

On the importance of pre- and post-conditioning procedures for speaker recognition systems based on total variability factors

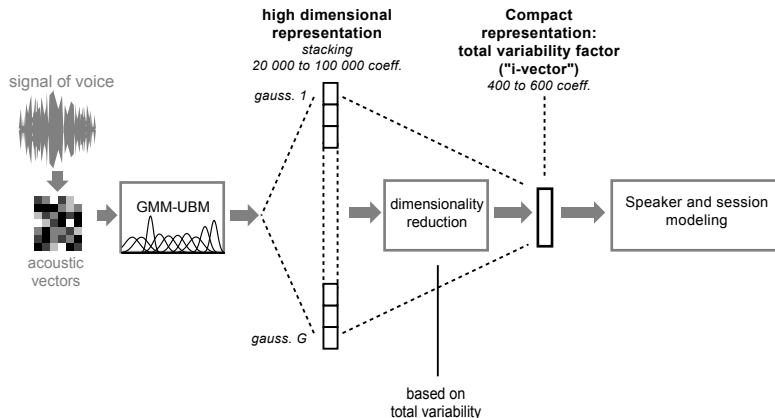
P.M. Bousquet, J.F. Bonastre



LIA
University of Avignon

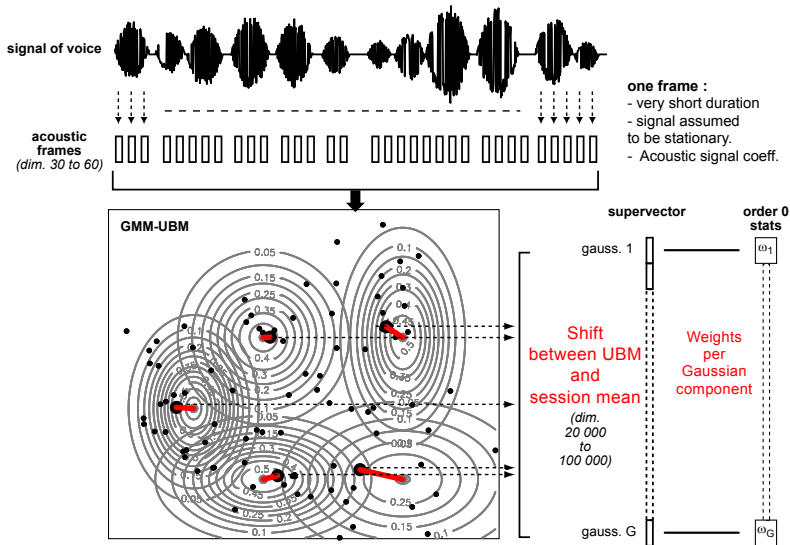


The different steps of the i-vector solution



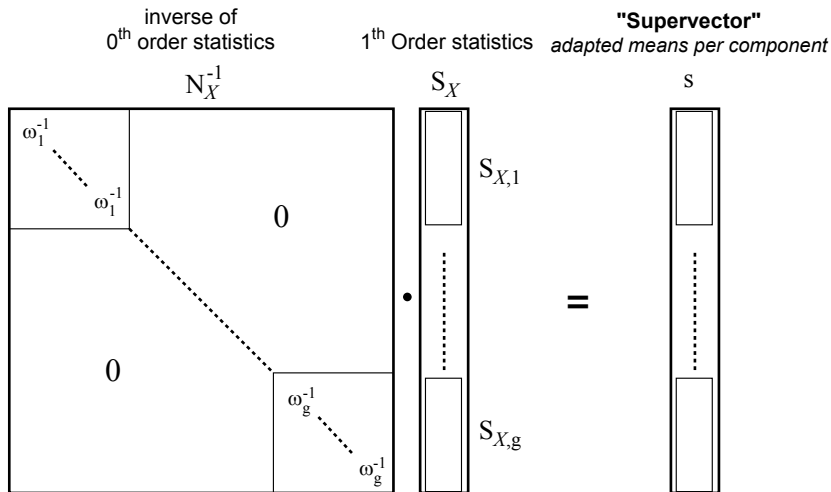
GMM-UBM high-dim. representation

i. With 0^{th} and 1^{th} order statistics



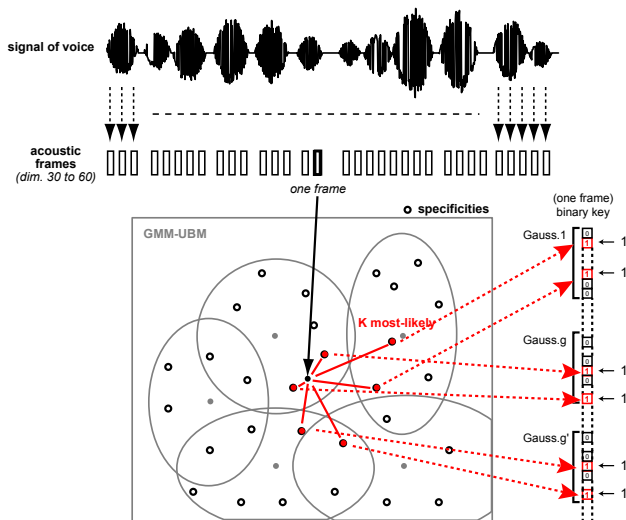
GMM-UBM high-dim. representation

i. With 0^{th} and 1^{th} order statistics



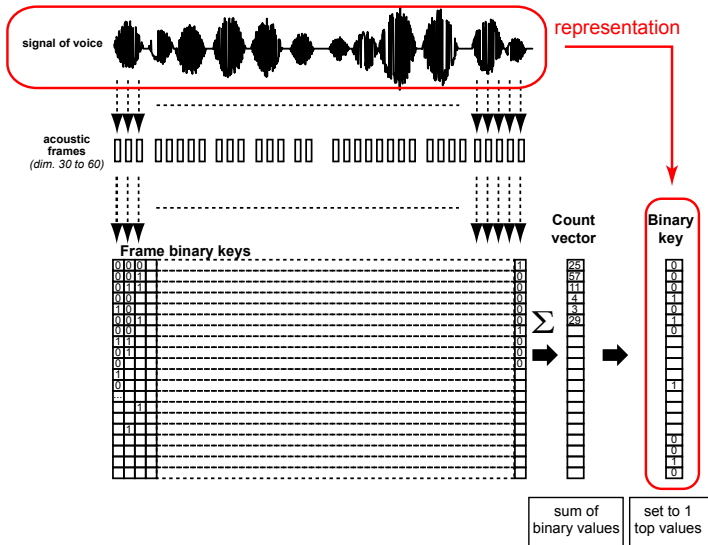
GMM-UBM high-dim. representation

ii. With speaker binary keys



GMM-UBM high-dim. representation

ii. With speaker binary keys



Dimensionality reduction (l-vector extractor)

i. With 0^{th} and 1^{th} order statistics

- A sequence \mathcal{X} of acoustic vectors \rightarrow supervector \mathbf{s}

Factor analysis : a probabilistic dimensionality reduction technique

- Supervector \mathbf{s} : dimension GF , where G number of GMM components and F feature dimension: **20 000 to 100 000 coefficients**.
- Factor Analysis Total Variability:

$$\mathbf{s} = \mathbf{m} + \mathbf{T}\mathbf{w} \quad (1)$$

- Factor \mathbf{w} (**i-vector**), dimension $p \ll GF$: **400 to 600 coefficients**.
- \mathbf{m} is the mean super-vector of the Universal Background Model (UBM),
- \mathbf{T} is the low-rank variability matrix $GF \times p$ ($p \ll GF$)
- Low-dimensional factor \mathbf{w} is assumed to have a standard normal prior distribution $\mathcal{N}(0, \mathbf{I})$.

Dehak, N., Kenny, P., Dehak, R., Dumouchel, P., and Ouellet, P. (2011). Front-End Factor Analysis for Speaker Verification. IEEE Transactions on Audio, Speech, and Language Processing.

Dimensionality reduction (l-vector extractor)

i. With 0^{th} and 1^{th} order statistics

Probabilistic determination of \mathbf{T}

by using EM-ML (*Expectation-Maximization - Maximum Likelihood*) procedure.

$$\mathbf{w} = (\mathbf{I} + \mathbf{T}^t \mathbf{\Sigma}^{-1} N_{\mathcal{X}} \mathbf{T})^{-1} \mathbf{T}^t \mathbf{\Sigma}^{-1} N_{\mathcal{X}} (\mathbf{s} - \mathbf{m}) \quad (2)$$

where

- $\mathbf{\Sigma}$ is the UBM diagonal covariance matrix,
- $N_{\mathcal{X}}$ is the $GF \times GF$ diagonal matrix composed of F blocks of $N_{\mathcal{X}}^{(g)} \mathbf{I}$ ($g = 1, \dots, G$) where $N_{\mathcal{X}}^{(g)}$ are the zero-order statistics estimated on the g -th Gaussian component of the UBM observing the set of feature vectors in the sequence \mathcal{X} .

Dimensionality reduction (I-vector extractor)

ii. With speaker binary keys

- Specificities: categorical variables (selected $\rightarrow 1$, otherwise $\rightarrow 0$).
 - ▶ *Multiple Correspondence Analysis* (MCA): equivalent-to-PCA dimensionality reduction technique for categorical variables.
 - ▶ Special case of binary variables (all variables have only two levels): MCA is equivalent to PCA on the binary coded vectors. Eigenvectors provided by both techniques are identical.
- Binary key \rightarrow dimensionality reduction \rightarrow “i-vector”.

Speaker and session modeling

Gaussian-PLDA

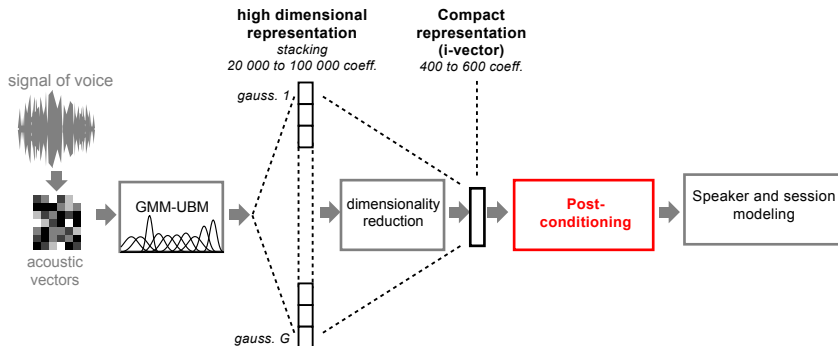
Given an i-vector \mathbf{w}

$$\mathbf{w} = \mu + \Phi \mathbf{y}_s + \varepsilon \quad \left\{ \begin{array}{l} \mathbf{y}_s \sim \mathcal{N}(0, \mathbf{I}_r) \text{ where } r \leq p \text{ and } \Phi \text{ is a } p \times r \text{ matrix} \\ \varepsilon \sim \mathcal{N}(0, \Lambda) \text{ where } \Lambda \text{ is a } p \times p \text{ matrix (full)} \\ \mathbf{y}_s \text{ and } \varepsilon \text{ statistically independent} \end{array} \right.$$

- Decision

$$\begin{aligned} \text{score}(\mathbf{w}_1, \mathbf{w}_2) &= \log \frac{P(\mathbf{w}_1, \mathbf{w}_2 | H0 : \text{same speaker})}{P(\mathbf{w}_1, \mathbf{w}_2 | H1 : \text{not the same})} \\ &= \log \frac{\mathcal{N}\left(\begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix} \mid \begin{bmatrix} \mu \\ \mu \end{bmatrix}, \begin{bmatrix} \Phi\Phi^t + \Lambda & \Phi\Phi^t \\ \Phi\Phi^t & \Phi\Phi^t + \Lambda \end{bmatrix}\right)}{\mathcal{N}\left(\begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix} \mid \begin{bmatrix} \mu \\ \mu \end{bmatrix}, \begin{bmatrix} \Phi\Phi^t + \Lambda & \mathbf{0} \\ \mathbf{0} & \Phi\Phi^t + \Lambda \end{bmatrix}\right)} \end{aligned} \quad (3)$$

Post-conditioning



Post-conditioning

$$\mathbf{w} \leftarrow \frac{\mathbf{A}^{-\frac{1}{2}} (\mathbf{w} - \mu)}{\left\| \mathbf{A}^{-\frac{1}{2}} (\mathbf{w} - \mu) \right\|} \quad (4)$$

where μ and \mathbf{A} denote the mean vector and a covariance matrix of a given i-vector corpus.

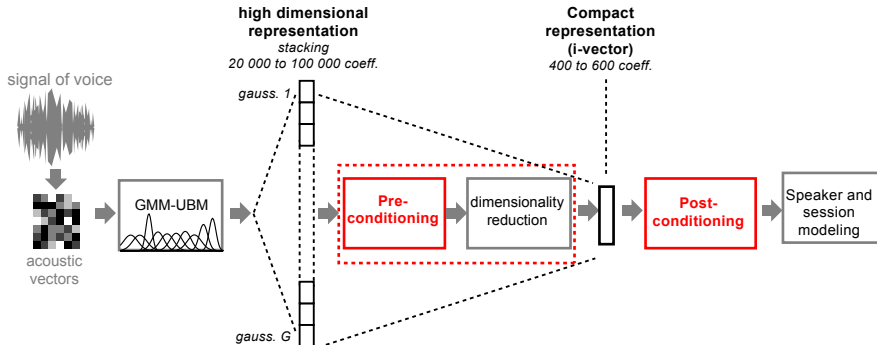
- $\mathbf{A} = \mathbf{\Sigma}$ (total covariance matrix) $\implies \mathbf{L}\mathbf{\Sigma}$ -conditioning
- $\mathbf{A} = \mathbf{W}$ (within-speaker covariance matrix) $\implies \mathbf{LW}$ -conditioning

... eventually iterated.

Parameters are computed for the i-vectors present in the training corpus and applied to test i-vectors.

- Improve Gaussianity of data and many other properties. See:
 - ▶ Garcia-Romero, D. and Espy-Wilson, C. Y. (2011). Analysis of i-vector length normalization in speaker recognition systems. In International Conference on Speech Communication and Technology, pages 249-252.
 - ▶ Bousquet, P.-M., Larcher, A., Matrouf, D., Bonastre, J.-F., and Plchot, O. (2012). Variance-Spectra based Normalization for I-vector Standard and Probabilistic Linear Discriminant Analysis. In Speaker and Language Recognition Workshop (IEEE Odyssey).
 - ▶ Bousquet, P.-M., Matrouf, D., and Bonastre, J.-F. (2011). Intersession compensation and scoring methods in the i-vectors space for speaker recognition. In International Conference on Speech Communication and Technology, pages 485-488.

Pre-conditioning



Pre-conditioning

With 0^{th} and 1^{th} order statistics

$$\mathbf{w} = (\mathbf{I} + \mathbf{T}^t \mathbf{\Sigma}^{-1} N_{\mathcal{X}} \mathbf{T})^{-1} \mathbf{T}^t \mathbf{\Sigma}^{-1} N_{\mathcal{X}} (\mathbf{s} - \mathbf{m}) \quad (5)$$

As $\mathbf{N}_{\mathcal{X}}$ and $\mathbf{\Sigma}$ are diagonal, this equation can be rewritten:

$$\mathbf{w} = \left[(\mathbf{I} + \tilde{\mathbf{T}}^t \mathbf{N}_{\mathcal{X}} \tilde{\mathbf{T}})^{-1} \tilde{\mathbf{T}}^t \mathbf{N}_{\mathcal{X}} \right] \mathbf{\Sigma}^{-\frac{1}{2}} (s - \mu) \quad (6)$$

where $\tilde{\mathbf{T}} = \mathbf{\Sigma}^{-\frac{1}{2}} \mathbf{T}$.

\Rightarrow Implicit pre-conditioning procedure

- $\mathbf{\Sigma}^{-\frac{1}{2}} (s - \mu)$
standardized version of supervector s (according to the world mean and covariance matrix).
- $(\mathbf{I} + \tilde{\mathbf{T}}^t \mathbf{N}_{\mathcal{X}} \tilde{\mathbf{T}})^{-1} \tilde{\mathbf{T}}^t \mathbf{N}_{\mathcal{X}}$
segment-dependent projection matrix (normalization depending on the amount of informations per Gaussian component, expressed in $\mathbf{N}_{\mathcal{X}}$).

Pre-conditioning

i. With 0^{th} and 1^{th} order statistics

To assess the importance of the pre-conditioning procedure:

- replacing F.A. extraction with the simple *Probabilistic Principal Component Analysis* PPCA on **non standardized** supervectors.
 - ▶ A **unique projection matrix** is estimated by using Gaussian EM-ML procedure then applied to any supervectors, without taking into account the amount of informations per Gaussian component.

Bishop, C. M. (2006). Pattern recognition and machine learning, volume 4. Springer.

Tipping, M. E. and Bishop, C. M. (1999). Probabilistic principal component analysis. Journal of the Royal Statistical Society: Series B (Statistical Methodology).

Pre-conditioning

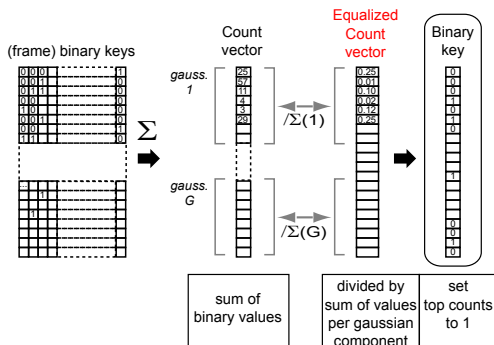
ii. With speaker binary keys (*Equalization per Gaussian Component*)

Let c denote the Gq -dimensional count vector of a segment.

Let $c_{g,k}$ denote the value of c for the k^{th} specificity of the g^{th} GMM component.

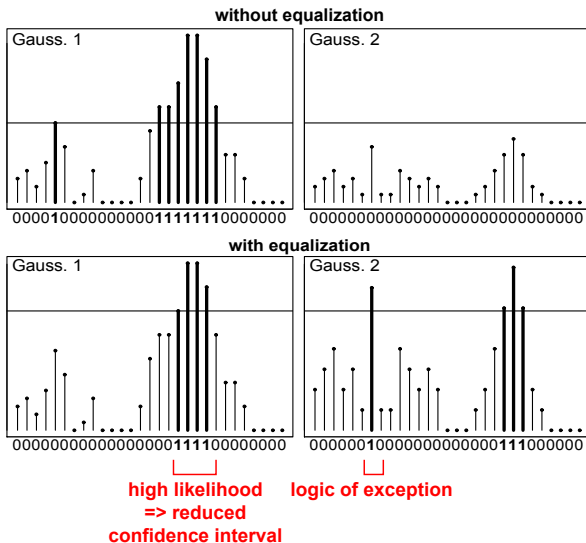
$$\hat{c}_{g,k} = \frac{c_{g,k}}{\sum_{k=1}^q c_{g,k}} \quad (7)$$

The count vector c becomes a local *frequency* vector \hat{c} (the sum of \hat{c} values for a Gaussian component is equal to 1).

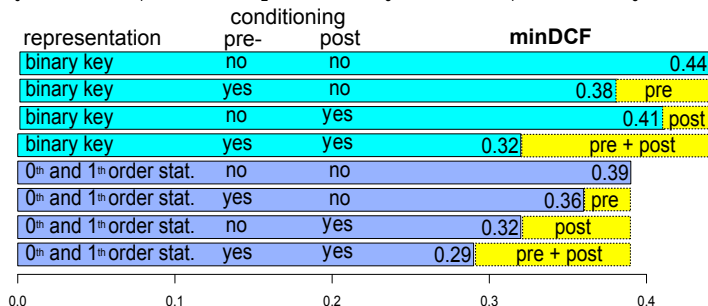
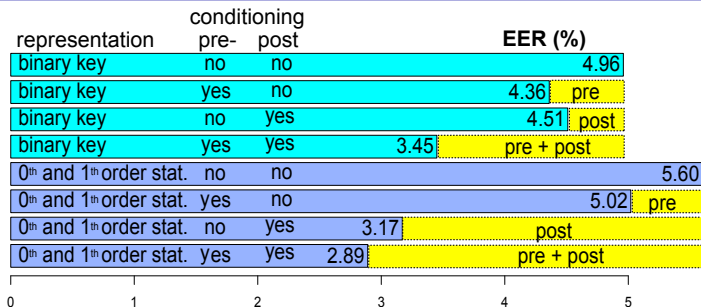


Pre-conditioning

ii. With speaker binary keys



Results



Conclusion

- Gaussian-PLDA model for total variability factors (*i-vectors*) achieves the best performance in speaker recognition, but these results are only obtained if transformations are applied to i-vectors, before modeling and after their generation by the dimensionality reduction technique.
- Today, the benefits of the **post-conditioning** procedures (comprised of standardization and length-normalization) are well known.
- The evaluation presented in this paper highlights the importance of **pre-conditioning** procedures, applied to high dimensional representations of utterance before their reduction.
 - ▶ With the binary key approach, the pre-conditioning procedure that we have designed for this representation turns out to be the most decisive of the two procedures, in terms of system accuracy.
 - ▶ With the Baum-Welch 0^{th} and 1^{th} statistics, we shows that FA dimensionality reduction includes an implicit pre-conditioning. Combining the two procedures (pre- and post-) is essential to achieve the state-of-the-art performance.
- This evaluation shows that the strategy of extending the GMM-UBM based speaker recognition systems through an additional stage of total variability factor is only relevant if **conditioning** procedures are applied **before** and **after** the dimensionality reduction stage.

Thank you ...