2^T, \dots, \mathbf{T}_M^T]$, $\bar{\mathbf{b}}_s^T = [\bar{\mathbf{b}}_{s1}^T, \bar{\mathbf{b}}_{s2}^T, \dots, \bar{\mathbf{b}}_{sM}^T]$, $\mathbf{n}_s^T = [\mathbf{n}_{s1}^T, \mathbf{n}_{s2}^T, \dots, \mathbf{n}_{sM}^T]$, respectively. In fact, also $\boldsymbol{\Sigma}$ may be updated in each iteration, for details see [10].

4 Probabilistic Linear Discriminant Analysis (PLDA)

Let us assume that the i-vector extractor (4) was already trained, and that for each feature set \mathbf{O}_s of a speaker s one i-vector \mathbf{w}_s was extracted. Further, let us assume that several sessions $\{\mathbf{O}_{sh}\}_{h=1}^{H_s}$ of a speaker s are available, and that for each set of feature vectors \mathbf{O}_{sh} of each session $h = 1, \dots, H_s$ one i-vector \mathbf{w}_{sh} was extracted. Since in the i-vector extraction phase no distinction between session space and speaker space were made a new model in the total variability space will be now described that is going to utilize also the session variabilities.

PLDA is a generative model of the form

$$\mathbf{w}_{sh} = \mathbf{m}_w + \mathbf{F} \mathbf{z}_s + \mathbf{G} \mathbf{r}_{sh} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{S}) \quad (6)$$

where \mathbf{m}_w is the mean of \mathbf{w}_{sh} , columns of \mathbf{F} span the speaker identity space, \mathbf{z}_s of dimension D_z are coordinates in this space and they do not change across sessions of one speaker, columns of \mathbf{G} span the channel space, \mathbf{r}_{sh} of dimension D_r are the session dependent speaker factors, and $\boldsymbol{\epsilon}$ is some residual noise with diagonal covariance \mathbf{S} and a zero mean. Further restrictions are put on distributions of latent variables \mathbf{z}_s and \mathbf{r}_{sh} , namely that both follow standard Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Hence, $\mathbf{w}_{sh} \sim \mathcal{N}(\mathbf{m}_w, \mathbf{F} \mathbf{F}^T + \mathbf{G} \mathbf{G}^T + \mathbf{S})$. It is common and reasonable assumption that $D_z \ll D_w$ and that $D_z + D_r \approx D_w$. To train the model parameters \mathbf{F} , \mathbf{G} and \mathbf{S} one has to solve the system of equations [7]

$$\begin{bmatrix} \mathbf{w}_{s1} - \mathbf{m}_w \\ \mathbf{w}_{s2} - \mathbf{m}_w \\ \vdots \\ \mathbf{w}_{sH_s} - \mathbf{m}_w \end{bmatrix} = \begin{bmatrix} \mathbf{F} & \mathbf{G} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{F} & \mathbf{0} & \mathbf{G} & \dots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{F} & \mathbf{0} & \mathbf{0} & \dots & \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{z}_s \\ \mathbf{r}_{s1} \\ \mathbf{r}_{s2} \\ \vdots \\ \mathbf{r}_{sH_s} \end{bmatrix} + \begin{bmatrix} \boldsymbol{\epsilon}_1 \\ \boldsymbol{\epsilon}_2 \\ \vdots \\ \boldsymbol{\epsilon}_{H_s} \end{bmatrix}, \quad \hat{\boldsymbol{\Sigma}}_{H_s} = \begin{bmatrix} \mathbf{S} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{S} & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{S} \end{bmatrix}, \quad (7)$$

and when rewritten to a compact form we get

$$\hat{\mathbf{w}}_s = \mathbf{A}_{H_s} \hat{\mathbf{z}}_s + \hat{\boldsymbol{\epsilon}}, \quad (8)$$

where $\hat{\boldsymbol{\epsilon}} \sim \mathcal{N}(\mathbf{0}, \hat{\boldsymbol{\Sigma}}_{H_s})$. Matrices \mathbf{A}_{H_s} , $\hat{\boldsymbol{\Sigma}}_{H_s}$ depend on s through the number of their row- and column-blocks given by the number of sessions H_s of speaker s .

Problem (8) is a standard FA problem, for details on how to solve it see the appendix in [7].

4.1 Verification

In the verification phase two hypotheses are tested [7], namely

- hypotheses \mathcal{H}_s that two i-vectors \mathbf{w}_1 and \mathbf{w}_2 share the same identity,
- hypotheses \mathcal{H}_d that the identity of two i-vectors \mathbf{w}_1 and \mathbf{w}_2 differs.

The log-likelihood ratio is given as

$$\begin{aligned} \text{LLR}(\mathbf{w}_1, \mathbf{w}_2) &= \log \frac{p(\mathbf{w}_1, \mathbf{w}_2 | \mathcal{H}_s)}{p(\mathbf{w}_1 | \mathcal{H}_d)p(\mathbf{w}_2 | \mathcal{H}_d)} = \\ &= \log \mathcal{N} \left(\begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix}; \begin{bmatrix} \mathbf{m}_w \\ \mathbf{m}_w \end{bmatrix}, \begin{bmatrix} \mathbf{C}_w & \mathbf{C}_F \\ \mathbf{C}_F & \mathbf{C}_w \end{bmatrix} \right) - \log \mathcal{N} \left(\begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix}; \begin{bmatrix} \mathbf{m}_w \\ \mathbf{m}_w \end{bmatrix}, \begin{bmatrix} \mathbf{C}_w & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_w \end{bmatrix} \right), \end{aligned} \quad (9)$$

where $\mathbf{C}_w = \mathbf{F}\mathbf{F}^T + \mathbf{G}\mathbf{G}^T + \mathbf{S}$ and $\mathbf{C}_F = \mathbf{F}\mathbf{F}^T$. Note that in this verification scenario we do not care about the form of the decomposition of \mathbf{w}_1 or \mathbf{w}_2 (latent variables \mathbf{z}_s , \mathbf{r}_{sh} stay unknown). The question stated is whether two vectors share the same identity given the subspaces generated by \mathbf{F} and \mathbf{G} .

5 Experiments

The question raised is whether all the available data from several distinct corpora should be pooled and used to train one PLDA model (thus find a decomposition of the total variability space based on all the available data), or if it would be more efficient to find a characteristic decomposition of the total variability space for each corpus individually (hence train several PLDA models), score each pair of vectors in relation to each total space decomposition, and finally fuse the obtained scores (we will use linear combination). We believe that the latter case makes the verification more robust since possible undesirable deviations in acoustic conditions of distinct corpora may become less evident. However, individual corpora still have to contain enough data to be able to train a reliable PLDA model.

5.1 Used Corpora

In order to be able to perform reliable tests we utilized corpora: NIST SRE 2004, NIST SRE 2005, NIST SRE 2006, Switchboard 1 Release 2 and Switchboard 2 Phase 3 for development purposes, and NIST SRE 2008, NIST SRE 2010 were used for calibration of Fusion Coefficients (FCs), and for the evaluation of generality of obtained FCs, respectively. We used only those speakers from development corpora who had more than 4 recorded sessions. Further, the development corpora were divided into 3 classes:

1. NIST040506 – containing 3787 recordings of 465 males of approximately 8 sessions for each male speaker,
2. SW1 – containing 2342 recordings of 211 males of approximately 11 sessions for each male speaker,
3. SW2 – containing 2183 recordings of 216 males of approximately 10 sessions for each male speaker,

Each of the recordings had approximately 5 minutes in duration including the silence. The division into the classes was made in relation to the similarity of corpora determined according to recording conditions given in the LDC Corpus Catalog¹.

In order to train the FCs "short2-short3 trials" from NIST SRE 2008 [11] were utilized, only telephone speech from males was used (648 target speakers and 1535 test speakers) yielding 16968 trials in total. To test the validity of learned FCs "core-core trials" from NIST SRE 2010 [12] were used, and again only telephone speech from males was used (1394 target speakers and 2474 test speakers) yielding 74762 trials in total. The duration of all the test and target recordings in both corpora was approximately 5 minutes including the silence.

5.2 Feature Extraction

The feature extraction was based on Linear Frequency Cepstral Coefficients (LFCCs), Hamming window of length 25 ms was used, the shift of the window was set to 10 ms. 25 triangular filter banks were spread linearly across the frequency spectrum, and 20 LFCCs were extracted, delta coefficients were added leading to a 40 dimensional feature vector. Also the Feature Warping (FW) normalization procedure was applied utilizing a sliding window of length 3 seconds. Right before the FW Voice Activity Detector (VAD), based on detection of energies in filter banks located in the frequency domain, was used in order to discard the non-speech frames. All the feature vectors were at the end down-sampled by a factor of 2.

The number of Gaussians in the UBM was set to 1024. The size of the total variability space matrix \mathbf{T} in the i-vector extraction was set to $1024 * 40 \times 800$, thus the latent dimension (dimension of i-vectors) was $D_w = 800$. At last, the dimension of the speaker identity space in the PLDA model was set to $D_z = 100$ and the dimension of the session/channel space was set to $D_r = 800$, thus \mathbf{F} was of size 800×100 , and the channel matrix was a square matrix of size 800×800 . The disproportion between dimensions of speaker and channel subspaces was adopted from [8].

5.3 Results and Analysis

UBM and the i-vector's extractor described in Section 3 were trained on the pooled dataset NIST040506 + SW1 + SW2. Next, three PLDA models were

¹ <http://www ldc upenn edu/Catalog/index.jsp>

trained utilizing subsequently each of the 3 corpora classes. Trials from NIST SRE 2008 were then scored using all 3 PLDA models, and the scores were used to train the fusion coefficients via the linear logistic regression from the FoCal toolkit [13]. Finally, in order to test the validity of learned FCs the same approach was performed with trials from NIST SRE 2010, but the already learned FCs were used in the linear combination of obtained scores. Results are shown in Figure 1 and Table 1, also minimum of the Decision Cost Function (DCF) is reported. In order to compute the value of DCF the cost of missing a target was set to 10, the cost of the false alarm was set to 1, and the probability of seeing a true trial was set to 0.01. These values are adopted from the NIST Speaker Recognition Evaluation (SRE) 2008 [11].

We have trained one PLDA model also from pooled corpora NIST040506 + SW1 + SW2 (this was not used in the fusion). Best results are obtained for the fused system in both NIST SRE 2008 and NIST SRE 2010. Note that PLDA trained only on SW2 outperforms all the other PLDA models trained on other corpora (even on the pooled corpora), but the fusion still increases the performance of the speaker verification system. However, in real conditions one can not rely only on one corpus (in this case it would be SW2) performing best on the development set.

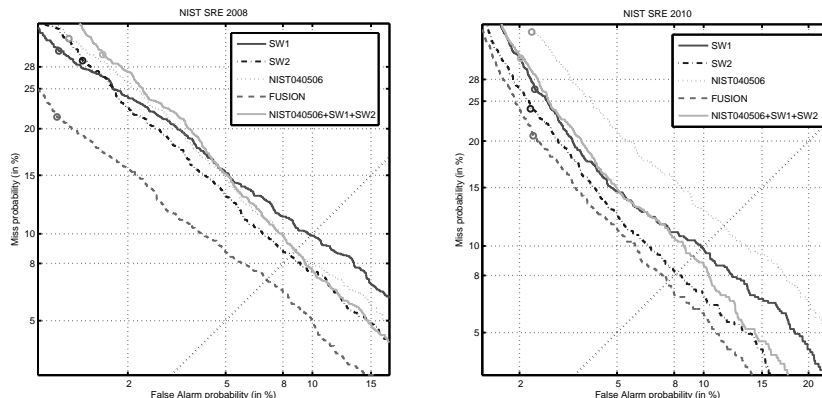


Fig. 1. DET curves for NIST SRE 2008 and NIST SRE 2010. Circles denote points where minDCF occurred.

Table 1. Results are given as EER [%] / minDCF. In the last column of the table also results for PLDA trained on pooled corpora NIST040506 + SW1 + SW2 is given.

	SW1	SW2	NIST040506	FUSION	pooled
NIST 2008	9.87/0.034	8.47/0.041	8.83/0.043	7.06/0.031	8.78/0.045
NIST 2010	9.89/0.050	8.11/0.046	11.58/0.057	7.58/0.043	9.26/0.051

6 Conclusion

Since often the verification conditions are unknown in advance (e.g. in the Speaker Recognition Evaluations (SREs) organized by NIST and other institutions) we cannot count on the use of one specific speech corpus performing best on the development set. It is more convenient to utilize several corpora. We have shown that if the utilized corpora have sufficient amount of data to train reliable PLDA models, it is preferable to train several PLDA models and fuse the results. The verification becomes more robust since the deviations in acoustic conditions of distinct corpora become less evident.

Acknowledgments. This research was supported by the grant of the University of West Bohemia, project No. SGS-2010-054.

References

1. Campbell, W., Sturim, D., Reynolds, D.: Support Vector Machines Using GMM Supervectors for Speaker Verification. In: IEEE Signal Processing Letters, vol. 13, pp. 308311 (2006)
2. Longworth, C., Gales, M.: Parametric and Derivative Kernels for Speaker Verification. In: Interspeech 2007 pp. 310313 (2007)
3. Solomonoff, A., Quillen, C., Campbell, W.: Channel compensation for SVM speaker recognition. In: Odyssey, pp. 5762 (2004)
4. Kenny, P.: Joint Factor Analysis of Speaker and Session Variability: Theory and Algorithms. Tech. report, Centre de Recherche Informatique de Montral (2006)
5. Dehak, N.: Discriminative and Generative Approaches for Long- and Short-term Speaker Characteristics Modeling: Application to Speaker Verification. Ph.D. thesis, École de Technologie Supérieure, Université du Québec (2009)
6. Dehak, N., Kenny, P., Dehak, R., Dumouchel, P., Ouellet, P.: Front-End Factor Analysis For Speaker Verification. In: IEEE Transactions on Audio, Speech and Language Processing (2010)
7. Prince, S., Elder, J.: Probabilistic Linear Discriminant Analysis for Inferences About Identity. In: IEEE 11th International Conference on Computer Vision, pp. 18 (2007)
8. Matějka, P., Glembek, O., Castaldo, F., Alam, J., Plchot, O., Kenny, P., Burget, L., Černocký, J.: Full-covariance UBM and Heavy-tailed PLDA in I-Vector Speaker Verification. In: ICASSP 2011, pp. 48284831 (2011)
9. Reynolds, D.A., Quatieri, T.F., Dunn, R.B.: Speaker Verification Using Adapted Gaussian Mixture Models. In: Digital Signal Processing, vol: 10, pp. 19-41 (2000)
10. Patrick, K., Pierre, O., Najim, D., Vishwa, G., Pierre, D.: A Study of Interspeaker Variability in Speaker Verification. In: IEEE Transactions on Audio, Speech and Language Processing, vol. 16, pp. 980-988 (2008)
11. The NIST Year 2008 Speaker Recognition Evaluation Plan, http://www.itl.nist.gov/iad/mig/tests/spk/2008/sre08_evalplan_release4.pdf
12. The NIST Year 2010 Speaker Recognition Evaluation Plan, http://www.itl.nist.gov/iad/mig/tests/spk/2010/NIST_SRE10_evalplan.r6.pdf
13. Brummer, N.: FoCal: Tools for fusion and calibration of automatic speaker detection systems (2006). <http://sites.google.com/site/nikobrummer/focal>