



# Automatic Relevance Determination in Nonnegative Matrix Factorization with the $\beta$ -Divergence

VINCENT Y. F. TAN<sup>†</sup>, CÉDRIC FÉVOTTE<sup>§</sup>

<sup>†</sup>Department of ECE, University of Wisconsin-Madison (email: vtan@wisc.edu)  
<sup>§</sup>CNRS LTCI, TELECOM ParisTech (email: fevotte@telecom-paristech.fr)



## Nonnegative Matrix Factorization

- Given a **nonnegative** data matrix  $\mathbf{V} \in \mathbb{R}_+^{F \times N}$ , find a decomposition

$$\mathbf{V} \approx \hat{\mathbf{V}} \triangleq \mathbf{W}\mathbf{H}$$

where  $\mathbf{W} \in \mathbb{R}_+^{F \times K}$  (basis) and  $\mathbf{H} \in \mathbb{R}_+^{K \times N}$  (coefficients)

- Common dimension**  $K$  is chosen such that  $FK + KN \ll FN$
- Alternating minimization** usually performed

$$\min_{\mathbf{H} \geq 0} D(\mathbf{V}|\mathbf{W}\mathbf{H}), \quad \min_{\mathbf{W} \geq 0} D(\mathbf{V}|\mathbf{W}\mathbf{H})$$

- The measure of fit  $D(\mathbf{V}|\mathbf{W}\mathbf{H})$  is separable, i.e.,

$$D(\mathbf{V}|\hat{\mathbf{V}}) = \sum_{f=1}^F \sum_{n=1}^N d(v_{fn}|\hat{v}_{fn})$$

- We take the scalar cost function  $d(v_{fn}|\hat{v}_{fn})$  to be the so-called  **$\beta$ -divergence**. Special cases include:

$\beta = 0$  : Itakura-Saito divergence

$\beta = 1$  : (Generalized) Kullback-Leibler divergence

$\beta = 2$  : (Squared) Euclidean distance

## Main Contributions

- In practical applications, model order  $K$  is hard to choose
- If  $K$  is too large  $\Rightarrow$  overfitting  $\ominus$
- If  $K$  is too small  $\Rightarrow$  data does not fit well to the model  $\ominus$
- Bayesian NMF** model based on **Automatic Relevance Determination** to estimate  $K$  and get a better decomposition  $\ominus$
- ARD has been previously employed in Bayesian PCA (Bishop 1999) and sparse Bayesian learning (Tipping 2001).

## Majorization-Minimization (MM) for $\beta$ -NMF

- Algorithms are based on the MM framework
- Let the cost function to be minimized be  $C(\mathbf{H})$
- Build an **auxiliary function**  $G(\mathbf{H}|\tilde{\mathbf{H}})$  such that

$$G(\mathbf{H}|\tilde{\mathbf{H}}) \geq C(\mathbf{H}), \quad G(\tilde{\mathbf{H}}|\tilde{\mathbf{H}}) = C(\tilde{\mathbf{H}})$$

- Then, optimizing  $G(\cdot|\mathbf{H}^{(i)})$  yields

$$C(\mathbf{H}^{(i+1)}) \leq G(\mathbf{H}^{(i+1)}|\mathbf{H}^{(i)}) \leq G(\mathbf{H}^{(i)}|\mathbf{H}^{(i)}) = C(\mathbf{H}^{(i)})$$

- Hence, the MM updates consists in performing

$$\mathbf{H}^{(i+1)} = \arg \min_{\mathbf{H} \geq 0} G(\mathbf{H}|\mathbf{H}^{(i)}).$$

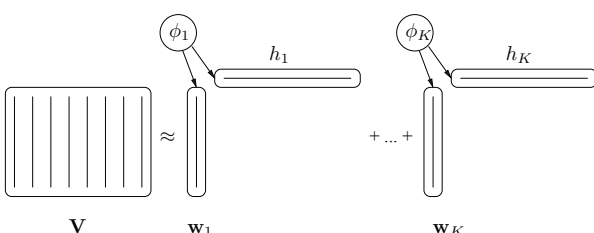
- For  $\beta$ -NMF,  $G(\mathbf{H}|\tilde{\mathbf{H}})$  can be found (Févotte and Idier 2011) and this amounts to the simple multiplicative update rule

$$h_{kn} = \tilde{h}_{kn} \left( \frac{p_{kn}}{q_{kn}} \right)^{\gamma(\beta)}$$

where  $p_{kn}, q_{kn}$  are simple functions of the previous iterate and the data and  $\gamma(\beta)$  is simple exponent in  $\beta$ .

## The Model for ARD in $\beta$ -NMF

- Main idea is to **tie** column  $\mathbf{w}_k$  and row  $h_k$  through their prior, a common (variance-like) **relevance weight**  $\phi_k$



- For  $\ell_2$ -ARD,  $\mathbf{W}$  and  $\mathbf{H}$  is assigned **Half-Normal** Priors

$$p(h_{kn}|\phi_k) = \sqrt{\frac{2}{\pi\phi_k}} \exp\left(-\frac{h_{kn}^2}{2\phi_k}\right)$$

- For  $\ell_1$ -ARD,  $\mathbf{W}$  and  $\mathbf{H}$  is assigned **Exponential** Priors

$$p(h_{kn}|\phi_k) = \frac{1}{\phi_k} \exp\left(-\frac{h_{kn}}{\phi_k}\right)$$

## The Model for ARD in $\beta$ -NMF (cont)

- Relevance weights**  $\phi_k$  assigned inverse-Gamma priors:

$$p(\phi_k; a, b) = \frac{b^a}{\Gamma(a)} \phi_k^{a-1} \exp\left(-\frac{b}{\phi_k}\right)$$

- Different  $\beta$ 's underlie different statistical noise models:

$$\text{IS-NMF : } \beta = 0 \quad v_{fn} \sim \mathcal{G}(v_{kn}|\alpha, \hat{v}_{fn}/\alpha)$$

$$\text{KL-NMF : } \beta = 1 \quad v_{fn} \sim \mathcal{P}(v_{kn}|\hat{v}_{fn})$$

$$\text{EUC-NMF : } \beta = 2 \quad v_{fn} \sim \mathcal{N}(v_{kn}|\hat{v}_{fn}, \sigma^2)$$

- The **likelihood** is given by

$$-\log p(\mathbf{V}|\mathbf{W}, \mathbf{H}) = \rho D_\beta(\mathbf{V}|\mathbf{W}\mathbf{H}) + \text{cst}$$

where  $\rho$  is some **regularization constant** reflecting our belief in the noise power, e.g.,  $\rho = 1/\sigma^2$  for  $\beta = 2$ .

## The Overall Cost Function (Posterior)

- Combining the likelihood and the priors gives the cost function (posterior):  $C(\mathbf{W}, \mathbf{H}, \phi) = -\log p(\mathbf{W}, \mathbf{H}, \phi|\mathbf{V}) =$

$$\rho D_\beta(\mathbf{V}|\mathbf{W}\mathbf{H}) + \sum_{k=1}^K \frac{1}{\phi_k} (f(\mathbf{w}_k) + f(h_k) + b) + c \log \phi_k + \text{cst.}$$

where  $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|_2^2$ ,  $c = (F + N)/2 + a + 1$  for  $\ell_2$ -ARD and  $f(\mathbf{x}) = \|\mathbf{x}\|_1$ ,  $c = F + N + a + 1$  for  $\ell_1$ -ARD

- MAP optimization**  $\Rightarrow$  some relevance weights converge to a **small constant** and corresponding components are **pruned**

- Optimizing only over  $\phi$  leads to the cost function  $C(\mathbf{W}, \mathbf{H}) =$

$$\rho D_\beta(\mathbf{V}|\mathbf{W}\mathbf{H}) + c \sum_{k=1}^N \log(f(\mathbf{w}_k) + f(h_k) + b) + \text{cst}$$

- This cost has connections to **reweighted  $\ell_1$  minimization** [Candès et al. 2008] and **group LASSO** [Yuan and Lin 2006]

## The Inference Algorithms

- Build **auxiliary functions** to optimize  $C(\mathbf{W}, \mathbf{H}, \phi)$  over  $\mathbf{H}$

- In the end, the **multiplicative updates** are

$$\ell_2 - \text{ARD : } h_{kn} = \tilde{h}_{kn} \left( \frac{p_{kn}}{q_{kn} + \tilde{h}_{kn}/(\rho\phi_k)} \right)^{\xi(\beta)}$$

$$\ell_1 - \text{ARD : } h_{kn} = \tilde{h}_{kn} \left( \frac{p_{kn}}{q_{kn} + 1/(\rho\phi_k)} \right)^{\gamma(\beta)}$$

- Updates of  $\phi_k$  proceed as follows  $\phi_k = \frac{f(\mathbf{w}_k) + f(h_k) + b}{c}$

- Cost function  $C(\mathbf{W}, \mathbf{H}, \phi)$  **decreases monotonically**

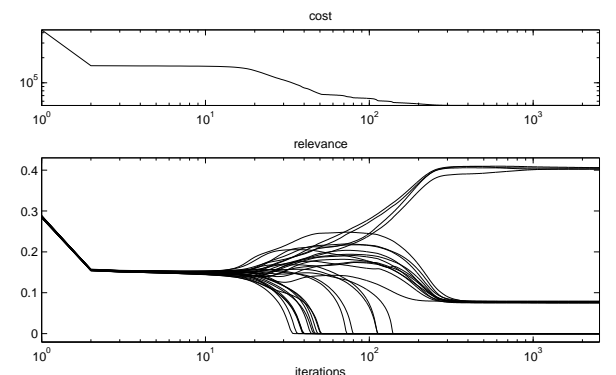
- Choose the hyperparameter  $b$  using the **method of moments**:

$$\ell_2 - \text{ARD : } \hat{b} = \frac{\pi(a-1)\hat{\mu}_V}{2K}$$

$$\ell_1 - \text{ARD : } \hat{b} = \left( \frac{(a-1)(a-2)\hat{\mu}_V}{K} \right)^{1/2}$$

## Noisy Swimmer Dataset (cont)

- Monotonic decrease in cost and evolution of relevances



## Music Decomposition Example

- Use the Itakura-Saito (IS) divergence ( $\beta = 0$ ) which is suited for audio applications (Févotte et al. 2009)

- Underlies a **generative statistical model** of superimposed Gaussian components in the squared STFT domain

- Sequence is composed of 4 piano notes, played all at once in the first measure and then played by pairs in all possible combinations in the subsequent measures.

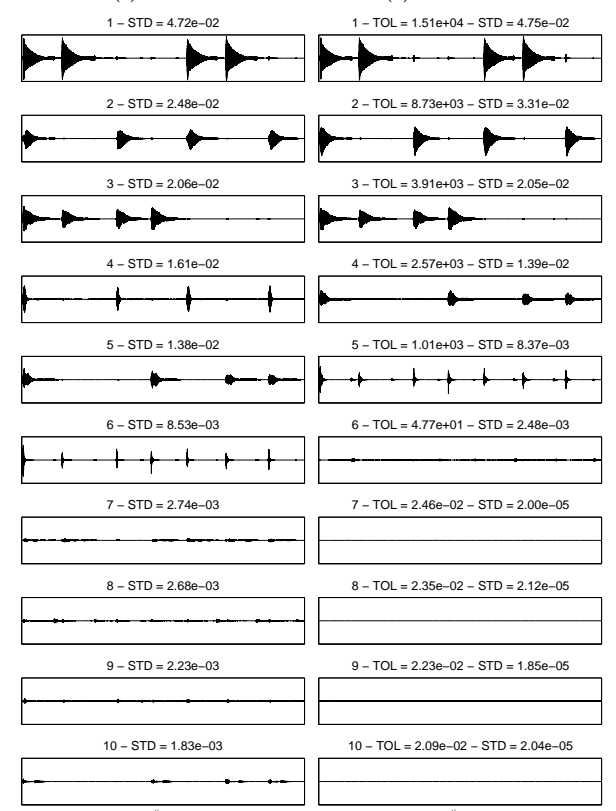


- Given an approximate factorization  $\mathbf{W}\mathbf{H}$  of the spectrogram  $v_{fn} = |x_{fn}|^2$ , the STFT estimate  $\hat{c}_{k,fn}$  of component  $k$  is

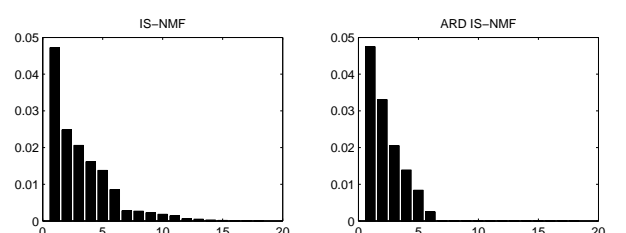
$$\hat{c}_{k,fn} = \frac{w_{fk} h_{kn}}{\sum_j w_{fj} h_{jn}} x_{fn}$$

(a) IS-NMF

(b) ARD IS-NMF



- Histogram shows that 6 components are retained by  $\ell_1$ -ARD



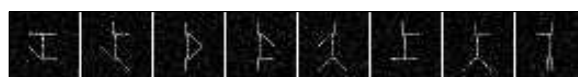
- First two components extract the **individual notes** and the next two components extract the **sound of hammer** hitting the strings and the sound produced by the **sustain pedal**

- Algorithms also applied to stock price prediction example. See <http://arxiv.org/abs/1111.6085>.

## Noisy Swimmer Dataset

- Synthetic dataset of  $N = 256$  images size  $F = 32 \times 32 = 1024$

- Each limb in four different positions. Example images:



- Dictionary learned using one run of  $\ell_1$ -ARD with  $a = 100$

