

Fisher Vectors in PLDA Speaker Verification System

Zbyněk Zajíc, Marek Hruží

University of West Bohemia

Faculty of Applied Sciences

NTIS - New Technologies for the Information Society

Univerzitní 8, 306 14 Plzeň, Czech Republic

Email: {zzajic, mhruz}@ntis.zcu.cz

Abstract—The goal of this paper is to examine the Fisher Vector and incorporate this vector in the PLDA based speaker verification system. The PLDA based system utilizes the Supervector of Statistics extracted from a Gaussian Mixture Model (adopted from the speaker adaptation task) to collect the information about a speaker from a dataset. We compare the efficiency of the PLDA based speaker verification system using Supervector of Statistics and the same system with Fisher vector. The experimental results of these two approaches to the verification task and the fusion of these two systems indicate that the Fisher Vector brings almost the same information to the PLDA verification process as the Supervector of Statistics when sufficient data are available.

Index Terms—Speaker Verification, PLDA, Fisher Vector, Supervector, iVector.

I. INTRODUCTION

General model based on Probabilistic Linear Discriminant Analysis (PLDA) [1] and EigenVectors (EVs) descriptors [2], used originally in image processing for face recognition, was successfully integrated into speaker verification system which is nowadays considered as a state-of-the-art approach [3]. A new method for automatic face verification utilizing Fisher Vectors (FVs) as high-dimensional descriptors was introduced in [4]. In this paper, we introduced FVs in a speaker verification system and compared it with a system that works with Gaussian Mixture Model (GMM) based supervectors of statistics [5].

Supervector is, in fact, a high-dimensional feature vector obtained by the concatenation of lower-dimensional vectors containing speaker dependent parameters - in our case the first and zeroth statistical moments of speaker data related to a Universal Background Model (UBM) based on GMM [6]. This Supervector of Statistics can be seen as the new Maximum Likelihood (ML) estimate of speaker identity and has roots in the task of speaker adaptation [7]. On the other hand, FV is based on the Fisher Information which measures the amount of information that an observable random variable O carries about an unknown parameter of a distribution that models O .

This paper is organized as follows: The PLDA based speaker verification system is described in Section II, where in Subsection II-A the Gaussian Mixture Model (GMM) based Supervector of Statistics is described. In Section III the FV and the replacement of the Supervector of Statistics in the verification process is introduced. The results of two speaker

verification systems based on Supervector of Statistics and the Fisher Vectors and the fusion of both system can be found in Section IV.

II. iVECTORS-PLDA SYSTEM

The iVector-PLDA framework, a state-of-the-art system for speaker verification [8], is based on the extraction of features from the speech and accumulation of the statistics of these features into supervectors. This supervector is high-dimensional (tens of thousands), and hence it is suitable to find a latent space of a much lower dimension which represents the speakers. Vectors from this space are called iVectors. The iVectors extraction (the dimensionality reduction of the supervector) is based on Factor Analysis (FA).

Two iVectors can be compared with each other using cosine distance. However, iVectors obtained by FA still contain some noisy information not relevant to the speaker identity (e.g. influence of the channel). Therefore a PLDA model trained on a huge amount of structured data (several representations of each speaker from different sources - sessions) is used for decomposing information from iVector into the speaker and session domain. Then, only the speaker domain is used for comparison of two speaker representations. Moreover, PLDA model itself can be used as a powerful tool for identity verification instead of cosine distance [9]. A diagram of this verification system can be seen in Figure 1. Each step is described in detail in the following subsections.

A. Statistics Extracted on GMM

Supervector of Statistics containing the first and zeroth statistical moments of speakers' data related to UBM has origins in the speaker adaptation process, where these statistics are used as a descriptor of a new speaker.

First, a GMM trained on a huge amount of data from different speakers is used as a UBM and consists of a set of parameters $\lambda_{\text{UBM}} = \{\omega_m, \mu_m, C_m\}_{m=1}^M$, where M is the number of Gaussians in the UBM, ω_m , μ_m , C_m are the weight, mean and covariance of the m^{th} Gaussian, respectively. In our case, the covariance matrix C_m is diagonal with vector σ_m on diagonal.

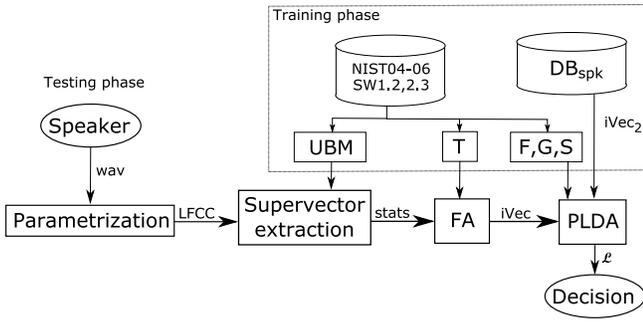


Fig. 1. Diagram of the verification process with four steps: parametrization, suprvector extraction, dimensionality reduction by FA and finally verification with another speaker (represented as iVector) using PLDA model. The output from the PLDA system is likelihood (\mathcal{L}) that two iVectors (new constructed iVec and iVec₂ from databases DB_{spk}) representing the same speaker.

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} . The soft count of the m^{th} Gaussian (zeroth statistical moments of feature vectors) is

$$n_m^s = \sum_{t=1}^{T_s} \gamma_m(\mathbf{o}_{st}), \quad (2)$$

and the sum of the first statistical moments of feature vectors with respect to the m^{th} Gaussian is

$$\mathbf{b}_m^s = \sum_{t=1}^{T_s} \gamma_m(\mathbf{o}_{st}) \mathbf{o}_{st}. \quad (3)$$

The speaker's suprvector for given data \mathbf{O}_s is a concatenation of the zeroth and first statistical moments of \mathbf{O}_s .

Note: the origin of these statistics can be seen if we rearrange the zeroth and first statistics into partial suprvecors (of size $DM \times 1$):

$$\begin{aligned} \mathbf{n}_s &= \sum_{t=1}^{T_s} \left([\gamma_1(\mathbf{o}_{st}), \dots, \gamma_m(\mathbf{o}_{st}), \dots, \gamma_M(\mathbf{o}_{st})]^T \otimes \mathbf{1}_D \right) \\ \mathbf{b}_s &= \sum_{t=1}^{T_s} [\gamma_1(\mathbf{o}_{st}) \mathbf{o}_{st}^T, \dots, \gamma_m(\mathbf{o}_{st}) \mathbf{o}_{st}^T, \dots, \gamma_M(\mathbf{o}_{st}) \mathbf{o}_{st}^T]^T \end{aligned} \quad (4)$$

where \otimes is the Kronecker product, and $\mathbf{1}_D$ is a D -dimensional vector of ones. If we denote \mathbf{N}_s a diagonal matrix containing \mathbf{n}_s as its diagonal, then

$$\mathbf{m}_s = \mathbf{N}_s^{-1} \mathbf{b}_s \quad (5)$$

can be seen as a new Maximum Likelihood (ML) estimation of $\mathbf{m}_0 = [\boldsymbol{\mu}_1^T, \dots, \boldsymbol{\mu}_m^T, \dots, \boldsymbol{\mu}_M^T]^T$ (a suprvector composed of UBM means) for given \mathbf{O}_s . The Maximum Aposteriori

Probability (MAP) adaptation [10] of UBM means (according to \mathbf{O}_s) is given by

$$\mathbf{m}_{\text{MAP}} = \tau \mathbf{m}_s + (1 - \tau) \mathbf{m}_0, \quad (6)$$

where τ is an empirically determined factor of data relevance.

B. iVectors extraction

For iVectors extraction the Factor Analysis (FA) approach [11] (or extended Joint Factor Analysis (JFA) [3] to handle more sessions of each speaker) is used for dimensionality reduction of the suprvector. The generative iVector model has the form

$$\boldsymbol{\psi}_s = \mathbf{m}_0 + \mathbf{T} \mathbf{w}_s + \boldsymbol{\epsilon}, \quad \mathbf{w}_s \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}) \quad (7)$$

where \mathbf{T} (of size $D \times D_w$) is called the total variability space matrix, \mathbf{w}_s is the s^{th} speaker's iVector 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 suprvector \mathbf{m}_0 is taken instead as a good approximation, 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 $\boldsymbol{\Sigma}$. The iVectors are also length-normalised [12]. Details about training of total variability space matrix \mathbf{T} can be seen in [13] or [14].

C. Probabilistic Linear Discriminant Analysis (PLDA)

In the iVector extraction phase by FA, no distinction between session space and speaker space were made (in contrast with JFA). If structured training data (more than one session - source for each speaker) are available, PLDA can be trained to model speaker and session variability separately. PLDA is a generative model [3] of the form:

$$\mathbf{w}_{sh} = \mathbf{m}_w + \mathbf{F} \mathbf{z}_s + \mathbf{G} \mathbf{r}_{sh} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{S}) \quad (8)$$

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 placed on distributions of latent variables \mathbf{z}_s and \mathbf{r}_{sh} , namely that both follow a 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 a common and reasonable assumption that $D_z \ll D_w$ and that $D_z + D_r \approx D_w$. To train the PLDA model parameters \mathbf{F} , \mathbf{G} and \mathbf{S} the system of equations must be solved [15] which leads to the standard FA problem (for more details see [1]).

III. FISHER VECTORS

Fisher Vectors based on Fisher Kernel [16] and Fisher information (which measures the amount of information that an observable random variable \mathbf{O} carries about an unknown parameter of a distribution that models \mathbf{O}) was recently used in face recognition [4] as an effective encoding of the feature space structure. If we assume that \mathbf{O} can be modeled by a

probability density function u_λ with parameter λ , then \mathbf{O} can be described by the gradient vector [17]:

$$\mathbf{G}_\lambda^O = \frac{1}{T} \nabla_\lambda \log(u_\lambda(\mathbf{O})), \quad (9)$$

where \mathbf{G}_λ^O describes the contribution of these parameters to the generation process. A natural kernel on these gradients is

$$K(\mathbf{O}, \mathbf{Q}) = \mathbf{G}_\lambda^{O^T} \mathbf{F}_\lambda \mathbf{G}_\lambda^Q, \quad (10)$$

where \mathbf{F}_λ is the Fisher information matrix defined as:

$$\mathbf{F}_\lambda = E_{\mathbf{o} \approx u_\lambda} [\nabla_\lambda \log(u_\lambda(\mathbf{o})) \nabla_\lambda \log(u_\lambda(\mathbf{O})^T)] \quad (11)$$

and this matrix has a Cholesky decomposition $\mathbf{F}_\lambda = \mathbf{L}_\lambda^T \mathbf{L}_\lambda$. From this, Fisher vector of \mathbf{O} can be defined as:

$$\phi_\lambda^O = \mathbf{L}_\lambda \mathbf{G}_\lambda^O. \quad (12)$$

We consider gradient with respect to the parameters of the UBM $\lambda_{UBM} = \{\omega_m, \mu_m, \mathbf{C}_m\}_{m=1}^M$. FVs encoding aggregates a large set of vectors into a high-dimensional supervector representation by fitting UBM to the features \mathbf{O} and encoding the derivatives of the log-likelihood of UBM. This representation captures the average first (and possibly second) order differences between the features and UBM components:

$$\phi_m^s = \frac{1}{N\sqrt{\omega_m}} \sum_{t=1}^T \gamma_m(\mathbf{o}_{st}) \left(\frac{\mathbf{o}_{st} - \mu_m}{\sigma_m} \right), \quad (13)$$

$$\phi_m^{s^2} = \frac{1}{N\sqrt{2\omega_m}} \sum_{t=1}^T \gamma_m(\mathbf{o}_{st}) \left(\frac{(\mathbf{o}_{st} - \mu_m)^2}{\sigma_m^2} - 1 \right), \quad (14)$$

where $\phi_m^{s^2}$ is the average second order differences of data \mathbf{O}_s dependent on UBM model.

FV is obtained by concatenating the differences of all UBM components into one supervector for each speaker s . In this paper, we use FV constructed only from the first-order differences (13) and soft count of occurrences (2). The goal of this paper is the comparison of the efficiency of the verification system with FVs and the system with Supervector of Statistics (where only zeroth and first moment is used in general). The dimensionality of FV is $M * (D + 1)$, where M is the number of components in UBM and D is the dimensionality of the feature vector \mathbf{o}_{st} .

IV. EXPERIMENTS

In this paper, we try to answer the question if the FVs can bring new information to the speaker verification system compared to the system which uses the Supervector of Statistics. The experiment was carried out on the Czech telephone corpus (cell phone or fixed line) consisting of 2005 speakers each with 2-4 min for training and 2-4 min for testing phase including the silence (which can be considered as sufficient amount of data in speaker recognition task). From all possible 4020025 trials, 10% was used for training the fusion coefficients via the linear logistic regression from the FoCal toolkit [18]. The rest was used for evaluating the verification systems.

TABLE I

COMPARISON OF THE STATE-OF-THE-ART SYSTEM USING SUPERVECTORS OF STATISTICS OR FISHERS VECTORS AND A COMBINATION OF THESE TWO SYSTEMS. RESULTS ARE GIVEN AS GIVEN AS EER [%] AND minDCF.

system	EER	minDCF
statistics	4.69%	0.3393
fisher vectors	6.23%	0.4789
combination	4.59%	0.3345

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

Speaker verification PLDA based system was trained using corpora: NIST SRE 2004, NIST SRE 2005, NIST SRE 2006, Switchboard 1 Release 2 and Switchboard 2 Phase 3. The number of Gaussians in the UBM was set to 512. The latent dimension (dimension of iVectors) in the the FA total variability space matrix \mathbf{T} in the iVector extraction was set to 400. At last, the dimension of the speaker identity space in the PLDA model was set to 200 and the dimension of the session/channel space was set to 400.

A. Results

The metrics for evaluation are Equal Error Rate (EER) and the Minimum Decision Cost Function (minDCF) [19]. The results are shown as Detection Error Tradeoff (DET) curve [20] in Figure 2 and in Table I.

Recently published paper [21] reported comparison on a similar verification system on NIST2010 with different results. They obtain slightly better results for system with Fisher Vectors than supervectors of statistics and the fusion of these two systems (FV and S) brings improvement. In their paper, authors used iVector model approach to speaker verification system with Linear discriminant analysis (LDA) and Within-Class Covariance Normalization (WCCN). The Fisher Vector consist of the first and second order moments while Supervector of Statistics doesn't. In our paper we used PLDA model instead and only the first and zeroth statistical moments for supervector (in both cases: Fisher Vectors and Supervector of Statistics). We assume that the comparison of supervectors of statistics and Fisher Vectors attribution to the verification system is more comparable if both contain the same amount of information (only zeroth and first moments).

The experimental results of these two approaches to the verification task and the fusion of these two systems indicates that the Fisher Vector brings almost the same information to

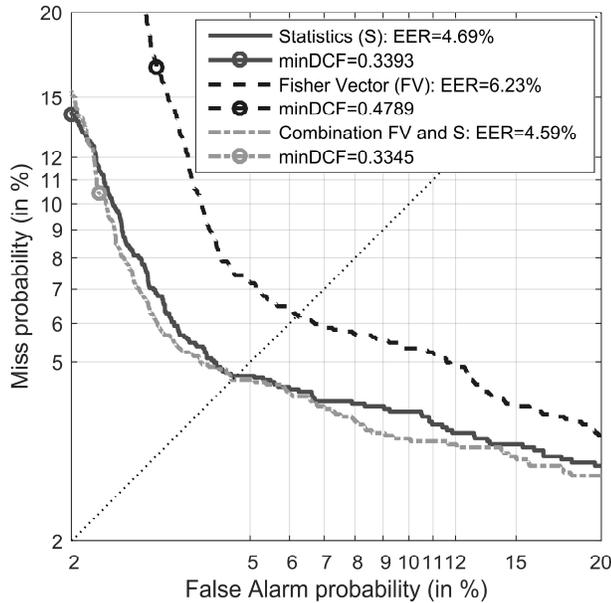


Fig. 2. DET curves of comparison systems: the system using Suprvector of Statistics (S) or FV and a combination of these two systems. Circles denote points where minDCF occurred. Dotted line indicates EER.

the PLDA verification process as the Suprvector of Statistics when sufficient data are available.

V. CONCLUSIONS

In this work, we compared two approaches to preserve the information about a speaker - as Suprvector of Statistics and as Fisher Vector - both containing the same amount of information (zeroth and first moments). These representations were used in a state-of-the-art verification system based on PLDA. The experimental results of these two approaches show only a small difference in the EER of these systems. Moreover, the fusion of these two systems indicates that the Fisher Vector brings almost the same information to the PLDA verification process as the Suprvector of Statistics. Although, when observing the shapes of the DET curves we can conclude that the fused system is a bit more robust on a larger scale of verification thresholds. Since our results are different from the ones reported recently, we assume that the difference is mainly in the use of PLDA and the amount of the information used for Fisher Vectors.

ACKNOWLEDGMENT

The work was supported by the Ministry of Education, Youth and Sports of the Czech Republic project No. LO1506. Access to computing and storage facilities owned by parties and projects contributing to the National Grid Infrastructure MetaCentrum, provided under the programme "Projects of Large Research, Development, and Innovations Infrastructures" (CESNET LM2015042), is greatly appreciated.

- [1] S. J. Prince and J. H. Elder, *Probabilistic Linear Discriminant Analysis for Inferences About Identity*. IEEE 11th International Conference on Computer Vision, 1–8, Rio de Janeiro, 2007.
- [2] M. Turk and A. Pentland, *Face Recognition Using Eigenfaces*. Journal of Cognitive Neuroscience 3(1), 72 – 86, 1991.
- [3] P. Kenny, *Joint factor analysis of speaker and session variability: Theory and algorithms*. Tech. rep., 2006.
- [4] K. Simonyan and O. Parkhi and A. Vedaldi and A. Zisserman, *Fisher Vector Faces in the Wild*. Proceedings of the British Machine Vision Conference 2013, 8.1–8.12. Bristol, 2013.
- [5] W. M. Campbell and D. E. Sturim and D. A. Reynolds, *Support vector machines using GMM supervectors for speaker verification*. IEEE Signal Processing Letters 13(5), 308–311, 2006.
- [6] J. Vaněk and J. Trmal and J. V. Psutka, *Optimization of the Gaussian Mixture Model Evaluation on GPU*. 12th Annual Conference of the International Speech Communication Association (INTERSPEECH), 1748–1751, Firenze, 2011.
- [7] Z. Zajíc and L. Machlica and L. Müller, *Robust Statistic Estimates for Adaptation in the Task of Speech Recognition*. Lecture Notes in Computer Science, vol. 6231, 464–471, 2010.
- [8] W. Rao and M. w. Mak and K. a. Lee, *Normalization of Total Variability Matrix for i-Vector/PLDA Speaker Verification*. Acoustics, Speech and Signal Processing (ICASSP), 4180–4184, 2015.
- [9] L. Machlica and Z. Zajíc and L. Müller, *On Complementarity of State-of-the-art Speaker Recognition Systems*. IEEE International Symposium on Signal Processing and Information Technology, 164–169, Ho Chi Minh City, 2012.
- [10] D. A. Reynolds and T. F. Quatieri and R. B. Dunn, *Speaker Verification Using Adapted Gaussian Mixture Models*. Digital Signal Processing 10(1-3), 19–41, 2000.
- [11] P. Kenny and P. Dumouchel, *Experiments in Speaker Verification using Factor Analysis Likelihood Ratios*. Odyssey - Speaker and Language Recognition Workshop, 219–226, Toledo, 2004.
- [12] D. Garcia-Romero and C. Y. Espy-Wilson, *Analysis of i-vector Length Normalization in Speaker Recognition Systems*. Interspeech 2011, 249–252, Florence, 2011.
- [13] L. Machlica and Z. Zajíc, *Factor Analysis and Nuisance Attribute Projection Revisited*. Interspeech 2012, 1570–1573, Portland, 2012.
- [14] P. Kenny and P. Ouellet and N. Dehak and V. Gupta and P. Dumouchel, *A Study of Interspeaker Variability in Speaker Verification*. IEEE Transactions on Audio, Speech, and Language Processing, vol. 16(5), 980–988, 2008.
- [15] L. Machlica and Z. Zajíc, *Analysis of the Influence of Speech Corpora in the PLDA Verification in the Task of Speaker Recognition*. Lecture Notes in Computer Science, vol. 7499, 464–471, 2012.
- [16] T. S. Jaakkola and D. Haussler, *Exploiting generative models in discriminative classifiers*. Proceedings of the 1998 conference on Advances in neural information processing systems II. vol. 1, 487–493, 1999.
- [17] F. Perronnin and J. Sánchez and T. Mensink, *Improving the Fisher kernel for large-scale image classification*. Lecture Notes in Computer Science, vol. 6314, 143–156, 2010.
- [18] N. Brummer *FoCal: Tools for fusion and calibration of automatic speaker detection systems*, 2006. <https://sites.google.com/site/nikobrunner/focal>
- [19] N. Brummer and E. de Villiers, *The BOSARIS Toolkit User Guide: Theory, Algorithms and Code for Binary Classifier Score Processing*. Tech. rep., 2011.
- [20] A. Martin and G. Doddington and T. Kamm and M. Ordowski and M. Przybocki, *The DET Curve in Assessment of Detection Task Performance*. Eurospeech. vol. 4, 1899–1903. Rhodes, 1997.
- [21] Y. Tian, and L. He and Z. Y. Li and W. L. Wu and W. Q. Zhang and J. Liu, *Speaker verification using Fisher vector*. In: Proceedings of the 9th International Symposium on Chinese Spoken Language Processing (ISCSLP), 419–422, Singapore, 2014.