

Total Variability Model

Parametric model for explaining changes in feature space distribution across utterances
Used in tasks like speaker and language identification

Model Formulation

Feature vectors : $\mathbf{X}_u = \{\mathbf{x}_{ut}\}_{t=1}^{T_u}$, **ivector** : w_u

$$\mathbf{W}_u \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \mathbf{x}_{ut} | w_u \sim \text{GMM} \{ p_c, \boldsymbol{\mu}_{uc} = \boldsymbol{\mu}_c + \mathbf{T}_c w_u, \Sigma_c \}_{c=1}^C$$

$$\mathbf{M}_0 = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \vdots \\ \boldsymbol{\mu}_C \end{bmatrix}, \quad \mathbf{M}_u = \begin{bmatrix} \boldsymbol{\mu}_{u1} \\ \vdots \\ \boldsymbol{\mu}_{uC} \end{bmatrix}, \quad \mathbf{T} = \begin{bmatrix} \mathbf{T}_1 \\ \vdots \\ \mathbf{T}_C \end{bmatrix} \quad \mathbf{M}_u = \mathbf{M}_0 + \mathbf{T} w_u$$

Conventional Training Method

Expectation Maximization

Expectation : Evaluate the posterior distribution for ivectors

$$w_u | \mathbf{X}_u, \Theta^{(n-1)} \sim \mathcal{N} \left((\mathbf{I} + \mathbf{T}^T \Sigma^{-1} \mathbf{N}_u \mathbf{T})^{-1} \mathbf{T}^T \Sigma^{-1} \mathbf{F}_u, (\mathbf{I} + \mathbf{T}^T \Sigma^{-1} \mathbf{N}_u \mathbf{T})^{-1} \right)$$

Maximization : Maximize the following objective function

$$\mathcal{L}(\Theta) = \sum_{u=1}^U \mathbb{E}_{w_u | \mathbf{X}_u, \Theta^{(n-1)}} \left[\log P(\mathbf{X}_u, w_u | \mathbf{C}_u, \Theta^{(n)}) \right]$$

Expensive to evaluate all ivectors every iteration

Key Result

Assuming

$$\mathbf{N}_{uc} \approx T_u p_c, \quad \mathbf{T}^T \Sigma^{-1} \mathbf{P} \mathbf{T} + \frac{1}{T_u} \mathbf{I} \approx \mathbf{T}^T \Sigma^{-1} \mathbf{P} \mathbf{T} + \frac{1}{T} \mathbf{I}$$

Given

Covariance matrices Σ_c , UBM posteriors \mathbf{C}_u

Objective Function : Marginalized Likelihood

$$\sum_{u=1}^U \log P(\mathbf{X}_u | \mathbf{C}_u, \Theta) = \sum_{u=1}^U \log \left(\int_{w_u} P(\mathbf{X}_u, w_u | \mathbf{C}_u, \Theta) \right)$$

Is maximized for \mathbf{T} = Partial SVD of normalized statistics

Proposed Training Method

Randomized SVD of Normalized Statistics

$$\Sigma_c^{-1} = \mathbf{L}_c \mathbf{L}_c^T, \quad \tilde{\mathbf{F}}_u = \begin{bmatrix} \frac{1}{\sqrt{N_{u1}}} \mathbf{L}_1^T \mathbf{F}_{u1} \\ \vdots \\ \frac{1}{\sqrt{N_{uC}}} \mathbf{L}_C^T \mathbf{F}_{uC} \end{bmatrix}, \quad \tilde{\mathbf{F}} = \begin{bmatrix} \tilde{\mathbf{F}}_1 & \dots & \tilde{\mathbf{F}}_U \end{bmatrix}$$

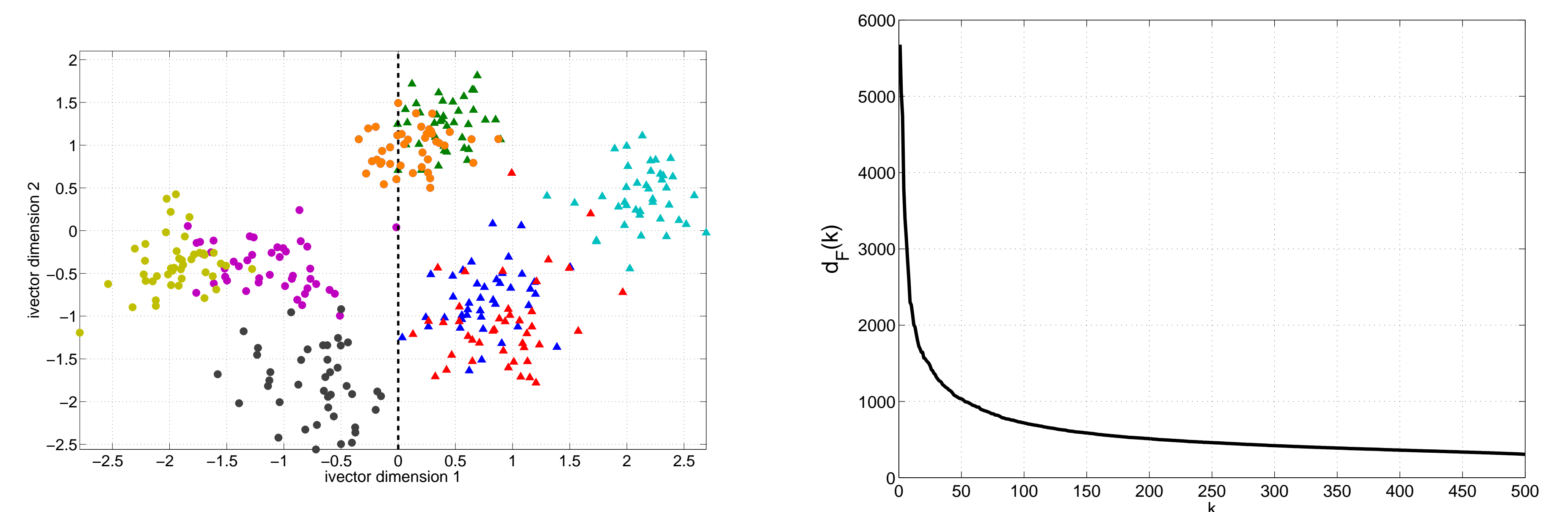
$$\tilde{\mathbf{F}} = \mathbf{U}_F \mathbf{D}_F \mathbf{V}_F^T, \quad \tilde{\mathbf{T}} = \mathbf{U}_F \mathbf{D}_F, \quad \mathbf{T} = \Sigma \mathbf{L} \mathbf{P}^{-\frac{1}{2}} \tilde{\mathbf{T}}$$

Key Advantages

Non-iterative, Efficient, Parallelizable, Interpretable

No need for iterations like in EM

i-vector dimensions ranked by singular values, no re-training for different i-vector dimensionality



Experimental Results

RATS language identification data : 6 languages, 16000 utterances per language, 30s each

	EM		SVD-1		SVD-2	
T_θ	33.5 hrs	30 min	30 min			
T_w	1.13 sec	0.12 sec	1.13 sec			
<i>Dur</i>	<i>EER</i>	<i>ACC</i>	<i>EER</i>	<i>ACC</i>	<i>EER</i>	<i>ACC</i>
10	8.30	83.50	8.78	83.60	9.26	81.15
5	10.04	80.80	10.40	80.55	10.94	78.05
3	13.41	72.75	13.59	74.25	14.97	70.55
1	21.59	56.45	22.13	58.35	22.61	54.65