# Automatic Relevance Determination in Nonnegative Matrix Factorization (NMF)

Vincent Y. F. Tan [†] and Cédric Févotte [*]

[†]Laboratory for Information and Decision Systems (LIDS),
Massachusetts Institute of Technology.

[*]Laboratoire Traitement et Communication de l'Information (LTCI) ,
CNRS - TELECOM ParisTech.

Massachusetts
Institute of
Technology

SPARS Workshop (Apr 6, 2009)

# Introduction and Motivation

- Nonnegative matrix factorization (Lee and Seung 1999) is a popular technique for:

  1. Data analysis.

  2. Dimensionality reduction.

- NMF $\equiv$ non-subtractive, parts-based representation of nonnegative data.

- Nonnegative matrix factorization (Lee and Seung 1999) is a popular technique for:

  1. Data analysis.

  2. Dimensionality reduction.

- NMF $\equiv$ non-subtractive, parts-based representation of nonnegative data.

- Often, number of latent dimensions (or components) is assumed. Usually, this is not provided *a-priori*.

- We propose a Bayesian approach to estimate the latent dimensionality or model order.

- This is achieved by performing Automatic Relevance Determination (Mackay 1995).

- This has been used in Bayesian PCA (Bishop 1999) and sparse linear regression (Tipping 2001).

- Generally little literature about model order selection in NMF.

- Generally little literature about model order selection in NMF.

- Variational Bayesian methods have been proposed (Winther & Petersen 2007, Cemgil 2008) for NMF but such methods are usually computationally demanding.

  - Authors compute an approximation to the model evidence for every model order.

# Related Work

- Generally little literature about model order selection in NMF.

- Variational Bayesian methods have been proposed (Winther & Petersen 2007, Cemgil 2008) for NMF but such methods are usually computationally demanding.

  - Authors compute an approximation to the model evidence for every model order.

- The work is somewhat similar to multiplicative sparse NMF algorithms with a sparsity $\ell_1$ or $\ell_0$ regularizer (Hoyer 2004, Mørup et al. 2008) added to the objective.

# Notation – The NMF Model

Given a nonnegative data matrix

$$\mathbf{V} \in \mathbb{R}_+^{F \times N}.$$

## Notation – The NMF Model

Given a nonnegative data matrix

$$\mathbf{V} \in \mathbb{R}_+^{F \times N}.$$

Task: Find two nonnegative matrices

1. Basis Matrix $\quad \mathbf{W} \in \mathbb{R}_+^{F \times K}$

2. Activation Matrix $\quad \mathbf{H} \in \mathbb{R}_+^{K \times N}$

such that

$$\mathbf{V} \approx \widehat{\mathbf{V}} = \mathbf{W}\mathbf{H} = \sum_{k=1}^{K} \mathbf{w}_k h_k.$$

## Notation – The NMF Model

Given a nonnegative data matrix

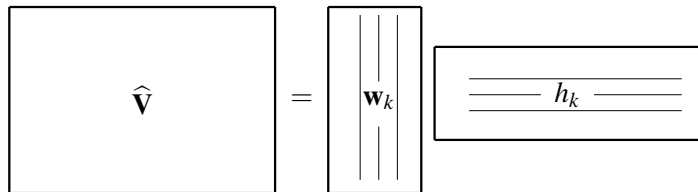$$\mathbf{V} \in \mathbb{R}_+^{F \times N}.$$
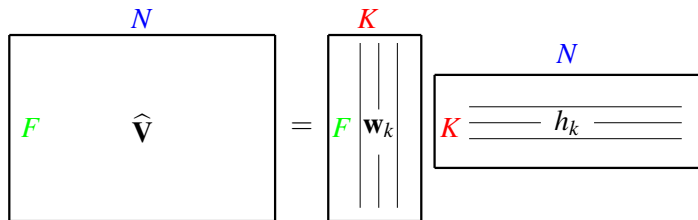
Task: Find two nonnegative matrices

      1. Basis Matrix      $\mathbf{W} \in \mathbb{R}_+^{F \times K}$

      2. Activation Matrix     $\mathbf{H} \in \mathbb{R}_+^{K \times N}$

such that

$$\mathbf{V} \approx \widehat{\mathbf{V}} = \mathbf{W}\mathbf{H} = \sum_{k=1}^{K} \mathbf{w}_k h_k.$$

Given a nonnegative data matrix

$$\mathbf{V} \in \mathbb{R}_+^{F \times N}.$$

Task: Find two nonnegative matrices

      1. Basis Matrix     $\mathbf{W} \in \mathbb{R}_+^{F \times K}$

      2. Activation Matrix     $\mathbf{H} \in \mathbb{R}_+^{K \times N}$

such that

$$\mathbf{V} \approx \widehat{\mathbf{V}} = \mathbf{WH} = \sum_{k=1}^{K} \mathbf{w}_k h_k.$$

# Nonnegative Matrix Factorization

Usually, a cost function $D(\cdot|\cdot)$ is minimized, i.e.,

$$\min_{\mathbf{W},\mathbf{H}} \; D(\mathbf{V}|\mathbf{W}\mathbf{H}) = \sum_{f=1}^{F} \sum_{n=1}^{N} d([\mathbf{V}]_{fn}|[\mathbf{W}\mathbf{H}]_{fn})$$

# Nonnegative Matrix Factorization

Usually, a cost function $D(\cdot|\cdot)$ is minimized, i.e.,

$$\min_{\mathbf{W},\mathbf{H}} \ D(\mathbf{V}|\mathbf{W}\mathbf{H}) = \sum_{f=1}^{F} \sum_{n=1}^{N} d([\mathbf{V}]_{fn}|[\mathbf{W}\mathbf{H}]_{fn})$$

where the cost (or divergence) $d$ can be

$$d_{EUC}(x|y) = \frac{1}{2}(x-y)^2, \qquad \text{(Euclidean cost)}$$

or

$$d_{KL}(x|y) = x \log\left(\frac{x}{y}\right) - x + y. \qquad \text{(KL-divergence)}$$

## Nonnegative Matrix Factorization

Usually, a cost function $D(\cdot|\cdot)$ is minimized, i.e.,

$$\min_{\mathbf{W}, \mathbf{H}} \ D(\mathbf{V}|\mathbf{W}\mathbf{H}) = \sum_{f=1}^{F} \sum_{n=1}^{N} d([\mathbf{V}]_{fn}|[\mathbf{W}\mathbf{H}]_{fn})$$

where the cost (or divergence) $d$ can be

$$d_{EUC}(x|y) = \frac{1}{2}(x-y)^2, \qquad \text{(Euclidean cost)}$$

or

$$d_{KL}(x|y) = x \log\left(\frac{x}{y}\right) - x + y. \qquad \text{(KL-divergence)}$$

Maximum Likelihood estimation of $\mathbf{W}$ and $\mathbf{H}$ corresponds to a particular noise model.

## Main Idea

- Set up a statistical model.
- Place precision-like scale parameters or relevance weights

$$\boldsymbol{\beta} \triangleq (\beta_1, \ldots, \beta_K) \in \mathbb{R}_+^K$$

on columns of $\mathbf{W}$ and rows of $\mathbf{H}$.

# Main Idea

- Set up a statistical model.

- Place precision-like scale parameters or relevance weights

$$\boldsymbol{\beta} \triangleq (\beta_1, \ldots, \beta_K) \in \mathbb{R}_+^K$$

  on columns of $\mathbf{W}$ and rows of $\mathbf{H}$.

- Automatic Relevance Determination.

## Main Idea

- Set up a statistical model.
- Place precision-like scale parameters or relevance weights

$$\boldsymbol{\beta} \triangleq (\beta_1, \ldots, \beta_K) \in \mathbb{R}_+^K$$

  on columns of $\mathbf{W}$ and rows of $\mathbf{H}$.

- Automatic Relevance Determination.



$$\widehat{\mathbf{V}} = \begin{array}{|c|} \hline \\ \mathbf{w}_1 \\ \\ \hline \end{array} \begin{array}{|c|} \hline \\ \text{---} h_1 \text{---} \\ \\ \hline \end{array} + \cdots + \begin{array}{|c|} \hline \\ \mathbf{w}_K \\ \\ \hline \end{array} \begin{array}{|c|} \hline \\ \text{---} h_K \text{---} \\ \\ \hline \end{array}$$

## Main Idea

- Set up a statistical model.

- Place precision-like scale parameters or relevance weights

$$\boldsymbol{\beta} \triangleq (\beta_1, \ldots, \beta_K) \in \mathbb{R}_+^K$$

  on columns of $\mathbf{W}$ and rows of $\mathbf{H}$.

- Automatic Relevance Determination.

# Main Idea

- Set up a statistical model.

- Place precision-like scale parameters or relevance weights

$$\boldsymbol{\beta} \triangleq (\beta_1, \ldots, \beta_K) \in \mathbb{R}_+^K$$

  on columns of $\mathbf{W}$ and rows of $\mathbf{H}$.

- Automatic Relevance Determination.

## Main Idea

- Set up a **statistical** model.

- Place **precision-like** scale parameters or **relevance weights**

$$\boldsymbol{\beta} \triangleq (\beta_1, \ldots, \beta_K) \in \mathbb{R}_+^K$$

  on columns of $\mathbf{W}$ and rows of $\mathbf{H}$.

- Automatic Relevance Determination.

$$\widehat{\mathbf{V}} = \begin{array}{|c|} \hline \\ \mathbf{w}_1 \\ \\ \hline \end{array} \quad \begin{array}{|c|} \hline \quad h_1 \quad \\ \hline \end{array} \quad + \quad \cdots \quad + \quad \begin{array}{|c|} \hline \\ \mathbf{w}_K \\ \\ \hline \end{array} \quad \begin{array}{|c|} \hline \quad h_K \quad \\ \hline \end{array}$$

$$\longleftarrow \beta_1 \qquad\qquad\qquad\qquad \longleftarrow \beta_K$$

- Perform **inference**.

- The number of $\beta_k$'s that remain below a certain threshold is the **model order**.

# Statistical Model – An Overview

- For the KL-divergence cost, find $\mathbf{W}^*, \mathbf{H}^*, \boldsymbol{\beta}^*$ such that the MAP criterion is optimized:

$$\min_{\mathbf{W},\mathbf{H},\boldsymbol{\beta}} \quad C_{\text{MAP}}(\mathbf{W},\mathbf{H},\boldsymbol{\beta}) \stackrel{\Delta}{=} - \underbrace{\log p(\mathbf{W},\mathbf{H},\boldsymbol{\beta}|\mathbf{V})}_{\text{posterior}}.$$

## Statistical Model – An Overview

- For the KL-divergence cost, find $\mathbf{W}^*, \mathbf{H}^*, \boldsymbol{\beta}^*$ such that the MAP criterion is optimized:

$$\min_{\mathbf{W},\mathbf{H},\boldsymbol{\beta}} \quad C_{\text{MAP}}(\mathbf{W},\mathbf{H},\boldsymbol{\beta}) \overset{\Delta}{=} -\underbrace{\log p(\mathbf{W},\mathbf{H},\boldsymbol{\beta}|\mathbf{V})}_{\text{posterior}}.$$

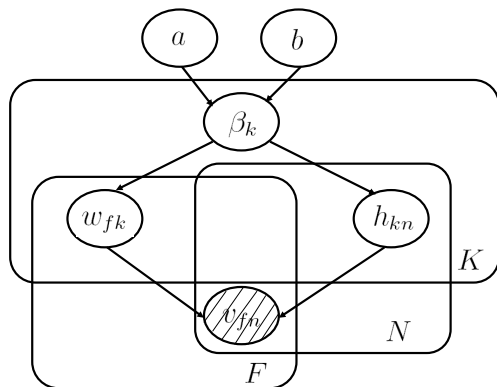where by Bayes' rule the posterior can be written as

$$-\underbrace{\log p(\mathbf{V}|\mathbf{W},\mathbf{H})}_{\text{likelihood}} -\underbrace{\log p(\mathbf{W}|\boldsymbol{\beta})}_{\text{prior on } \mathbf{W}} -\log \underbrace{p(\mathbf{H}|\boldsymbol{\beta})}_{\text{prior on } \mathbf{H}} -\underbrace{\log p(\boldsymbol{\beta}|a,b)}_{\text{prior on } \boldsymbol{\beta}}.$$

## Statistical Model – An Overview

- For the KL-divergence cost, find $\mathbf{W}^*, \mathbf{H}^*, \boldsymbol{\beta}^*$ such that the MAP criterion is optimized:

$$\min_{\mathbf{W}, \mathbf{H}, \boldsymbol{\beta}} \quad C_{\mathrm{MAP}}(\mathbf{W}, \mathbf{H}, \boldsymbol{\beta}) \triangleq - \underbrace{\log p(\mathbf{W}, \mathbf{H}, \boldsymbol{\beta} | \mathbf{V})}_{\text{posterior}}.$$

where by Bayes' rule the posterior can be written as

$$- \underbrace{\log p(\mathbf{V} | \mathbf{W}, \mathbf{H})}_{\text{likelihood}} - \underbrace{\log p(\mathbf{W} | \boldsymbol{\beta})}_{\text{prior on } \mathbf{W}} - \log \underbrace{p(\mathbf{H} | \boldsymbol{\beta})}_{\text{prior on } \mathbf{H}} - \underbrace{\log p(\boldsymbol{\beta} | a, b)}_{\text{prior on } \boldsymbol{\beta}}.$$

- Define the likelihood and priors and optimize $C_{\mathrm{MAP}}(\mathbf{W}, \mathbf{H}, \boldsymbol{\beta})$ efficiently.
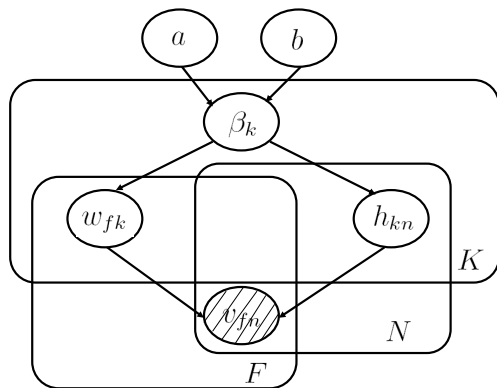
A Bayesian Network that describes our NMF statistical model.

# Dependences between Variables

A Bayesian Network that describes our NMF statistical model.



Need to specify:

- $p(\mathbf{V}|\mathbf{W}, \mathbf{H})$.
- $p(\mathbf{W}|\boldsymbol{\beta})$.
- $p(\mathbf{H}|\boldsymbol{\beta})$.
- $p(\boldsymbol{\beta}|a, b)$.

- $p(\mathbf{V}|\mathbf{W}, \mathbf{H})$.

- Assume that the likelihood of an element of the matrix $\mathbf{V}$, denoted $p(v_{fn}|\hat{v}_{fn})$, is given by a Poisson with rate $\hat{v}_{fn}$.

- $p(\mathbf{V}|\mathbf{W}, \mathbf{H})$.

- Assume that the likelihood of an element of the matrix $\mathbf{V}$, denoted $p(v_{fn}|\hat{v}_{fn})$, is given by a Poisson with rate $\hat{v}_{fn}$.

- This corresponds to the log-likelihood of $\mathbf{V}$ given $\mathbf{W}, \mathbf{H}$ as:

$$-\log p(\mathbf{V}|\mathbf{W}, \mathbf{H}) \overset{c}{=} D_{KL}(\mathbf{V}|\mathbf{WH}).$$

- Maximizing log-likelihood $\equiv$ Minimizing KL-divergence.

- Free from hyperparameters.

- $p(\mathbf{W}|\boldsymbol{\beta})$ and $p(\mathbf{H}|\boldsymbol{\beta})$.

- Independent half-normal priors over each column $k$ of $\mathbf{W}$ and row $k$ of $\mathbf{H}$.

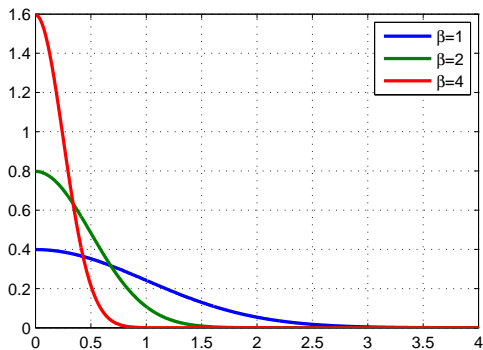# Bayesian NMF Model – Prior Models on $\mathbf{W}$ and $\mathbf{H}$

- $p(\mathbf{W}|\boldsymbol{\beta})$ and $p(\mathbf{H}|\boldsymbol{\beta})$.

- Independent half-normal priors over each column $k$ of $\mathbf{W}$ and row $k$ of $\mathbf{H}$.

- The priors are tied together through a single, common precision parameter $\beta_k$.

$$
\begin{aligned}
p(w_{fk}|\beta_k) &= \mathcal{HN}(w_{fk}|0, \beta_k^{-1}), \\
p(h_{kn}|\beta_k) &= \mathcal{HN}(h_{kn}|0, \beta_k^{-1}),
\end{aligned}
$$

where $\mathcal{HN}(w_{fk}|0, \beta_k^{-1})$ is the half-normal density with precision $\beta_k$.

- Least informative, High level of entropy.

# Half-Normal Densities



- Half-normal densities with different precision parameters $\beta$.
- The larger the $\beta$, the "peakier" the density $\Rightarrow$ Less relevant components will be sparse.

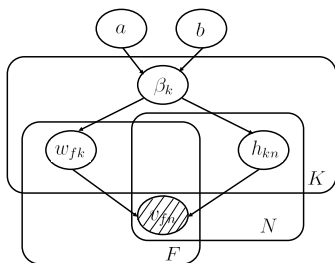- Each precision parameter $\beta_k$ is given a Gamma density:

$$p(\beta_k|a, b) = \mathcal{G}(\beta_k|a, b) = \frac{b^a}{\Gamma(a)} \beta_k^{a-1} \exp(-\beta_k b).$$
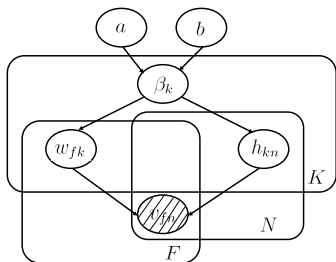
- This is the conjugate prior for $\beta_k$.

- $p(\mathbf{V}|\mathbf{W}, \mathbf{H})$ – KL-div.
- $p(\mathbf{W}|\boldsymbol{\beta})$ – Half-normal.
- $p(\mathbf{H}|\boldsymbol{\beta})$ – Half-normal.
- $p(\boldsymbol{\beta}|a, b)$ – Gamma.

# Recap: Dependences between Variables



- $p(\mathbf{V}|\mathbf{W}, \mathbf{H})$ – KL-div.
- $p(\mathbf{W}|\boldsymbol{\beta})$ – Half-normal.
- $p(\mathbf{H}|\boldsymbol{\beta})$ – Half-normal.
- $p(\boldsymbol{\beta}|a, b)$ – Gamma.

The MAP objective $C_{\mathrm{MAP}}$ is:

$$- \log p(\mathbf{W}, \mathbf{H}, \boldsymbol{\beta}|\mathbf{V}) \stackrel{\mathsf{c}}{=} D_{KL}(\mathbf{V}|\mathbf{W}, \mathbf{H}) +$$

$$+ \frac{1}{2} \sum_k \left[ \left( \sum_f w_{fk}^2 + \sum_n h_{kn}^2 + 2b \right) \beta_k - (F + N - 2(a-1)) \log \beta_k \right].$$

Tradeoff involving the size of the $\beta_k$'s.

- We have fully specified the Bayesian NMF statistical model.

# Inference

- We have fully specified the Bayesian NMF statistical model.

- Inference is done using efficient multiplicative updates, which ensures positivity.

- To update a parameter $\theta$ (e.g. an element of $\mathbf{W}$)

$$\theta \leftarrow \theta \, \frac{[\nabla_\theta C_{\mathrm{MAP}}(\theta)]_+}{[\nabla_\theta C_{\mathrm{MAP}}(\theta)]_-} \, .$$

where

$$\nabla_\theta C_{\mathrm{MAP}}(\theta) = [\nabla_\theta C_{\mathrm{MAP}}(\theta)]_+ - [\nabla_\theta C_{\mathrm{MAP}}(\theta)]_- .$$

Please refer to our paper for the details.

# Inference

- We have fully specified the Bayesian NMF statistical model.

- Inference is done using efficient multiplicative updates, which ensures positivity.

- To update a parameter $\theta$ (e.g. an element of $\mathbf{W}$)

$$\theta \leftarrow \theta \, \frac{[\nabla_\theta C_{\text{MAP}}(\theta)]_+}{[\nabla_\theta C_{\text{MAP}}(\theta)]_-}.$$

  where

$$\nabla_\theta C_{\text{MAP}}(\theta) = [\nabla_\theta C_{\text{MAP}}(\theta)]_+ - [\nabla_\theta C_{\text{MAP}}(\theta)]_-.$$

  Please refer to our paper for the details.

- The model order is

$$K_{\text{eff}} \triangleq |\{\beta_k \; : \; \beta_k < L\}|,$$

  and $L$ can be found analytically.

## ARD for NMF with KL-divergence cost

**Input** : Nonnegative data $\mathbf{V}$, fixed hyperparameters $a$, $b$.

**Output** : $\boldsymbol{\beta}$, $K_{\text{eff}}$, $\mathbf{W}$ and $\mathbf{H}$ s.t. $\mathbf{V} \approx \hat{\mathbf{V}} = \mathbf{W}\mathbf{H}$.

- Initialize $\mathbf{W}$ and $\mathbf{H}$ to nonnegative values.

- For $i = 1 : n_{iter}$

  - $\mathbf{H} \leftarrow \frac{\mathbf{H}}{\mathbf{W}^T \mathbf{1}_{F \times N} + \text{diag}(\boldsymbol{\beta})\mathbf{H}} \cdot \left[ \mathbf{W}^T \left( \frac{\mathbf{V}}{\mathbf{W}\mathbf{H}} \right) \right]$

  - $\mathbf{W} \leftarrow \frac{\mathbf{W}}{\mathbf{1}_{F \times N}\mathbf{H}^T + \mathbf{W}\text{diag}(\boldsymbol{\beta})} \cdot \left[ \left( \frac{\mathbf{V}}{\mathbf{W}\mathbf{H}} \right) \mathbf{H}^T \right]$

  - $\boldsymbol{\beta} \leftarrow \frac{F + N + 2(\mathbf{a} - 1)}{\mathbf{1}_{1 \times F}(\mathbf{W} \cdot \mathbf{W}) + (\mathbf{H} \cdot \mathbf{H})\mathbf{1}_{N \times 1} + 2\mathbf{b}}$

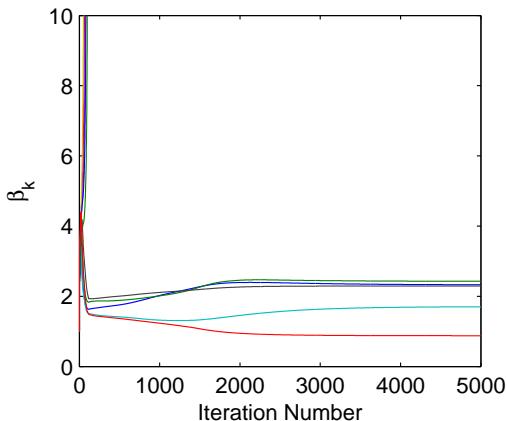- End For.

- Compute $K_{\text{eff}}$.

Linear in $F, N, K$.

# A Synthetic Dataset

Generated $\mathbf{V} \in \mathbb{R}_+^{100 \times 1000}$ with effective dimensionality $K_{\text{eff}} = 5$.

Set $K = 10$ and ran inference.

# A Synthetic Dataset

Generated $\mathbf{V} \in \mathbb{R}_+^{100 \times 1000}$ with effective dimensionality $K_{\text{eff}} = 5$.
Set $K = 10$ and ran inference.



- $K_{\text{eff}} = 5$ relevant components.

- $K - K_{\text{eff}} = 5$ irrelevant components.

# The Swimmer Dataset

- Swimmer dataset (Donoho and Stodden 2003).

- $N = 256$ images each of size $F = 32 \times 32$.

- A figure with four moving parts (limbs), each able to exhibit four articulations.

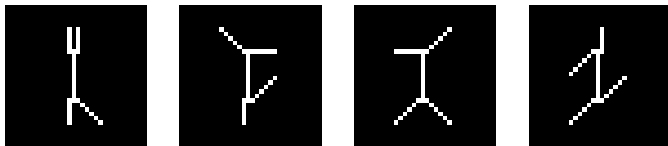- Fixed the shape parameter $a = 2$ and varied the scale $b$.



Figure: Sample images from the Swimmer dataset.

Figure: $K_{\mathrm{eff}}$ against $b$, the scale parameter of the Gamma prior.

Figure: The 16 limb positions are correctly recovered.

Figure: $K_{\mathrm{eff}}$ against $b$, the scale parameter of the Gamma prior.

# The MIT CBCL Dataset – Basis Images



$\beta_1 = 9.88$  $\beta_2 = 12.1$  $\beta_3 = 16.5$  $\beta_4 = 18.1$

$\beta_5 = 20.8$  $\beta_6 = 30.2$  $\beta_7 = 32.8$  $\beta_8 = 34.0$

$\beta_9 = 51.6$  $\beta_{10} = 59.7$  $\beta_{11} = 84.7$  $\beta_{12} = 85.5$
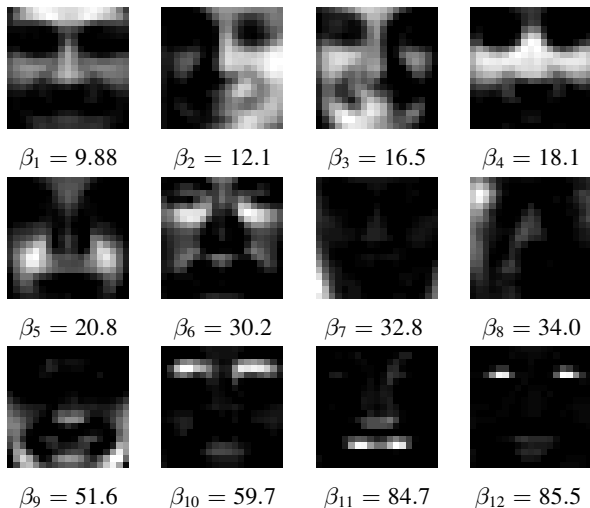
Figure: Basis images and corresponding relevance weights $\{\beta_k\}$.

- Bayesian approach that performs model order selection for NMF by borrowing ideas from ARD.

- Computationally cheap.

# Conclusions

- Bayesian approach that performs model order selection for NMF by borrowing ideas from ARD.

- Computationally cheap.

- Identify components that are 'relevant' for modeling the data.

- Experiments show that we are able to recover the latent dimensionality of synthetic and real data.

# Extensions

- Different cost functions (Euclidean, Itakura-Saito).

- Different prior models.

- Nonnegative Tensor Factorization.

$$\widehat{\mathbf{V}} = \sum_{k=1}^{K} \mathbf{w}_k^{(1)} \circ \mathbf{w}_k^{(2)} \circ h_k.$$

## Questions and Comments

- Thank you for your kind attention.

- Matlab<sup>©</sup> Code can be found online at
  `http://web.mit.edu/vtan/www/spars09`

- Authors can be reached at

  - Vincent Tan: `http://web.mit.edu/vtan/www/`

  - Cédric Févotte: `http://www.tsi.enst.fr/~fevotte/`