

# Supervised group nonnegative matrix factorisation with similarity constraints and applications to speaker identification

Romain Serizel<sup>1</sup>, Victor Bisot<sup>2</sup>, Slim Essid<sup>2</sup>, Gaël  
Richard<sup>2</sup>

<sup>1</sup>LORIA, Univeristé de Lorraine, Inria, CNRS (France)

<sup>2</sup>LTCl, Télécom ParisTech, Université Paris-Saclay (France)

Monday 6<sup>th</sup>, March 2017



# Outline

- 1 Speaker identification : What ? Why ?
- 2 Task-driven group NMF
- 3 Conclusions

# Speaker identification

## Main goal

Identify a person from an audio recording



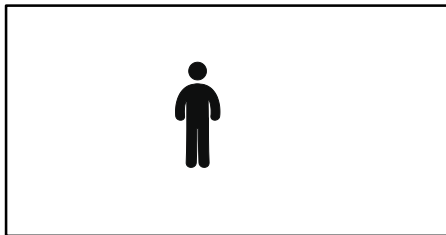
# Speaker identification

## Main goal

Identify a person from an audio recording



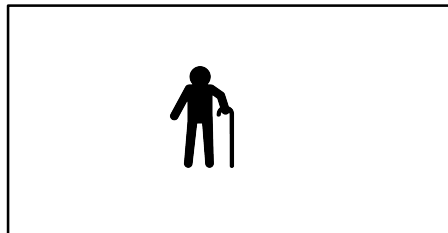
# Session concept



Recording variability

Icons made by Freepik from [www.flaticon.com](http://www.flaticon.com)

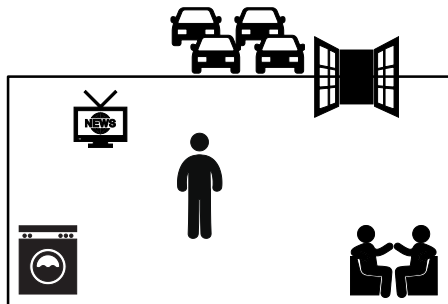
# Session concept



## Recording variability

- Aging,

## Session concept

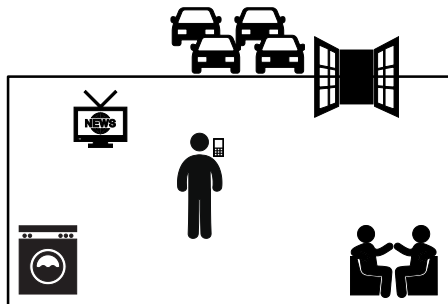


### Recording variability

- Aging,
- Perturbations,

Icons made by Freepik from [www.flaticon.com](http://www.flaticon.com)

## Session concept



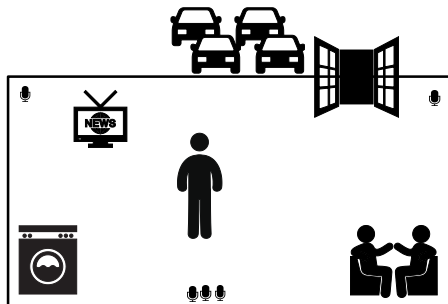
### Recording variability

- Aging,
- Perturbations,
- Microphones...

Icons made by Freepik from [www.flaticon.com](http://www.flaticon.com)



# Session concept



## Recording variability

- Aging,
- Perturbations,
- Microphones...

# Applications

## Audio indexing (broadcast show, conferences,..)

- Content retrieval
- Rich-text transcription

## Robust speech transcription

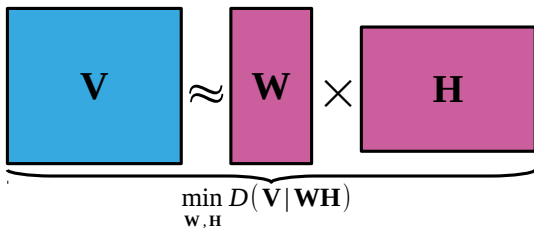
- Speaker adaptive training
- Speaker related feature/model adaptation

## Voice-based identification

- Soft biometrics

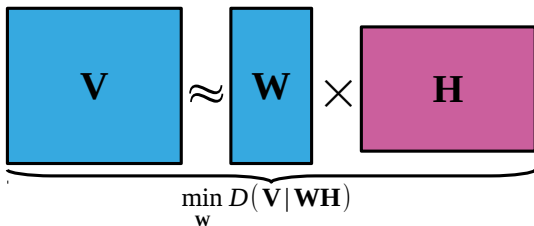
# Standard classification chain (1)

- Train the dictionary



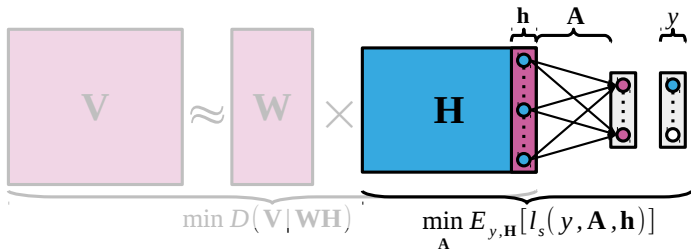
## Standard classification chain (2)

- Project data on the new space



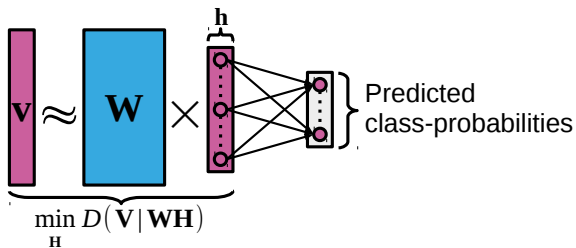
## Standard classification chain (3)

- Train the classifier on projected data



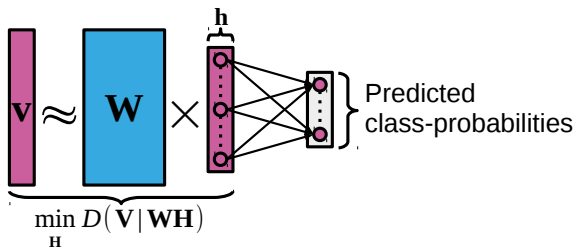
## Standard classification chain (4)

- Use the dictionary and classifier



## Standard classification chain (4)

- Use the dictionary and classifier



### Problem

$W$  and  $H$  are optimised according to a **reconstruction** criterion

# Task-driven NMF (1)

## General idea

Learn the dictionaries together with classifier parameters :

- Nested optimisation problem

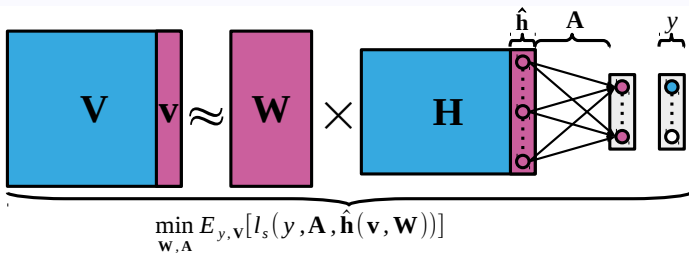
## Dictionary divergence

- Euclidean norm : **closed form solution** for dictionaries
  - Task driven dictionary learning (Mairal et al., 2012)
  - Application to audio scene analysis (Bisot et al., 2016)
- General  $\beta$ -divergence :
  - Application to source separation (Sprechmann et al., 2014)
  - Application to event detection (Bisot et al., 2017)



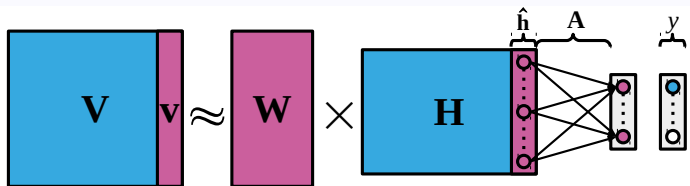
## Task-driven NMF (2)

### Task-driven NMF : General idea



## Task-driven NMF (2)

### Task-driven NMF : General idea

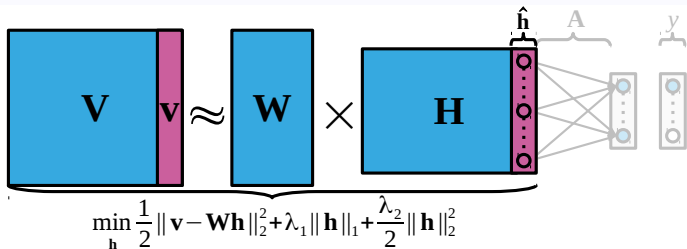


$$\min_{W, A} E_{y, v} [l_s(y, A, \hat{h}(v, W))] + \frac{\nu}{2} \|A\|_2^2$$

# Task-driven NMF : algorithm (1)

- For each new sample ( $\mathbf{v}$ )

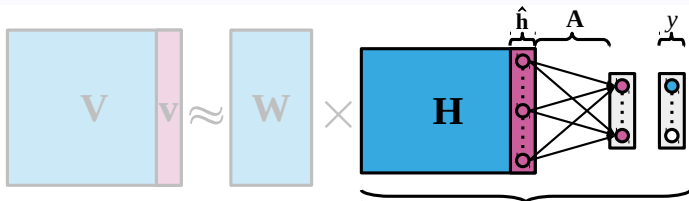
Project  $\mathbf{v}$  on  $\hat{\mathbf{h}}$



## Task-driven NMF : algorithm (2)

- For each new sample ( $\mathbf{v}$ )

Update the classifier parameters  $\mathbf{A}$

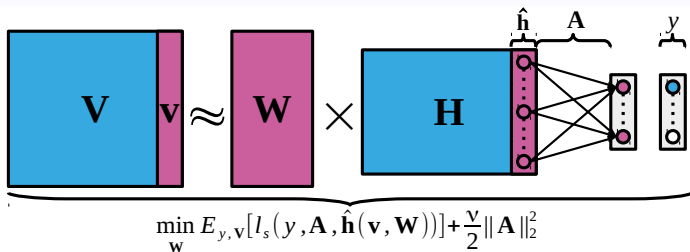


$$\min_{\mathbf{A}} E_{y, \mathbf{v}} [l_s(y, \mathbf{A}, \hat{\mathbf{h}}(\mathbf{v}, \mathbf{W}))] + \frac{\mathbf{v}}{2} \|\mathbf{A}\|_2^2$$

## Task-driven NMF : algorithm (3)

- For each new sample ( $\mathbf{v}$ )

Update the dictionary  $\mathbf{W}$



# Task-driven NMF in practice

## Implementation details

- Can be applied to sample or mini-batch
- Supports nonnegativity constraints for  $\mathbf{W}$  and  $\mathbf{H}$
- Dictionary ( $\mathbf{W}$ ) initialisation :
  - Random
  - NMF
  - Concatenated group NMF dictionaries (Serizel et al., 2016)

# Task-driven group NMF (1)

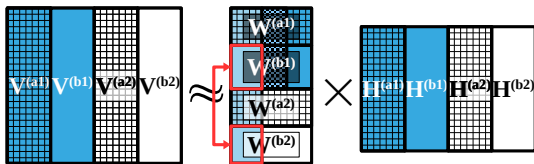
Include group NMF in a task-driven framework (Serizel et al., 2016)

The diagram shows a matrix  $V$  on the left, partitioned into four vertical blocks:  $V^{(a1)}$  (blue grid),  $V^{(b1)}$  (solid blue),  $V^{(a2)}$  (white grid), and  $V^{(b2)}$  (white). This matrix is approximately equal to the product of a subdictionary  $W$  and a coefficient matrix  $H$ . The subdictionary  $W$  is a 2x2 grid of blocks:  $W^{(a1)}$  (blue grid),  $W^{(b1)}$  (solid blue),  $W^{(a2)}$  (white grid), and  $W^{(b2)}$  (white). The coefficient matrix  $H$  is a 2x2 grid of blocks:  $H^{(a1)}$  (blue grid),  $H^{(b1)}$  (solid blue),  $H^{(a2)}$  (white grid), and  $H^{(b2)}$  (white). The approximation is indicated by a tilde symbol  $\approx$  and a multiplication symbol  $\times$ .

- Subdictionaries (related to a single speaker/session)

# Task-driven group NMF (1)

Include group NMF in a task-driven framework (Serizel et al., 2016)

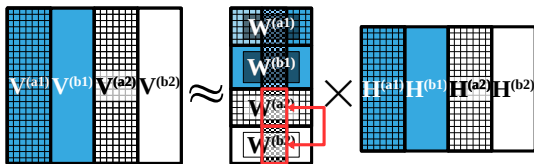


- Subdictionaries (related to a single speaker/session)
- Impose speaker/session similarity constraints



# Task-driven group NMF (1)

Include group NMF in a task-driven framework (Serizel et al., 2016)

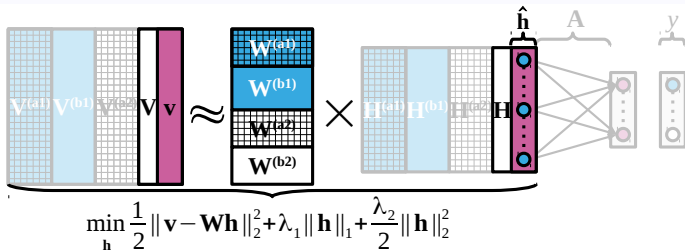


- Subdictionaries (related to a single speaker/session)
- Impose speaker/session similarity constraints

# Task-driven group NMF : algorithm (1)

- For each new sample ( $\mathbf{v}$ )

Project  $\mathbf{v}$  on  $\hat{\mathbf{h}}$

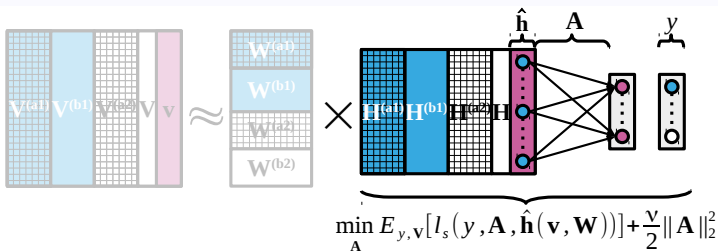


Source code is available at <https://github.com/rserizel/TGNMF>

## Task-driven group NMF : algorithm (2)

- For each new sample ( $\mathbf{v}$ )

Update the classifier parameters  $\mathbf{A}$

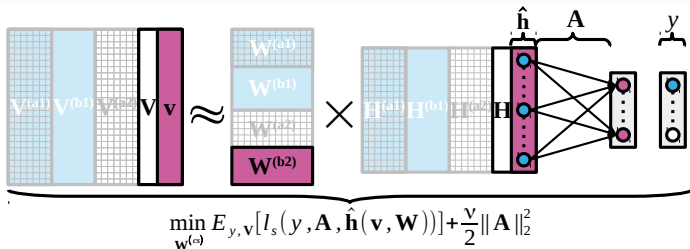


Source code is available at <https://github.com/rserizel/TGNMF>

# Task-driven group NMF : algorithm (3)

- For each new sample ( $\mathbf{v}$ )

Update the corresponding dictionary  $\mathbf{W}^{(cs)}$

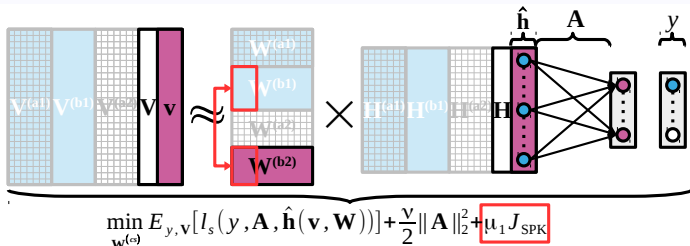


Source code is available at <https://github.com/rserizel/TGNMF>

## Task-driven group NMF : algorithm (3)

- **Speaker** similarity constraint

Update the corresponding dictionary  $W^{(cs)}$

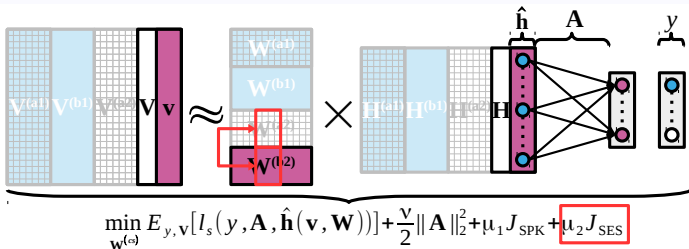


Source code is available at <https://github.com/rserizel/TGNMF>

# Task-driven group NMF : algorithm (3)

- **Session** similarity constraint

Update the corresponding dictionary  $W^{(cs)}$



Source code is available at <https://github.com/rserizel/TGNMF>

# Experiments

## Experiment setup

- Subset of the ESTER corpus ( $\approx$  6 hours training data)
- 132 constant-Q transform coefficients
- Initial dictionary obtained with (group-)NMF : 100 iterations
- Projection on  $\mathbf{h}$  with SPAMS toolbox<sup>a</sup>
- Classifier : multinomial logistic regression
- After 5 epochs : fix  $\mathbf{W}$ , train  $\mathbf{A}$  alone for 50 epochs

---

a. <http://spams-devel.gforge.inria.fr/>

## Results (1)

### Weighted F1-scores

	Initialisations			
	I-vector	NMF	GNMF <sub>0</sub>	GNMF <sub>c</sub>
Unsupervised	76.1%	75.6%	80.7%	81.7%
TNMF Tuning	—	79.9%	81.1%	81.9%

- GNMF<sub>0</sub> : group NMF **without** similarity constraints
- GNMF<sub>c</sub> : group NMF **with** similarity constraints (speaker and session)



## Results (2)

### Weighted F1-scores

		Initialisations	
		GNMF <sub>0</sub>	GNMF <sub>c</sub>
Unsupervised		80.7%	81.7%
Tuning	TNMF	81.1%	81.9%
	TGNMF <sub>0</sub>	81.7%	82.1%
	TGNMF <sub>c</sub>	82.0%	<b>82.2%</b>

- (T)GNMF<sub>0</sub> : (task-driven) group NMF **without** similarity constraints
- (T)GNMF<sub>c</sub> : (task-driven) group NMF **with** similarity constraints (speaker and session)

# Conclusions and future work

## NMF for speaker identification

- Can be competitive with I-vectors

## Task-driven NMF

- Large improvements for small dictionaries
- TGNMF<sub>c</sub> best performance to date on the corpus

## Future work

- Experiment with  $\beta$ -divergence
- Extend the framework to deep learning. . .

## Further readings

- V. Bisot, R. Serizel, S. Essid, and G. Richard. Feature Learning with Matrix Factorization Applied to Acoustic Scene Classification. HAL-archives ouvertes : working paper or preprint (hal-01362864), September 2016. URL <https://hal.archives-ouvertes.fr/hal-01362864>.
- V. Bisot, S. Essid, and G. Richard. Overlapping sound event detection with supervised nonnegative matrix factorization. In *Proc. of ICASSP*, 2017.
- J. Mairal, F. Bach, and J. Ponce. Task-driven dictionary learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(4) :791–804, 2012.
- R. Serizel, S. Essid, and G. Richard. Group nonnegative matrix factorisation with speaker and session variability compensation for speaker identification. In *Proc. of ICASSP*, 2016.
- Pablo Sprechmann, Alex M Bronstein, and Guillermo Sapiro. Supervised non-euclidean sparse nmf via bilevel optimization with applications to speech enhancement. In *Hands-free Speech Communication and Microphone Arrays (HSCMA), 2014 4th Joint Workshop on*, pages 11–15. IEEE, 2014.