

TEXT-INDEPENDENT SPEAKER VERIFICATION: A COMPARATIVE ANALYSIS STUDY

Mohamed SOLTANE, Nouredine DOGHMANE and Nouredine GUERSI

Electronics Department, Faculty of engineering science

Badji Mokhtar University of Annaba, 23000 Annaba - Algeria

soltane@lri-annaba.net & xor99@hotmail.com, ndoghmane@univ-annaba.org, guersi54@yahoo.fr

Abstract: Gaussian mixture models (GMMs) remain the state of the art technique for modeling spectral envelope features for speech recognition systems. This paper presents a comparative analysis of the performance of three estimation algorithms Expectation Maximization (EM), Greedy EM Algorithm (GEM) and Figueiredo-Jain Algorithm (FJ) based Gaussian mixture models (GMMs) for text-independent speech biometrics verification. The simulation results are showed significant performance achievements. The test performance of, EER=0.26 % for "EM", EER=0.21 % for "GEM" and EER=0.16 % for "FJ", show that the behavioral information scheme of speech biometrics is more robust and have a discriminating power, which can be explored for identity authentication.

Key words: Biometric authentication, behavioral, speech, soft decision and Gaussian Mixture Modal, EM, GEM and FJ.

I. INTRODUCTION

BIOMETRIC is a Greek composite word stemming from the synthesis of bio and metric, meaning life measurement. In this context, the science of biometrics is concerned with the accurate measurement of unique biological characteristics of an individual in order to securely identify them to a computer or other electronic system. Biological characteristics measured usually include fingerprints, voice patterns, retinal and iris scans, face patterns, and even the chemical composition of an individual's DNA [1]. Biometrics authentication (BA) (*Am I whom I claim I am?*) involves confirming or denying a person's *claimed identity* based on his/her physiological or behavioral characteristics [2]. BA is becoming an important alternative to traditional authentication methods such as keys ("something one has", i.e., by possession) or PIN numbers ("something one knows", i.e., by knowledge) because it is essentially "who one is", i.e., by biometric information. Therefore, it is not susceptible to misplacement or forgetfulness [3]. These biometric systems for personal authentication and identification are based upon physiological or behavioral features which are typically distinctive, although time varying, such as fingerprints, hand geometry, face, voice, lip movement, gait, and iris patterns. An identity verification system has to deal with two kinds of events: either the person claiming a given identity is the one who he claims to be (in which case, he is called a *client*), or he is not (in which case, he is called an *impostor*). Moreover, the system may generally take two decisions: either *accept* the *client* or *reject* him and decide he is an *impostor*.

Some works based on biometric speech identity verification systems has been reported in literature. **Fortuna J. et al.** [16] presents a comparative analysis of the performance of decoupled and adapted Gaussian mixture models (GMMs) for open-set, text-independent speaker identification (OSTISI) and concluded that the speaker identification performance is noticeably better

with adapted-GMMs than with decoupled-GMMs. Their analysis is based on a set of experiments using an appropriate subset of the NIST-SRE 2003 database and various score normalization methods. They included a detailed description of the experiments and discuss how the OSTI-SI performance is influenced by the characteristics of each of the two modeling techniques and the normalization approaches adopted. **Eduardo Sanchew-Soto et al.** [17] present a new adaptation technique for speaker verification of models built using Bayesian Networks tested using the NIST 2002 data base and showed improvement in the verification performances. The adaptation problem of parameters of the conditional probability tables (CPTs) is treated in a specific manner. The model adaptation involves estimating the new vectors with a transformation that includes vectors in the world model and the speaker model and the combination of both models is based on a value computed using a measure of distance between vectors of both CPTs. **Arnon Cohen et al.** [18] describes an HMM based speaker verification system evaluated on a text-dependent database, which verifies speakers in their own specific feature space. The user feature space is determined by a Dynamic Programming (DP) feature selection algorithm, in which a suitable criterion, correlated with Equal Error Rate (EER) was developed and is used for this feature selection algorithm. A significant improvement in verification results was demonstrated with the DP selected individual feature space. An EER of 4.8% was achieved when the feature set was the "almost standard" Mel Frequency Cepstrum Coefficients (MFCC) space (12 MFCC + 12 Δ MFCC). Under the same conditions, a system based on the selected feature space yielded an EER of only 2.7%. **Mijail Arcienega et al.** [19] present a Bayesian network approach for modeling the pitch and spectral envelope and showed an increase in the performance of the speaker recognition system. In which the conditional statistical distributions (represented by GMMs) of the features are

simultaneously exploited for increasing the recognition score within the approach, and in particularly in noisy conditions. **Driss Matrouf et al. [20]** investigates the effect of voice transformation on automatic speaker recognition systems performance and showed an increase of about 2.7 time of the likelihood ratio, without a degradation of the natural aspect of the voice. It focuses on increasing the impostor acceptance rate, by modifying the voice of an impostor in order to target a specific speaker. Their work is inspired from the idea that in several forensic situations, it is reasonable to think that some organizations have a knowledge on the speaker recognition method used by the police department and could impersonate a given, well known speaker.

II. BIOMETRIC SPEECH VERIFICATION

1. Speech Analysis and Feature Extraction

Gaussian Mixture Models (GMMs), is the main tool used in text-independent speaker verification, in which can be trained using the Expectation Maximization (EM) algorithm [4]. In this work the speech modality, is authenticated with a multi-lingual text-independent speaker verification system. The speech trait is comprised of two main components as shown in figure 1: speech feature extraction and a Gaussian Mixture Model (GMM) classifier. The speech signal is analyzed on a frame by frame basis, with a typical frame length of 20 ms and a frame advance of 10 ms [5]. For each frame, a dimensional feature vector is extracted, the discrete Fourier spectrum is obtained via a fast Fourier transform from which magnitude squared spectrum is computed and put it through a bank of filters. The critical band warping is done following an approximation to the Mel-frequency scale which is linear up to 1000 Hz and logarithmic above 1000 Hz. The Mel-scale cepstral coefficients are computed from the outputs of the filter bank [6]. The state of the art speech feature extraction schemes (Mel frequency cepstral coefficients (MFCC) is based on auditory processing on the spectrum of speech signal and cepstral representation of the resulting features [7]. One of the powerful properties of cepstrum is the fact that any periodicities, or repeated patterns, in a spectrum will be mapped to one or two specific components in the cepstrum. If a spectrum contains several harmonic series, they will be separated in a way similar to the way the spectrum separates repetitive time patterns in the waveform. The description of the different steps to exhibit features characteristics of an audio sample with MFCC is showed in figure 2.

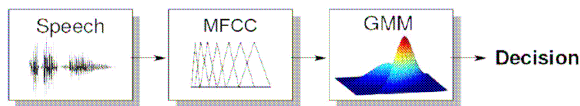


Figure 1. Acoustic Speech Analysis.

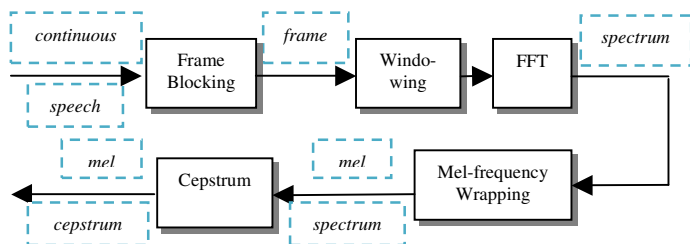


Figure 2. MFCC calculation Block diagram [6].

The distribution of feature vectors for each person is modeled by a GMM. The parameters of the Gaussian mixture probability density function are estimated using three different estimation algorithms. The Expectation Maximization (EM) algorithm [8], Greedy algorithm (GEM) [8] and Figueiredo-Jain (FJ) algorithm [8].

Given a claim for person C 's identity and a set of feature vectors $X = \{\vec{x}_i\}_{i=1}^{N_C}$ supporting the claim, the average log likelihood of the claimant being the true claimant is calculated using:

$$\mathcal{L}(X|\lambda_C) = \frac{1}{N_C} \sum_{i=1}^{N_C} \log p(\vec{x}_i|\lambda_C) \quad (1)$$

where $p(\vec{x}|\lambda) = \sum_{j=1}^{N_M} m_j \mathcal{N}(\vec{x}; \vec{\mu}_j, \Sigma_j)$ (2)

and $\lambda = \{m_j, \vec{\mu}_j, \Sigma_j\}_{j=1}^{N_M}$ (3)

Here λ_C is the model for person C . N_M is the number of mixtures, m_j is the weight for mixture j (with constraint $\sum_{j=1}^{N_M} m_j = 1$), and $\mathcal{N}(\vec{x}; \vec{\mu}, \Sigma)$ is a multi-variate Gaussian function with mean $\vec{\mu}$ and diagonal covariance matrix Σ . Given a set $\{\lambda_b\}_{b=1}^B$ of B background person models for person C , the average log likelihood of the claimant being an impostor is found using:

$$\mathcal{L}(X|\lambda_C) = \log \left[\frac{1}{B} \sum_{b=1}^B \exp \mathcal{L}(X|\lambda_b) \right] \quad (4)$$

The set of background person models is found using the method described in [15]. An opinion on the claim is found using:

$$o = \mathcal{L}(X|\lambda_C) - \mathcal{L}(X|\lambda_{\bar{C}}) \quad (5)$$

The opinion reflects the likelihood that a given claimant is the true claimant (i.e., a low opinion suggests that the claimant is an impostor, while a high opinion suggests that the claimant is the true claimant).

2. Maximum Likelihood Parameter Estimation

Given a set of observation data in a matrix X and a set of observation parameters θ the ML parameter estimation aims at maximizing the likelihood $L(\theta)$ or log likelihood of the observation data $X = \{X_1, \dots, X_n\}$

$$\hat{\theta} = \arg \max_{\theta} L(\theta). \quad (6)$$

Assuming that it has independent, identically distributed data, it can write the above equations as:

$$L(\theta) = p(X|\theta) = p(X_1, \dots, X_n|\theta) = \prod_{i=1}^n p(X_i|\theta). \quad (7)$$

The maximum for this function can be find by taking the derivative and set it equal to zero, assuming an analytical function.

$$\frac{\partial}{\partial \theta} L(\theta) = 0. \quad (8)$$

The incomplete-data log-likelihood of the data for the mixture model is given by:

$$\mathcal{L}(X|\theta) = \log p(X|\theta) = \sum_{i=1}^N \log p(x_i|\theta) \quad (9)$$

which is difficult to optimize because it contains the log of the sum. If it considers X as incomplete, however, and posits the existence of unobserved data items $Y = \{y_i\}_{i=1}^N$ whose values inform us which component density generated each data item, the likelihood expression is significantly simplified. That is, it assume that $y_i \in \{1 \dots K\}$ for each i , and $y_i = k$ if the i -th sample was generated by the k -th mixture component. If it knows the values of Y , it obtains the complete-data log-likelihood, given by:

$$L(\theta, Y) = \log p(X, Y | \theta) \tag{10}$$

$$= \sum_{i=1}^N \log p(x_i, y_i | \theta) \tag{11}$$

$$= \sum_{i=1}^N \log(p(y_i | \theta) p(x_i | y_i, \theta)) \tag{12}$$

$$= \sum_{i=1}^N (\log p_{y_i} + \log g(x_i | \mu_{y_i}, \Sigma_{y_i})) \tag{13}$$

Which, given a particular form of the component densities, can be optimized using a variety of techniques [9].

2.1 EM algorithm:

The expectation-maximization (EM) algorithm [8,10, 11, 12] is a procedure for maximum-likelihood (ML) estimation in the cases where a closed form expression for the optimal parameters is hard to obtain. This iterative algorithm guarantees the monotonic increase in the likelihood L when the algorithm is run on the same training database.

The probability density of the Gaussian mixture of k components in \mathbb{R}^d can be described as follows:

$$\phi(x) = \sum_{i=1}^k \pi_i \mathcal{O}(x | \theta_i) \quad \forall x \in \mathbb{R}^d, \tag{14}$$

where $\mathcal{O}(x | \theta_i)$ is a Gaussian probability density with the parameters $\theta_i = (m_i, \Sigma_i)$, m_i is the mean vector and Σ_i is the covariance matrix which is assumed positive definite given by:

$$\mathcal{O}(x | \theta_i) = \mathcal{O}(x | m_i, \Sigma_i) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_i|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-m_i)^T \Sigma_i^{-1} (x-m_i)}, \tag{15}$$

and $\pi_i \in [0, 1]$ ($i = 1, 2, \dots, k$) are the mixing proportions under the constraint $\sum_{i=1}^k \pi_i = 1$. If it encapsulate all the parameters into one vector: $\theta_k = (\pi_1, \pi_2, \dots, \pi_k, \theta_1, \theta_2, \dots, \theta_k)$, then, according to Eq. (23), the density of Gaussian mixture can be rewritten as:

$$\phi(x | \theta_k) = \sum_{i=1}^k \pi_i \mathcal{O}(x | \theta_i) = \sum_{i=1}^k \pi_i \mathcal{O}(x | m_i, \Sigma_i). \tag{16}$$

For the Gaussian mixture modeling, there are many learning algorithms. But the EM algorithm may be the most well-known one. By alternatively implementing the E-step to estimate the probability distribution of the unobservable random variable and the M-step to increase the log-likelihood function, the EM algorithm can finally lead to a local maximum of the log-likelihood function of the model. For the Gaussian mixture model, given a sample data set $S = \{x_1, x_2, \dots, x_N\}$ as a special

incomplete data set, the log-likelihood function can be expressed as follows:

$$\log p(S | \theta_k) = \log \prod_{i=1}^N \mathcal{O}(x_i | \theta_k) = \sum_{i=1}^N \log \sum_{j=1}^k \pi_j \mathcal{O}(x_i | \theta_j). \tag{17}$$

Which can be optimized iteratively via the EM algorithm as follows:

$$P(j | x_i) = \frac{\pi_j \mathcal{O}(x_i | \theta_j)}{\sum_{l=1}^k \pi_l \mathcal{O}(x_i | \theta_l)} \tag{18}$$

$$\pi_j^* = \frac{1}{N} \sum_{i=1}^N P(j | x_i), \tag{19}$$

$$\mu_j^* = \frac{1}{\sum_{i=1}^N P(j | x_i)} \sum_{i=1}^N P(j | x_i) x_i, \tag{20}$$

$$\Sigma_j^* = \frac{1}{\sum_{i=1}^N P(j | x_i)} \sum_{i=1}^N P(j | x_i) (x_i - \mu_j^*)(x_i - \mu_j^*)^T. \tag{21}$$

Although the EM algorithm can have some good convergence properties in certain situations, it certainly has no ability to determine the proper number of the components for a sample data set because it is based on the maximization of the likelihood.

2.2 Greedy EM Algorithm:

The greedy algorithm (GEM) [8, 10, 12, 13] starts with a single component and then adds components into the mixture one by one. The optimal starting component for a Gaussian mixture is trivially computed, optimal meaning the highest training data likelihood. The algorithm repeats two steps: insert a component into the mixture, and run EM until convergence. Inserting a component that increases the likelihood the most is thought to be an easier problem than initializing a whole near-optimal distribution. Component insertion involves searching for the parameters for only one component at a time. Recall that EM finds a local optimum for the distribution parameters, not necessarily the global optimum which makes it initialization dependent method.

Given p_C a C -component Gaussian mixture with parameters θ_C , the general greedy algorithm for Gaussian mixture is as follows:

1. Compute the optimal (in the ML sense) one-component mixture p_1 and set $C \leftarrow 1$.
2. Find a new component $\mathcal{N}(x; \mu', \Sigma')$ and corresponding mixing weight α' that increase the likelihood the most:

$$(\mu', \Sigma', \alpha') = \arg \max_{\mu, \Sigma, \alpha} \sum_{i=1}^N \ln[(1-\alpha)p_C(x_i) + \alpha \mathcal{N}(x_i; \mu, \Sigma)] \tag{22}$$

while keeping p_C fixed.

3. Set $p_{C+1}(x) \leftarrow (1-\alpha')p_C(x) + \alpha' \mathcal{N}(x; \mu', \Sigma')$ and then $C \leftarrow C + 1$.
4. Update p_C using EM (or more other method) until convergence.
5. Evaluate some stopping criterion; go to step 2 or quit.

The stopping criterion in Step 5 can be for example any kind of model selection criterion, wanted number of components, or the minimum message length criterion.

The crucial point is of course Step 2. Finding the optimal new component requires a global search, which is performed by creating CN_{cand} candidate components. The number of candidates will increase linearly with the number of components C , having N_{cand} candidates per each existing component. The candidate resulting in the highest likelihood when inserted into the (previous) mixture is selected. The parameters and weight of the best candidate are then used in Step 3 instead of the truly optimal values.

The candidates for executing Step 2 are initialized as follows: the training data set X is partitioned into C disjoint data sets $\{A_c\}, c = 1 \dots C$, according to the posterior probabilities of individual components; the data set is Bayesian classified by the mixture components. From each A_c number of N_{cand} candidates are initialized by picking uniformly randomly two data points x_l and x_r in A_c . The set A_c is then partitioned into two using the smallest distance selection with respect to x_l and x_r . The mean and covariance of these two new subsets are the parameters for two new candidates. The candidate weights are set to half of the weight of the component that produced the set A_c . Then new x_l and x_r are drawn until N_{cand} candidates are initialized with A_c . The partial EM algorithm is then used on each of the candidates. The partial EM differs from the EM and CEM algorithms by optimizing (updating) only one component of a mixture; it does not change any other components. In order to reduce the time complexity of the algorithm a lower bound on the log-likelihood is used instead of the true log-likelihood. The lower-bound log-likelihood is calculated with only the points in the respective set A_c . The partial EM update equations are as follows:

$$w_{l,c+1} = \frac{\alpha N(x_l, \mu, \Sigma)}{(1-\alpha) p_c(x) + \alpha N(x_l, \mu, \Sigma)} \quad (23)$$

$$\alpha = \frac{1}{N(A_c)} \sum_{l \in A_c} w_{l,c+1} \quad (24)$$

$$\mu = \frac{\sum_{l \in A_c} w_{l,c+1} x_l}{\sum_{l \in A_c} w_{l,c+1}} \quad (25)$$

$$\Sigma = \frac{\sum_{l \in A_c} w_{l,c+1} (x_l - \mu)(x_l - \mu)^T}{\sum_{l \in A_c} w_{l,c+1}} \quad (26)$$

where $N(A_c)$ is the number of training samples in the set A_c . These equations are much like the basic EM update equations in Eqs. (19) - (21). The partial EM iterations are stopped when the relative change in log-likelihood of the resulting $C + 1$ -component mixture drops below threshold or maximum number of iterations is reached. When the partial EM has converged the candidate is ready to be evaluated.

2.3 Figueiredo-Jain Algorithm:

The Figueiredo-Jain (FJ) [8,10,12,13] algorithm tries to overcome three major weaknesses of the basic EM algorithm. The EM algorithm presented previous section requires the user to set the number of components and the number will be fixed during the estimation process. The FJ algorithm adjusts the number of components during estimation by annihilating components that are not supported by the data. This leads to the other EM failure point, the boundary of the parameter space. FJ avoids the

boundary when it annihilates components that are becoming singular. FJ also allows starting with an arbitrarily large number of components, which tackles the initialization issue with the EM algorithm. The initial guesses for component means can be distributed into the whole space occupied by training samples, even setting one component for every single training sample.

The classical way to select the number of mixture components is to adopt the "model-class/model" hierarchy, where some candidate models (mixture pdf's) are computed for each model-class (number of components), and then select the "best" model. The idea behind the FJ algorithm is to abandon such hierarchy and to find the "best" overall model directly. Using the minimum message length criterion and applying it to mixture models leads to the objective function:

$$A(\theta, X) = \frac{V}{2} \sum_{c: \alpha_c > 0} \ln \left(\frac{N \alpha_c}{2} \right) + \frac{c_{nz}}{2} \ln \frac{N}{2} + \frac{c_{nz}(V+1)}{2} - \ln \mathcal{L}(X, \theta) \quad (27)$$

Where N is the number of training points, V is the number of free parameters specifying a component, and c_{nz} is the number of components with nonzero weight in the mixture ($\alpha_c > 0$). θ in the case of Gaussian mixture is the same as in (Eq. 3) the last term $\ln \mathcal{L}(X, \theta)$ is the log-likelihood of the training data given the distribution parameters (Eq. 13).

The EM algorithm can be used to minimize Eq. 27 with a fixed c_{nz} . It leads to the M-step with component weight updating formula:

$$\alpha_c^{i+1} = \frac{\max \left\{ 0, \left(\sum_{n=1}^N w_{n,c} - \frac{V}{2} \right) \right\}}{\sum_{j=1}^C \max \left\{ 0, \left(\sum_{n=1}^N w_{n,j} - \frac{V}{2} \right) \right\}} \quad (28)$$

This formula contains an explicit rule of annihilating components by setting their weights to zero.

The above M-steps are not suitable for the basic EM algorithm though. When initial C is high, it can happen that all weights become zero because none of the components have enough support from the data. Therefore a component-wise EM algorithm (CEM) is adopted. CEM updates the components one by one, computing the E-step (updating W) after each component update, where the basic EM updates all components "simultaneously". When a component is annihilated its probability mass is immediately redistributed strengthening the remaining components.

When CEM converges, it is not guaranteed that the minimum of $A(\theta, X)$ is found, because the annihilation rule (Eq. 28) does not take into account the decrease caused by decreasing c_{nz} . After convergence the component with the smallest weight is removed and the CEM is run again, repeating until $c_{nz} = 1$. Then the estimate with the smallest $A(\theta, X)$ is chosen. The implementation of the FJ algorithm uses a modified cost function instead of $A(\theta, X)$.

$$A'(\theta, X) = \frac{V}{2} \sum_{c: \alpha_c > 0} \ln \alpha_c + \frac{c_{nz}(V+1)}{2} \ln N - \ln \mathcal{L}(X, \theta). \quad (29)$$

III. EXPERIMENTS AND RESULTS

The experiments were performed using audio database extracted from video, which is encoded in raw UYVY. AVI 640 x 480, 15.00 fps with uncompressed 16bit PCM audio; mono, 32000 Hz little endian. The capturing devices for recording the video and audio data were: Allied Vision Technologies AVT marlin MF-046C 10 bit ADC, 1/2" (8mm) Progressive scan SONY IT CCD; and Shure SM58 microphone. Frequency response 50 Hz to 15000 Hz. Unidirectional (Cardiod) dynamic vocal microphones. The extracted 16 bit PCM audio files (with wav header), were sampled at 16000 Hz, mono little endian. Thirty subjects were used for the experiments in which twenty-six are males and four are females. For each subject, six multi-lingual (.wav files) of one minute each recording were used for each subject. The database obtained from eNTERFACE 2005 [14]. For the experts, four speech recording samples of one minute each one were used for the modeling (training); two samples were used for the subsequent validation and testing. Three sessions of the speech database were used separately. Session one was used for training the speech experts. Each expert used ten mixture client models. To find the performance, Sessions two and three were used for obtaining expert opinions of known impostor and true claims.

Performance Criteria:

The basic error measure of a verification system is false rejection rate (*FRR*) and false acceptance rate (*FAR*) as defined in the following equations:

False Rejection Rate (*FRR_i*): is an average of number of falsely rejected transactions. If *n* is a transaction and *x(n)* is the verification result where 1 is falsely rejected and 0 is accepted and *N* is the total number of transactions then the personal False Rejection Rate for user *i* is

$$FRR_i = \frac{1}{N} \sum_{n=1}^N x(n) \tag{30}$$

False Acceptance rate (*FAR_i*): is an average of number of falsely accepted transactions. If *n* is a transaction and *x(n)* is the verification result where 1 is a falsely accepted transaction and 0 is genuinely accepted transaction and *N* is the total number of transactions then the personal False Acceptance Rate for user *i* is

$$FAR_i = \frac{1}{N} \sum_{n=1}^N x(n) \tag{31}$$

Both *FRR_i* and *FAR_i* are usually calculated as averages over an entire population in a test. If *P* is the size of populations then these averages are

$$FRR = \frac{1}{P} \sum_i FRR_i \tag{32}$$

$$FAR = \frac{1}{P} \sum_i FAR_i \tag{33}$$

Equal Error Rate (*EER*), is an intersection where *FAR* and *FRR* are equal at an optimal threshold value. This threshold value shows where the system performs at its best (see Figure 3).

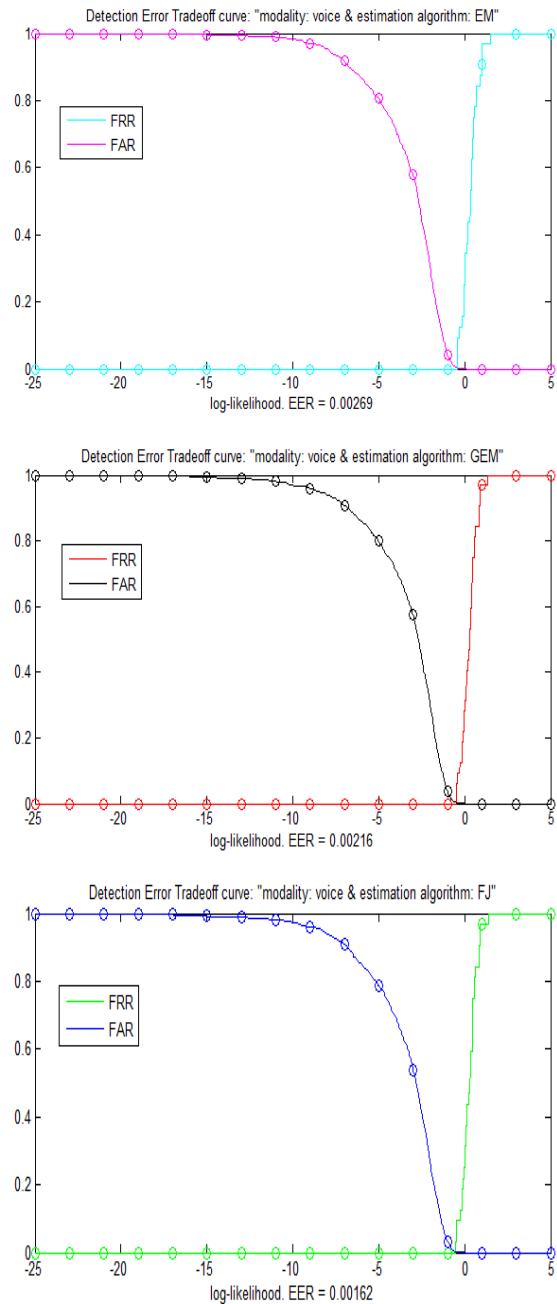


Figure 3. Detection error tradeoff curves

As a common starting point, classifier parameters were selected to obtain performance as close as possible to *EER* on clean test data (following the standard practice in the biometric verification area of using *EER* as a measure of expected performance). A good decision is to choose the decision threshold such as the false accept equal to the false reject rate. In this paper it uses the Detection Error Tradeoff (DET) curve to visualize and compare the performance of the system.

IV. CONCLUSIONS

The paper has presented a human authentication method of behavioural biometrics speech information. Simulation results show that state-of-the art finite mixture modal

(GMM) is quite effective in modelling the genuine and impostor score densities. The (EM), (GEM) and (FJ) estimation algorithms achieve a significant performance rates, EER=0.26 % for "EM", EER=0.21 % for "GEM" and EER=0.16 % for "FJ". Hence, the behavioral information scheme based speech biometrics is robust and have a discriminating power, which can be explored for identity authentication.

REFERENCES

- [1] Sofia Gleni & Panagiotis Petratos, "DNA Smart Card for Financial Transactions" The ACM Student Magazine 2004, <http://www.acm.org>
- [2] Girija Chetty and Michael Wagner, "Audio-Visual Multimodal Fusion for Biometric Person Authentication and Liveness Verification", Copyright © 2006, Australian Computer Society, Inc. This paper appeared at the *NICTA-HCSNet Multimodal UserInteraction Workshop (MMUI2005)*, Sydney, Australia.
- [3] Norman Poh and Samy Bengio, "Database, Protocol and Tools for Evaluating Score-Level Fusion Algorithms in Biometric Authentication", IDIAP RR 04-44, August 2004, a IDIAP, CP 592, 1920 Martigny, Switzerland.
- [4] Conrad Sanderson, Samy Bengio, Herve Boulard, Johnny Mariéthoz, Ronan Collobert, Mohamed F. BenZeghiba, Fabien Cardinaux, and Sébastien Marcel, "SPEECH & FACE BASED BIOMETRIC AUTHENTICATION AT IDIAP", Dalle Molle Institute for Perceptual Artificial Intelligence (IDIAP). Rue du Simplon 4, CH-1920 Martigny, Switzerland.
- [5] Claus Vielhauer^a, Sascha Schimke^a, Valsamakis Thanassis^b, Yannis Stylianou^b,^a Otto-von-Guericke University Magdeburg, Universitaetsplatz 2, D-39106, Magdeburg, Germany,^b University of Crete, Department of Computer Science, Heraklion, Crete, Greece, "Fusion Strategies for Speech and Handwriting Modalities in HCP", Multimedia on Mobile Devices, edited by Reiner Creutzburg, Jarmo H. Takala, Proc. of SPIE-IS&T Electronic Imaging, Vol. 5684 © 2005
- [6] Lasse L. Mølgaard and Kasper W. Jørgensen, "Speaker Recognition: Special Course", IMM_DTU December 14, 2005
- [7] S. Davis and P. Mermelstein. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE Transactions on Acoustics, Speech and Signal Processing*, (4):357–366, 1980.
- [8] Pekka Paalanen, "Bayesian classification using gaussian mixture model and EM estimation: implementation and comparisons", Information Technology Project, 2004, Lappeenranta, June 23, 2004, <http://www.it.lut.fi/project/gmmbayes/>
- [9] Pham Minh Tri, "On estimating the parameters of Gaussian mixtures using EM", School of Computer Engineering, Nanyang Technological University
- [10] P. Paalanen, J.-K. Kamarainen, J. Ilonen, H. Kälviäinen, "Feature Representation and Discrimination Based on Gaussian Mixture Model Probability Densities: Practices and Algorithms", Department of Information Technology, Lappeenranta University of Technology, P.O.Box 20, FI-53851 Lappeenranta, Finland 2005
- [11] Lei Li and Jinwen Ma*, "A BYY Split-and-Merge EM Algorithm for Gaussian Mixture Learning", F. Sun et al. (Eds.): ISNN 2008, Part I, LNCS 5263, pp. 600–609, 2008. ©Springer-Verlag Berlin Heidelberg 2008. Department of Information Science, School of Mathematical, Sciences and LAMA, Peking University, Beijing, 100871, China
- [12] J.J. Verbeek N. Vlassis B. Krose, "Efficient Greedy Learning of Gaussian Mixture Models", Informatics Institute, University of Amsterdam- Kruislaan 403, 1098 SJ Amsterdam, The Netherlands. Published in *Neural Computation* 15(2), pages 469-485, 2003.
- [13] Jan R.J. Nunnink Jakob J. Verbeek Nikos Vlassis, "Accelerated Greedy Mixture Learning", Informatics Institute, Faculty of Science, University of Amsterdam Kruislaan 403, 1098 SJ Amsterdam, The Netherlands
- [14] Yannis Stylianou, Yannis Pantazis, Felipe Calderero, Pedro Larroy, Francois Severin, Sascha Schimke, Rolando Bonal, Federico Matta, and Athanasios Valsamakis, "GMM-Based Multimodal Biometric Verification", eINTERFACE 2005 The summer Workshop on Multimodal Interfaces July 18th – August 12th, Faculté Polytechnique de Mons, Belgium.
- [15] D.A. Reynolds, "Experimental Evaluation of Features for Robust Speaker Identification", *IEEE Trans. Speech and Audio Processing* 2 (4), 1994, 639-643.
- [16] J. Fortuna, A. Malesonkar, A. Ariyaceenia and P. Sivasubramanian*, "ON THE USE OF DECOUPLED AND ADAPTED GAUSSIAN MIXTURE MODELS FOR OPEN-SET SPEAKER IDENTIFICATION", University of Hertfordshire, Hatfield, UK, *Canon Research Centre Europe Ltd., Bracknell, UK. THE THIRD COST 275 WORKSHOP, Biometrics on the Internet. University of Hertfordshire Hatfield, UK 27-28 October 2005
- [17] Eduardo Sánchez-Soto, Raphaël Blouet, marc Sigelle and Gérard Chollet, "Model Adaptation for Speaker Verification using Conditional probability tables in Bayesian Networks", Ecole Nationale Supérieure des Télécommunications – Département de Traitement de Signal et des Images. LTCI/CNRS URA 820, 46 rue Barrault 75634 Paris Cedex 13 France. 2nd COST 275 Workshop – Biometrics on the Internet Vigo, 25-26 March 2004.
- [18] Arnon Cohen and Yaniv Zigel, "ON FEATURE SELECTION FOR SPEAKER VERIFICATION", Electrical and Computer Engineering Department, Ben-Gurion University, Beer-Sheva, Israel. COST 275 Workshop – The Advent Biometrics on the Internet Rome, Italy November 7-8, 2002.
- [19] Mijail Arcienega and Andrzej Drygajlo, "SPECTRAL ENVELOPE FEATURES FOR SPEAKER VERIFICATION", Signal Processing Institute Swiss Federal Institute of Technology, Lausanne. COST 275 Workshop – The Advent Biometrics on the Internet Rome, Italy November 7-8, 2002.
- [20] Driss Matrouf, Jean-François Bonastre, Jean-Pierre Costa, "EFFECT OF IMPOSTOR SPEECH TRANSFORMATION ON AUTOMATIC SPEAKER RECOGNITION", LIA, Université d'Avignon Agroparc, BP 1228- 84911 Avignon CEDEX 9, France. THE THIRD COST 275 WORKSHOP, Biometrics on the Internet. University of Hertfordshire Hatfield, UK 27-28 October 2005.