

---

**Bayesian Inference for Nonnegative Matrix Factorisation Models**

Ali Taylan Cemgil

**CUED/F-INFENG/TR.609**

July 2008

University of Cambridge  
Department of Engineering  
Trumpington Street  
Cambridge CB2 1PZ  
United Kingdom

Email: [atc27@cam.ac.uk](mailto:atc27@cam.ac.uk)

# Bayesian Inference for Nonnegative Matrix Factorisation Models\*

Ali Taylan Cemgil

Signal Processing and Comms. Lab, University of Cambridge  
Department of Engineering, Trumpington Street, Cambridge, CB2 1PZ, UK  
atc27@eng.cam.ac.uk

July 2008

## Abstract

We describe non-negative matrix factorisation (NMF) with a Kullback-Leibler error measure in a statistical framework, with a hierarchical generative model consisting of an observation and a prior component. Omitting the prior leads to standard NMF algorithms as special cases, where maximum likelihood parameter estimation is carried out via the Expectation-Maximisation (EM) algorithm. Starting from this view, we develop Bayesian extensions that facilitate more powerful modelling and allow full Bayesian inference via variational Bayes or Monte Carlo. Our construction retains conjugacy and enables us to develop models that fit better to real data while retaining attractive features of standard NMF such as fast convergence and easy implementation. We illustrate our approach on model order selection and image reconstruction.

## 1 Introduction

Nonnegative Matrix Factorisation (NMF) was introduced by Lee and Seung [16] as an alternative to k-means clustering and principal component analysis (PCA) for data analysis and compression. In NMF, given a  $W \times K$  nonnegative matrix  $X = \{x_{\nu,\tau}\}$  where  $\nu = 1:W$ ,  $\tau = 1:K$ , we seek positive matrices  $T$  and  $V$  such that

$$x_{\nu,\tau} \approx [TV]_{\nu,\tau} = \sum_i t_{\nu,i} v_{i,\tau}$$

where  $i = 1:I$ . We will refer to the  $W \times I$  matrix  $T$  as the *template matrix*, and  $I \times K$  matrix  $V$  the *excitation matrix*. The key property of NMF is that  $T$  and  $V$  are constrained to be positive matrices. This is in contrast with PCA, where there are no positivity constraints or k-means clustering where each column of  $V$  is constrained to be a unit vector. Subject to the positivity constraints, we seek a solution to the following minimisation problem

$$(T, V)^* = \arg \min_{T, V} D(X||TV) \quad (1)$$

Here, the function  $D$  is a suitably chosen error function. One particular choice for  $D$ , on which we will focus here, is the information (Kullback-Leibler) divergence, which we write as

$$D(X||\Lambda) = - \sum_{\nu,\tau} \left( x_{\nu,\tau} \log \frac{\lambda_{\nu,\tau}}{x_{\nu,\tau}} - \lambda_{\nu,\tau} + x_{\nu,\tau} \right) \quad (2)$$

Using Jensen's inequality [2] and concavity of  $\log x$ , it can be shown, that  $D(\cdot)$  is nonnegative and  $D(X||\Lambda) = 0$  if and only if  $X = \Lambda$ . The objective in (1) could be minimised by any suitable optimisation algorithm. Lee and Seung [16] have proposed a very efficient variational bound minimisation algorithm that has attractive convergence properties and which has been successfully applied in various applications in signal analysis and source separation, e.g. [23, 13, 1].

The interpretation of NMF, like SVD (singular value decomposition), as a low rank matrix approximation is sufficient for the derivation of a useful inference algorithm, yet this view arguably does not provide the complete picture about assumptions underlying the statistical properties of  $X$ . Therefore, we describe NMF from a statistical perspective as a hierarchical model. In our framework, the original nonnegative multiplicative update equations of NMF appear as an expectation-maximisation (EM) algorithm for maximum likelihood estimation of a conditionally Poisson model via *data augmentation*. Starting from this view, we develop Bayesian extensions that facilitate more powerful modelling and allow more sophisticated inference, such as Bayesian model selection. Inference in the resulting models can be carried out easily using variational (structured mean field) or Markov Chain Monte Carlo (Gibbs sampler). The resulting algorithms outperform existing NMF strategies and open up the way for a full Bayesian treatment for model selection via computation of the marginal likelihoods (the evidence), such as estimating the dimensions of the template matrix or regularising overcomplete representations via automatic relevance determination.

---

\*This research is funded by the Engineering and Physical Sciences Research Council (EPSRC) under the grant EP/D03261X/1).

## 2 The statistical perspective

The interpretation of NMF as a low rank matrix approximation is sufficient for the derivation of an inference algorithm, yet this view arguably does not provide the complete picture. In this section, we describe NMF from a statistical perspective. This view will pave the way for developing extensions that facilitate more realistic and flexible modelling as well as more sophisticated inference, such as Bayesian model selection.

Our first step is the derivation of the information divergence error measure from a maximum likelihood principle. We consider the following hierarchical model:

$$T \sim p(T|\Theta^t) \qquad V \sim p(V|\Theta^v) \qquad (3)$$

$$s_{\nu,i,\tau} \sim \mathcal{PO}(s_{\nu,i,\tau}; t_{\nu,i}v_{i,\tau}) \qquad x_{\nu,\tau} = \sum_i s_{\nu,i,\tau} \qquad (4)$$

Here,  $\mathcal{PO}(s; \lambda)$  denotes the Poisson distribution of the random variable  $s \in \mathbb{N}_0$  with nonnegative intensity parameter  $\lambda$  where

$$\mathcal{PO}(s; \lambda) = \exp(s \log \lambda - \lambda - \log \Gamma(s + 1))$$

and  $\Gamma(s + 1) = s!$  is the gamma function. The priors  $p(T|\cdot)$  and  $p(V|\cdot)$  will be specified later. We call the variables  $S_i = \{s_{\nu,i,\tau}\}$  *latent sources*. We can analytically marginalise out the latent sources  $S = \{S_1 \dots S_I\}$  to obtain the marginal likelihood

$$\begin{aligned} \log p(X|T, V) &= \log \sum_S p(X|S)p(S|T, V) = \log \prod_{\nu,\tau} \mathcal{PO}(x_{\nu,\tau}; \sum_i t_{\nu,i}v_{i,\tau}) \\ &= \sum_{\nu} \sum_{\tau} (x_{\nu,\tau} \log [TV]_{\nu,\tau} - [TV]_{\nu,\tau} - \log \Gamma(x_{\nu,\tau} + 1)) \end{aligned} \qquad (5)$$

This result follows from the well known *superposition property* of Poisson random variables [14], namely when  $s_i \sim \mathcal{PO}(s_i; \lambda_i)$  and  $x = s_1 + s_2 + \dots + s_I$  then the marginal probability is given by  $p(x) = \mathcal{PO}(x; \sum_i \lambda_i)$ . The maximisation of this objective in  $T$  and  $V$  is equivalent to the minimisation of the information divergence in (2). In the derivation of original NMF in [17], this objective is stated first; the  $S$  variables are introduced implicitly later during the optimisation on  $T$  and  $V$ . In the sequel, we show that this algorithm is actually equivalent to EM, ignoring the priors  $p(T|\cdot)$  and  $p(V|\cdot)$ .

### 2.1 Maximum Likelihood and the EM algorithm

The loglikelihood of the observed data  $X$  can be written as

$$\mathcal{L}_X(T, V) \equiv \log \sum_S p(X|S)p(S|T, V) \geq \sum_S q(S) \log \frac{p(X, S|T, V)}{q(S)} \equiv \mathcal{B}_{EM}[q] \qquad (6)$$

where  $q(S)$  is an instrumental distribution, that is arbitrary provided that the sum on the right exists;  $q$  can only vanish at a particular  $S$  only when  $p$  does so. Note that this defines a lower bound to the loglikelihood. It can be shown via functional derivatives and imposing the normalisation condition  $\sum_S q(S) = 1$  via Lagrange multipliers that the lower bound is tight for the exact posterior of the latent sources, i.e.

$$\arg \max_{q(S)} \mathcal{B}_{EM}[q] = p(S|X, T, V)$$

Hence the loglikelihood can be maximised iteratively

$$\begin{aligned} \text{E Step} \quad & q(S)^{(n)} = p(S|X, T^{(n-1)}, V^{(n-1)}) \\ \text{M Step} \quad & (T^{(n)}, V^{(n)}) = \arg \max_{T, V} \langle \log p(S, X|T, V) \rangle_{q(S)^{(n)}} \end{aligned}$$

Here,  $\langle f(x) \rangle_{p(x)} = \int p(x)f(x)dx$ , the expectation of some function  $f(x)$  with respect to  $p(x)$ . In the E step, we compute the posterior distribution of  $S$ . This defines a lower bound on the likelihood

$$\mathcal{B}^{(n)}(T, V|T^{(n-1)}, V^{(n-1)}) = \langle \log p(S, X|T, V) \rangle_{q(S)^{(n)}}$$

For many models in the exponential family, which includes (4), the expectation on the right depends on the sufficient statistics of  $q(S)^{(n)}$  and is readily available; in fact ‘‘calculating  $q(S)$ ’’ should be literally taken as calculating the sufficient statistics of  $q(S)$ . The lower bound is readily obtained as a function of these sufficient statistics and maximisation in the M Step yields a fixed point equation.

### 2.1.1 The E Step

To derive the posterior of the latent sources, we observe that

$$p(S|X, T, V) = p(S, X|T, V)/p(X|T, V) \quad (7)$$

For the model in (4), we have

$$\begin{aligned} \log p(S, X|T, V) = & \sum_{\nu} \sum_{\tau} \left( \sum_i (-t_{\nu,i} v_{i,\tau} + s_{\nu,i,\tau} \log(t_{\nu,i} v_{i,\tau})) \right. \\ & \left. - \log \Gamma(s_{\nu,i,\tau} + 1) \right) + \log \delta(x_{\nu,\tau} - \sum_i s_{\nu,i,\tau}) \end{aligned} \quad (8)$$

It follows from (4), (7), (8) and (5)

$$\begin{aligned} \log p(S|X, T, V) = & \sum_{\nu} \sum_{\tau} \left( \sum_i \left( s_{\nu,i,\tau} \log(t_{\nu,i} v_{i,\tau} / \sum_{i'} t_{\nu,i'} v_{i',\tau}) - \log \Gamma(s_{\nu,i,\tau} + 1) \right) \right. \\ & \left. + \log \Gamma(x_{\nu,\tau} + 1) + \log \delta(x_{\nu,\tau} - \sum_i s_{\nu,i,\tau}) \right) \\ = & \sum_{\nu} \sum_{\tau} \log \mathcal{M}(s_{\nu,1,\tau}, \dots, s_{\nu,I,\tau}; x_{\nu,\tau}, p_{\nu,1,\tau}, \dots, p_{\nu,I,\tau}) \end{aligned} \quad (9)$$

where  $p_{\nu,i,\tau} \equiv t_{\nu,i} v_{i,\tau} / \sum_{i'} t_{\nu,i'} v_{i',\tau}$  are the cell probabilities. Here,  $\mathcal{M}$  denotes a multinomial distribution defined by

$$\mathcal{M}(\mathbf{s}; x, \mathbf{p}) = \binom{x}{s_1 \ s_2 \ \dots \ s_I} p_1^{s_1} p_2^{s_2} \dots p_I^{s_I} \delta(x - \sum_i s_i) = \delta(x - \sum_i s_i) x! \prod_{i=1}^I \frac{p_i^{s_i}}{s_i!}$$

where  $\mathbf{s} = \{s_1, s_2, \dots, s_I\}$  and  $\mathbf{p} = \{p_1, p_2, \dots, p_I\}$  and  $p_1 + p_2 + \dots + p_I = 1$ . Here,  $p_i, i = 1 \dots I$  are the cell probabilities and  $x$  is the index parameter where  $s_1 + s_2 + \dots + s_I = x$ . The Kronecker delta function is defined by  $\delta(x) = 1$  when  $x = 0$ , and  $\delta(x) = 0$  otherwise. It is a standard result that the marginal mean is

$$\langle s_i \rangle = x p_i$$

i.e., the expected value each source  $s_i$  is a fraction of the observation, where the fraction is given by the corresponding cell probability.

### 2.1.2 The M Step

It is indeed a good news that the posterior has an analytic form since now the M step can be calculated easily

$$\begin{aligned} \langle \log p(S, X|T, V) \rangle_{p(S|X,T,V)} = & \sum_{\nu} \sum_{\tau} \left( \sum_i (-t_{\nu,i} v_{i,\tau} + \langle s_{\nu,i,\tau} \rangle \log(t_{\nu,i} v_{i,\tau})) \right. \\ & \left. - \langle \log \Gamma(s_{\nu,i,\tau} + 1) \rangle \right) + \left\langle \log \delta(x_{\nu,\tau} - \sum_i s_{\nu,i,\tau}) \right\rangle \end{aligned}$$

Fortunately, for maximisation w.r.t.  $T$  and  $V$ , the last two difficult terms are merely constant and we need only to maximise the simpler objective

$$Q(T, V) = \sum_{\nu} \sum_{\tau} \left( \sum_i (-t_{\nu,i} v_{i,\tau} + \langle s_{\nu,i,\tau} \rangle^{(n)} \log(t_{\nu,i} v_{i,\tau})) \right)$$

where we only need the expected value of the sources given by the previous values of the templates and excitations

$$\langle s_{\nu,i,\tau} \rangle^{(n)} = x_{\nu,\tau} \frac{t_{\nu,i}^{(n)} v_{i,\tau}^{(n)}}{\sum_{i'} t_{\nu,i'}^{(n)} v_{i',\tau}^{(n)}}$$

Maximisation of the objective  $Q$  and substituting  $\langle s_{\nu,i,\tau} \rangle^{(n)}$  gives the following fixed point equations:

$$\begin{aligned} \frac{\partial Q}{\partial t_{\nu,i}} = & - \sum_{\tau} v_{i,\tau}^{(n)} + \sum_{\tau} \langle s_{\nu,i,\tau} \rangle^{(n)} / t_{\nu,i} \\ t_{\nu,i}^{(n+1)} = & \sum_{\tau} \langle s_{\nu,i,\tau} \rangle^{(n)} / \sum_{\tau} v_{i,\tau}^{(n)} = t_{\nu,i}^{(n)} \frac{\sum_{\tau} x_{\nu,\tau} v_{i,\tau}^{(n)} / \sum_{i'} t_{\nu,i'}^{(n)} v_{i',\tau}^{(n)}}{\sum_{\tau} v_{i,\tau}^{(n)}} \end{aligned} \quad (10)$$

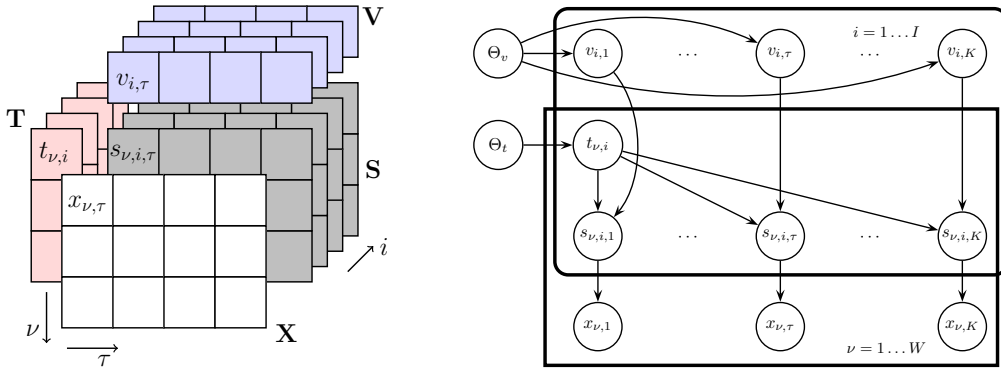


Figure 1: (Left) A schematic description of the NMF model with data augmentation. (Right) Graphical model with hyperparameters. Each source element  $s_{\nu,i,\tau}$  is Poisson distributed with intensity  $t_{\nu,i}v_{i,\tau}$ . The observations are given by  $x_{\nu,\tau} = \sum_i s_{\nu,i,\tau}$ . In matrix notation, we write  $X = \sum_i S_i$ . We can analytically integrate out over  $S$ . Due to superposition property of Poisson distribution, intensities add up and we obtain  $\langle X \rangle = TV$ . Given  $X$ , the NMF algorithm is shown to seek the maximum likelihood estimates of the templates  $T$  and excitations  $V$ . In our Bayesian treatment, we further assume elements of  $T$  and  $V$  are Gamma distributed with hyperparameters  $\Theta$ .

$$\begin{aligned} \frac{\partial Q}{\partial v_{i,\tau}} &= -\sum_{\nu} t_{\nu,i}^{(n)} + \sum_{\nu} \langle s_{\nu,i,\tau} \rangle^{(n)} / v_{i,\tau} \\ v_{i,\tau}^{(n+1)} &= \sum_{\nu} \langle s_{\nu,i,\tau} \rangle^{(n)} / \sum_{\nu} t_{\nu,i}^{(n)} = v_{i,\tau}^{(n)} \frac{\sum_{\nu} t_{\nu,i}^{(n)} x_{\nu,\tau} / \sum_{i'} t_{\nu,i'}^{(n)} v_{i',\tau}^{(n)}}{\sum_{\nu} t_{\nu,i}^{(n)}} \end{aligned} \quad (11)$$

The equations (10) and (11) are identical to the multiplicative update rules of [17]. However, our derivation via data augmentation obtains the same result as an EM algorithm. It is interesting to note that in the literature, NMF is often described as EM-like; here, we show that it is actually just an EM algorithm. We see that the efficiency of NMF is due to the fact that the  $W \times I \times K$  object  $\langle S \rangle$  need not be explicitly calculated as we only need its marginal statistics (sums across  $\tau$  or  $\nu$ ).

We note that our model is valid when  $X$  is integer valued. See [15] for a detailed discussion about consequences of this issue. Here, we assume that for nonnegative real valued  $\tilde{X}$  we only consider the integer part, i.e. we let  $\tilde{X} = X + E$  where  $E$  is a noise matrix with entries uniformly drawn in  $[0, 1)$ . In practice, this is not an obstacle when the entries of  $X$  are large.

The interpretation of NMF as a maximum likelihood method in a Poisson model is mentioned in the original NMF paper [16] and discussed in more detail by [13, 22]. Kameoka in [13] focuses on the optimisation and gives an equivalent description using auxiliary function maximisation. In contrast, the auxiliary variables can be viewed as model variables (the sources  $s$ ) that are analytically integrated out [22]. However, none of these approaches provided a full Bayesian treatment.

## 2.2 Hierarchical Prior Structure

Given the probabilistic interpretation, it is possible to propose various hierarchical prior structures to fit the requirements of an application. Here we will describe a simple choice where we have a conjugate prior

$$t_{\nu,i} \sim \mathcal{G}(t_{\nu,i}; a_{\nu,i}^t, b_{\nu,i}^t / a_{\nu,i}^t) \quad v_{i,\tau} \sim \mathcal{G}(v_{i,\tau}; a_{i,\tau}^v, b_{i,\tau}^v / a_{i,\tau}^v) \quad (12)$$

Here,  $\mathcal{G}$  denotes the density of a gamma random variable  $x \in \mathbb{R}_+$  with shape  $a \in \mathbb{R}_+$  and scale  $b \in \mathbb{R}_+$  defined by

$$\mathcal{G}(x; a, b) = \exp((a-1) \log x - x/b - \log \Gamma(a) - a \log b)$$

The primary motivation for choosing a Gamma distribution is computational convenience: Gamma distribution is the conjugate prior to Poisson intensity. The indexing highlights the most general case where there are individual parameters for each element  $t_{\nu,i}$  and  $v_{i,\tau}$ . Typically, we don't allow many free hyperparameters but tie them depending upon the requirements of an application. See figure 1 for an example. As an example, consider a model where we tie the hyperparameters such as  $a_{\nu,i}^t = a^t$ ,  $b_{\nu,i}^t = b^t$ ,  $a_{i,\tau}^v = a^v$  and  $b_{i,\tau}^v = b^v$  for  $i = 1 \dots I$ ,  $\nu = 1 \dots W$  and  $\tau = 1 \dots K$ . This model is simple to interpret, where each component of the templates and the excitations is drawn independently from the Gamma family shown in Figure 2. Qualitatively, the shape parameter  $a$  controls the *sparsity* of the representation. Remember that  $\mathcal{G}(x; a, b/a)$  has the mean  $b$  and standard deviation  $b/\sqrt{a}$ . Hence, for large  $a$ , all coefficients will have more or less the same magnitude  $b$  and typical representations will be full. In contrast, for small  $a$ , most of the coefficients will be

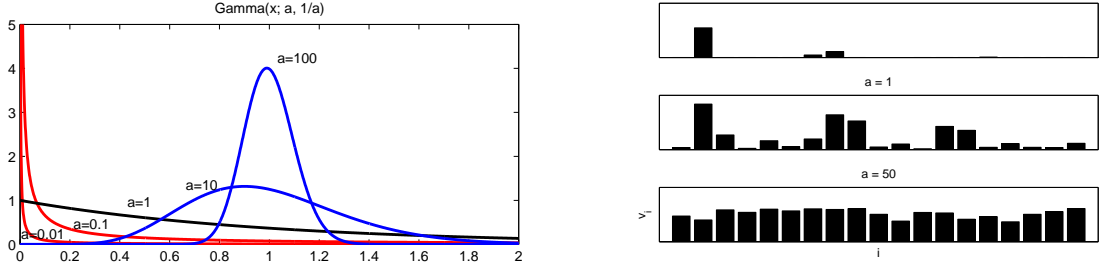


Figure 2: (Left) The family of densities  $p(v; a, b = 1) = \mathcal{G}(v; a, b/a)$  with the same mean  $\langle v \rangle = b = 1$ . (Right) Typical independent draws  $v_i \sim p(v; a, 1)$  for  $a = [0.05, 1, 50]$ . Small values of  $a$  enforce (sparser representations) and large values tie all values to be close to a nonzero mean (nonsparse representation).

very close to zero and only very few will be dominating, hence favoring a sparse representation. The scale parameter  $b$  is adapted to give the expected magnitude of each component.

To model missing data, i.e., when some of the  $x_{\nu, \tau}$  are not observed, we define a *mask* matrix  $M = \{m_{\nu, \tau}\}$ , the same size as  $X$  where,  $m_{\nu, \tau} = 0$ , if  $x_{\nu, \tau}$  is missing and 1 otherwise (see appendix A.4 for details). Using the mask variables, the observation model with missing data can be written as

$$p(X|S)p(S|T, V) = \prod_{\nu, \tau} (p(x_{\nu, \tau} | s_{\nu, 1:I, \tau}) p(s_{\nu, 1:I, \tau} | t_{\nu, 1:I}, v_{1:I, \tau}))^{m_{\nu, \tau}}$$

The hierarchical model in (12) is potentially more flexible than the basic model of (4), in that it allows a lot of freedom for more realistic modelling. First of all, the hyperparameters can be estimated from examples of a certain class of source to capture the invariant features. Another possibility is Bayesian model selection, where we can compare alternative models in terms of their marginal likelihood. This enables one to estimate the model order, for example, the optimum number of templates to represent a source.

### 3 Full Bayesian Inference

Below, we describe various interesting problems that can be cast to Bayesian inference problems. In signal analysis and feature extraction with NMF we may wish to calculate the posterior of templates and excitations given data and hyperparameters  $\Theta \equiv (\Theta^t, \Theta^v)$ . Another important quantity is the marginal likelihood (also known as the *evidence*) where

$$p(X|\Theta) = \int dT dV \sum_S p(X|S) p(S|T, V) p(T, V|\Theta)$$

The marginal likelihood can be used to estimate the hyperparameters, given examples of a source class

$$\Theta^* = \arg \max_{\Theta} p(X|\Theta)$$

or to compare two given models via *Bayes factors*

$$l(\Theta_1, \Theta_2) = \frac{p(X|\Theta_1)}{p(X|\Theta_2)}$$

This latter quantity is particularly useful for comparing different classes of models. Unfortunately, the integrations required can not be computed in closed form. In the sequel, we will describe the Gibbs sampler and variational Bayes as approximate inference strategies.

#### 3.1 Variational Bayes

We sketch here the Variational Bayes (VB) [9, 2] method to bound the marginal loglikelihood as

$$\mathcal{L}_X(\Theta) \equiv \log p(X|\Theta) \geq \sum_S \int d(T, V) q \log \frac{p(X, S, T, V|\Theta)}{q} = \langle \log p(X, S, V, T|\Theta) \rangle_q + H[q] \equiv \mathcal{B}_{VB}[q] \quad (13)$$

where,  $q = q(S, T, V)$  is an instrumental distribution and  $H[q]$  is its entropy. The bound is tight for the exact posterior  $q(S, T, V) = p(S, T, V|X, \Theta)$ , but as this distribution is complex we assume a factorised form for the instrumental distribution by ignoring some of the couplings present in the exact posterior

$$q(S, T, V) = q(S)q(T)q(V) = \left( \prod_{\nu, \tau} q(s_{\nu, 1:I, \tau}) \right) \left( \prod_{\nu, i} q(t_{\nu, i}) \right) \left( \prod_{i, \tau} q(v_{i, \tau}) \right) \equiv \prod_{\alpha \in \mathcal{C}} q_{\alpha}$$

where  $\alpha \in \mathcal{C} = \{\{S\}, \{T\}, \{V\}\}$  denotes set of disjoint clusters. Hence, we are no longer guaranteed to attain the exact marginal likelihood  $\mathcal{L}_X(\Theta)$ . Yet, the bound property is preserved and the strategy of VB is to optimise the bound. Although the best  $q$  distribution respecting the factorisation is not available in closed form, it turns out that a local optimum can be attained by the following fixed point iteration:

$$q_{\alpha}^{(n+1)} \propto \exp \left( \langle \log p(X, S, T, V|\Theta) \rangle_{q_{-\alpha}^{(n)}} \right) \quad (14)$$

where  $q_{-\alpha} = q/q_{\alpha}$ . This iteration monotonically improves the individual factors of the  $q$  distribution, i.e.  $\mathcal{B}[q^{(n)}] \leq \mathcal{B}[q^{(n+1)}]$  for  $n = 1, 2, \dots$  given an initialisation  $q^{(0)}$ . The order is not important for convergence, one could visit blocks in arbitrary order. However, in general, the attained fixed point depends upon the order of the updates as well as the starting point  $q^{(0)}(\cdot)$ . We choose the following update order in our derivations,

$$q(S)^{(n+1)} \propto \exp \left( \langle \log p(X, S, T, V|\Theta) \rangle_{q(T)^{(n)}q(V)^{(n)}} \right) \quad (15)$$

$$q(T)^{(n+1)} \propto \exp \left( \langle \log p(X, S, T, V|\Theta) \rangle_{q(S)^{(n+1)}q(V)^{(n)}} \right) \quad (16)$$

$$q(V)^{(n+1)} \propto \exp \left( \langle \log p(X, S, T, V|\Theta) \rangle_{q(S)^{(n+1)}q(T)^{(n+1)}} \right) \quad (17)$$

### 3.2 Variational update equations and sufficient statistics

The expectations of  $\langle \log p(X, S, T, V|\Theta) \rangle$  are functions of the sufficient statistics of  $q$  (see the expression in the appendix A.2). The update equation for the latent sources (15) leads to the following:

$$\begin{aligned} q(s_{\nu, 1:I, \tau}) &\propto \exp \left( \sum_i (s_{\nu, i, \tau} (\langle \log t_{\nu, i} \rangle + \langle \log v_{i, \tau} \rangle) - \log \Gamma(s_{\nu, i, \tau} + 1)) \right) \delta(x_{\nu, \tau} - \sum_i s_{\nu, i, \tau}) \\ &\propto \mathcal{M}(s_{\nu, 1, \tau}, \dots, s_{\nu, i, \tau}, \dots, s_{\nu, I, \tau}; x_{\nu, \tau}, p_{\nu, 1, \tau}, \dots, p_{\nu, i, \tau}, \dots, p_{\nu, I, \tau}) \\ p_{\nu, i, \tau} &= \exp(\langle \log t_{\nu, i} \rangle + \langle \log v_{i, \tau} \rangle) / \sum_i \exp(\langle \log t_{\nu, i} \rangle + \langle \log v_{i, \tau} \rangle) \end{aligned} \quad (18)$$

$$\langle s_{\nu, i, \tau} \rangle = x_{\nu, \tau} p_{\nu, i, \tau} \quad (19)$$

These equations are analogous to the multinomial posterior of EM given in 9; only the computation of cell probabilities is different. The excitation and template distributions and their sufficient statistics follow from properties of the gamma distribution.

$$\begin{aligned} q(t_{\nu, i}) &\propto \exp \left( (a_{\nu, i}^t + \sum_{\tau} \langle s_{\nu, i, \tau} \rangle - 1) \log(t_{\nu, i}) - \left( \frac{a_{\nu, i}^t}{b_{\nu, i}^t} + \sum_{\tau} \langle v_{i, \tau} \rangle \right) t_{\nu, i} \right) \\ &\propto \mathcal{G}(t_{\nu, i}; \alpha_{\nu, i}^t, \beta_{\nu, i}^t) \end{aligned}$$

$$\begin{aligned} \alpha_{\nu, i}^t &\equiv a_{\nu, i}^t + \sum_{\tau} \langle s_{\nu, i, \tau} \rangle & \beta_{\nu, i}^t &\equiv \left( \frac{a_{\nu, i}^t}{b_{\nu, i}^t} + \sum_{\tau} \langle v_{i, \tau} \rangle \right)^{-1} \\ \exp(\langle \log t_{\nu, i} \rangle) &= \exp(\Psi(\alpha_{\nu, i}^t)) \beta_{\nu, i}^t & \langle t_{\nu, i} \rangle &= \alpha_{\nu, i}^t \beta_{\nu, i}^t \end{aligned}$$

$$\begin{aligned} q(v_{i, \tau}) &\propto \exp \left( (a_{i, \tau}^v + \sum_{\nu} \langle s_{\nu, i, \tau} \rangle - 1) \log v_{i, \tau} - \left( \frac{a_{i, \tau}^v}{b_{i, \tau}^v} + \sum_{\nu} \langle t_{\nu, i} \rangle \right) v_{i, \tau} \right) \\ &\propto \mathcal{G}(v_{i, \tau}; \alpha_{i, \tau}^v, \beta_{i, \tau}^v) \end{aligned}$$

$$\begin{aligned} \alpha_{i, \tau}^v &\equiv a_{i, \tau}^v + \sum_{\nu} \langle s_{\nu, i, \tau} \rangle & \beta_{i, \tau}^v &\equiv \left( \frac{a_{i, \tau}^v}{b_{i, \tau}^v} + \sum_{\nu} \langle t_{\nu, i} \rangle \right)^{-1} \\ \exp(\langle \log v_{i, \tau} \rangle) &= \exp(\Psi(\alpha_{i, \tau}^v)) \beta_{i, \tau}^v & \langle v_{i, \tau} \rangle &= \alpha_{i, \tau}^v \beta_{i, \tau}^v \end{aligned}$$

### 3.3 Efficient implementation

One of the attractive features of NMF is easy and efficient implementation. In this section, we derive that the update equations of section 3.2 in compact matrix notation to illustrate that these attractive properties are retained for the full Bayesian treatment. A subtle but key point in the efficiency of the algorithm is that we can avoid explicitly storing and computing the  $W \times I \times K$  object  $\langle S \rangle$ , as we only need the marginal statistics during optimisation. Consider Eq. 18 and 19. We can write

$$\begin{aligned} \sum_{\tau} \langle s_{\nu,i,\tau} \rangle &= \sum_{\tau} x_{\nu,\tau} p_{\nu,i,\tau} \\ &= \exp(\langle \log t_{\nu,i} \rangle) \sum_{\tau} \left( x_{\nu,\tau} / \left( \sum_{i'} \exp(\langle \log t_{\nu,i'} \rangle) \exp(\langle \log v_{i',\tau} \rangle) \right) \right) \exp(\langle \log v_{i,\tau} \rangle) \\ \Sigma_t &= L_t .* ((X ./ (L_t L_v)) L_v^\top) \end{aligned}$$

Here, the denominator has to be nonzero. In the last line, we have represented the expression in compact notation where we define the following matrices:

$$\begin{aligned} E_t &= \{ \langle t_{\nu,i} \rangle \}, \quad L_t = \{ \exp(\langle \log t_{\nu,i} \rangle) \}, \quad \Sigma_t = \left\{ \sum_{\tau} \langle s_{\nu,i,\tau} \rangle \right\}, \quad A_t = \{ \alpha_{\nu,i}^t \}, \quad B_t = \{ b_{\nu,i}^t \}, \quad \alpha_t = \{ \alpha_{\nu,i}^t \}, \quad \beta_t = \{ \beta_{\nu,i}^t \} \\ E_v &= \{ \langle v_{i,\tau} \rangle \}, \quad L_v = \{ \exp(\langle \log v_{i,\tau} \rangle) \}, \quad \Sigma_v = \left\{ \sum_{\nu} \langle s_{\nu,i,\tau} \rangle \right\}, \quad A_v = \{ a_{i,\tau}^v \}, \quad B_v = \{ b_{i,\tau}^v \}, \quad \alpha_v = \{ \alpha_{i,\tau}^v \}, \quad \beta_v = \{ \beta_{i,\tau}^v \} \end{aligned}$$

The matrices subscripted with  $t$  are in  $\mathbb{R}_+^{W \times I}$  and with  $v$  are in  $\mathbb{R}_+^{I \times K}$ . For notational convenience, we define  $*$  and  $./$  as element wise matrix multiplication and division respectively and  $\mathbf{1}_W$  as a  $W \times 1$  vector of ones. After straightforward substitutions, we obtain the *variational nonnegative matrix factorisation* algorithm, that can compactly be expressed as in panel Algorithm 1).

---

#### Algorithm 1 Variational Nonnegative Matrix Factorisation

---

1: Initialise :

$$L_t^{(0)} = E_t^{(0)} \sim \mathcal{G}(\cdot; A_t, B_t ./ A_t) \quad L_v^{(0)} = E_v^{(0)} \sim \mathcal{G}(\cdot; A_v, B_v ./ A_v)$$

2: **for**  $n = 1 \dots$  MAXITER **do**

3: Source sufficient statistics

$$\begin{aligned} \Sigma_t^{(n)} &:= L_t^{(n-1)} .* ((X .* M) ./ (L_t^{(n-1)} L_v^{(n-1)})) L_v^{(n-1)\top} \\ \Sigma_v^{(n)} &:= L_v^{(n-1)} .* (L_t^{(n-1)\top} ((X .* M) ./ (L_t^{(n-1)} L_v^{(n-1)}))) \end{aligned}$$

4: Means

$$\begin{aligned} E_t^{(n)} &:= \alpha_t^{(n)} .* \beta_t^{(n)} & \alpha_t^{(n)} &= A_t + \Sigma_t^{(n)} & \beta_t^{(n)} &= 1. / (A_t ./ B_t + M E_v^{(n-1)\top}) \\ E_v^{(n)} &:= \alpha_v^{(n)} .* \beta_v^{(n)} & \alpha_v^{(n)} &= A_v + \Sigma_v^{(n)} & \beta_v^{(n)} &= 1. / (A_v ./ B_v + E_t^{(n)\top} M) \end{aligned}$$

5: Optional: Compute Bound (See appendix, (28))

6: Means of Logs

$$L_t^{(n)} = \exp(\Psi(\alpha_t^{(n)})) .* \beta_t^{(n)} \quad L_v^{(n)} = \exp(\Psi(\alpha_v^{(n)})) .* \beta_v^{(n)}$$

7: Optional: Update Hyperparameters (See appendix, section A.5)

8: **end for**

---

Similarly, an iterative conditional modes (ICM) algorithm can be derived to compute the maximum a-posteriori (MAP) solution (see section A.4)

$$V := (A_v + V .* (T^\top ((M .* X) ./ (TV)))) ./ (A_v ./ B_v + T^\top M) \quad (20)$$

$$T := (A_t + T .* (((M .* X) ./ (TV)) V^\top)) ./ (A_t ./ B_t + M V^\top) \quad (21)$$

Note that when the shape parameters go to zero, i.e.  $A_t, A_v \rightarrow \mathbf{0}$ , we obtain the maximum likelihood NMF algorithm.

### 3.4 Markov Chain Monte Carlo, the Gibbs sampler

Monte Carlo methods [10, 18] are powerful computational techniques to estimate expectations of form

$$E = \langle f(x) \rangle_{p(x)} \approx \frac{1}{N} \sum_{n=1}^N f(x^{(i)}) = \tilde{E}_N \quad (22)$$



where  $x^{(i)}$  are independent samples drawn from  $p(x)$ . Under mild conditions on  $f$ , the estimate  $\tilde{E}_N$  converges to the true expectation for  $N \rightarrow \infty$ . The difficulty here is obtaining independent samples  $\{x^{(i)}\}_{i=1\dots N}$  from complicated distributions.

The Markov Chain Monte Carlo (MCMC) techniques generate subsequent samples from a Markov chain defined by a *transition kernel*  $\mathcal{T}$ , that is one generates  $x^{(i+1)}$  conditioned on  $x^{(i)}$

$$x^{(i+1)} \sim \mathcal{T}(x|x^{(i)})$$

Note that the transition kernel  $\mathcal{T}$  is not needed explicitly in practice; all is needed is a procedure to sample a new configuration given the previous one. Perhaps surprisingly, even though subsequent samples are correlated, provided that  $\mathcal{T}$  satisfies certain ergodicity conditions, (22) remains still valid and estimated expectations converge to their true values when number of samples  $N$ , goes to infinity [10]. To design a transition kernel  $\mathcal{T}$  such that the desired distribution is the stationary distribution, i.e.  $p(x) = \int dx' \mathcal{T}(x|x')p(x')$ , many alternative strategies can be employed [18]. One particularly convenient and simple procedure is the Gibbs sampler where one samples each block of variables from *full conditional distributions*. For the NMF model, a possible Gibbs sampler is

$$S^{(n+1)} \sim p(S|T^{(n)}, V^{(n)}, X, \Theta) \quad (23)$$

$$T^{(n+1)} \sim p(T|V^{(n)}, S^{(n+1)}, X, \Theta) \quad (24)$$

$$V^{(n+1)} \sim p(V|S^{(n+1)}, T^{(n+1)}, X, \Theta) \quad (25)$$

Note that this procedure implicitly defines a transition kernel  $\mathcal{T}(\cdot|\cdot)$ . It can be shown [10] that the stationary distribution of  $\mathcal{T}$  is the exact posterior  $p(S, T, V|X, \Theta)$ . Eventually, the Gibbs sampler converges regardless of the order that the blocks are visited, provided that each block is visited infinitely often in the limit  $n \rightarrow \infty$ . However the rate of convergence is very difficult to assess as it depends upon the order of the updates as well as the starting configuration  $(T^{(0)}, V^{(0)}, S^{(0)})$ . It is instructive to contrast above equations (23)-(25) with the variational update equations(15)-(17) : algorithmically the two approaches are quite similar. The pseudo-code is given in Algorithm 2.

### 3.4.1 Marginal Likelihood estimation with Chib's method

The marginal likelihood can be estimated from the samples generated by the Gibbs sampler using a method proposed by Chib [6]. Suppose we have run the block Gibbs sampler until convergence and have  $N$  samples

$$\{T^{(n)}\}_{n=1:N} \quad \{V^{(n)}\}_{n=1:N} \quad \{S^{(n)}\}_{n=1:N}$$

The marginal likelihood is (omitting hyperparameters  $\Theta$ )

$$p(X) = \frac{p(V, T, S, X)}{p(V, T, S|X)} \quad (26)$$

This equation holds for *all* points  $(V, T, S)$ . We choose a point in the configuration space; provided that the distribution is unimodal, a good candidate is a configuration near the mode  $(\tilde{T}, \tilde{V}, \tilde{S})$ . The numerator in (26) is easy to evaluate. The denominator is

$$\begin{aligned} p(\tilde{V}, \tilde{T}, \tilde{S}|X) &= p(\tilde{V}|\tilde{T}, \tilde{S}, X)p(\tilde{T}|\tilde{S}, X)p(\tilde{S}|X) \\ &= p(\tilde{V}|\tilde{T}, \tilde{S})p(\tilde{T}|\tilde{S})p(\tilde{S}|X) \end{aligned}$$

The first term is the full conditional so it is available for the Gibbs sampler. The third term is

$$p(\tilde{S}|X) = \int dV dT p(\tilde{S}|V, T, X)p(V, T|X) \approx \frac{1}{N} \sum_{n=1}^N p(\tilde{S}|V^{(n)}, T^{(n)}, X) \quad (27)$$

The second term is trickier

$$p(\tilde{T}|\tilde{S}) = \int dV p(\tilde{T}|V, \tilde{S})p(V|\tilde{S})$$

The first term here is the full conditional. However, the original Gibbs run gives us only samples from  $p(V|X)$ , not  $p(V|\tilde{S})$ . The idea is to run the Gibbs sampler for  $M$  further iterations where we sample from  $(V_{\tilde{S}}^{(m)}, T_{\tilde{S}}^{(m)}) \sim p(V, T|S = \tilde{S})$ , i.e. with  $S$  clamped at  $\tilde{S}$ . The resulting estimate is

$$p(\tilde{T}|\tilde{S}) \approx \frac{1}{M} \sum_{m=1}^M p(\tilde{T}|V_{\tilde{S}}^{(m)}, \tilde{S})$$

Chib’s method estimates the marginal likelihood as follows:

$$\begin{aligned}
\log p(X|\Theta) &= \log p(\tilde{V}, \tilde{T}, \tilde{S}, X|\Theta) - \log p(\tilde{V}, \tilde{T}, \tilde{S}|X, \Theta) \\
&\approx \log p(\tilde{V}, \tilde{T}, \tilde{S}, X|\Theta) - \log p(\tilde{V}|\tilde{T}, \tilde{S}, \Theta) \\
&\quad - \log \sum_{m=1}^M p(\tilde{T}|V_{\tilde{S}}^{(m)}, \tilde{S}, \Theta) - \log \sum_{n=1}^N p(\tilde{S}|V^{(n)}, T^{(n)}, X, \Theta) + \log(MN)
\end{aligned}$$

---

**Algorithm 2** Gibbs sampler for Nonnegative Matrix Factorisation

---

1: Initialize :

$$T^{(0)} = \sim \mathcal{G}(\cdot; A_t, B_t) \quad V^{(0)} \sim \mathcal{G}(\cdot; A_v, B_v)$$

2: **for**  $n = 1 \dots \text{MAXITER}$  **do**

3:   Sample Sources

4:   **for**  $\tau = 1 \dots K, \nu = 1 \dots W$  **do**

5:      $p_{\nu,1:I,\tau}^{(n)} = T^{(n-1)}(\nu, 1:I) \cdot * V^{(n-1)}(1:I, \tau)^\top ./ (T^{(n-1)}(\nu, 1:I) V^{(n-1)}(1:I, \tau))$

6:      $S^{(n)}(\nu, 1:I, \tau) \sim \mathcal{M}(s_{\nu,1:I,\tau}; x_{\nu,\tau}, p_{\nu,1:I,\tau}^{(n)})$

7:   **end for**

$$\Sigma_t^{(n)} = \sum_{\tau} S_{\nu,i,\tau}^{(n)} \quad \Sigma_v^{(n)} = \sum_{\nu} S_{\nu,i,\tau}^{(n)}$$

8:   Sample Templates

$$\begin{aligned}
\alpha_t^{(n)} &= A_t + \Sigma_t^{(n)} & \beta_t^{(n)} &= 1./ \left( A_t ./ B_t + \mathbf{1}_W (V^{(n-1)} \mathbf{1}_K)^\top \right) \\
T^{(n)} &\sim \mathcal{G}(T; \alpha_t^{(n)}, \beta_t^{(n)})
\end{aligned}$$

9:   Sample Excitations

$$\begin{aligned}
\alpha_v^{(n)} &= A_v + \Sigma_v^{(n)} & \beta_v^{(n)} &= 1./ \left( A_v ./ B_v + (\mathbf{1}_W^\top T^{(n-1)})^\top \mathbf{1}_K^\top \right) \\
V^{(n)} &\sim \mathcal{G}(V; \alpha_v^{(n)}, \beta_v^{(n)})
\end{aligned}$$

10: **end for**

---

## 4 Simulations

Our goal is to illustrate our approach in a model selection context. We first illustrate that the variational approximation to the marginal likelihood is close to the one obtained from the Gibbs sampler via Chib’s method. Then, we compare the quality of solutions we obtain via Variational NMF and compare them to the original NMF on a prediction task. Finally, we focus on reconstruction quality in the overcomplete case where the standard NMF is subject to overfitting.

**Model Order Determination:** To test our approach, we generate synthetic data from the hierarchical model in (12) with  $W = 16, K = 10$  and the number of sources  $I_{\text{true}} = 5$ . The inference task is to find the correct number of sources, given  $X$ . The hyperparameters of the true model are set to  $a_{\nu,i}^t = a^t = 10, b_{\nu,i}^t = b^t = 1, a_{i,\tau}^v = a^v = 1, b_{i,\tau}^v = b^v = 100$ . In the first experiment the hyperparameters are assumed to be known and in the second are jointly estimated from data, using hyperparameter adaptation. We evaluate the marginal likelihood for models with the number of templates  $I = 1 \dots 10$ , with the Gibbs sampler using Chib’s method and variational lower bound  $\mathcal{B}$  via variational Bayes. We run the Gibbs sampler for MAXITER = 10000 steps following a burn-in period of 5000 steps; then we clamp the sources  $S$  and continue the simulation for a further 10000 steps to estimate quantities required by Chib’s method. We run the variational algorithm until convergence of the bound or 10000 iterations, whichever occurs first. In Figure 3-top, we show a comparison of the variational estimate with the average of 5 independent runs obtained via Chib’s method. We observe, that both methods give consistent results. In Figure 4, we show the lower bound as a function of model order  $I$ , where for each  $I$ , the bound is optimised independently by jointly optimising hyperparameters  $a_t, b_t, a_v$  and  $b_v$  using the equations derived in the Appendix. We observe, that the correct model order can be inferred even when the hyperparameters are unknown a-priori. This is potentially useful for estimation of model order from real data.

As real data, we use a version of the Olivetti face image database ( $K = 400$  images of  $64 \times 64$  pixels available at <http://www.cs.toronto.edu/~roweis/data/olivettifaces.mat>). We further downsampled to  $16 \times 16$  or  $32 \times 32$  pixels, hence our data matrix  $X$  is  $16^2 \times 400$  or  $32^2 \times 400$ . We use a model with tied hyperparameters as  $a_{\nu,i}^t = a^t, b_{\nu,i}^t = b^t, a_{i,\tau}^v = a^v$  and  $b_{i,\tau}^v = b^v$ , where all hyperparameters are jointly estimated. In Figure 4, bottom, we show results of model order determination for this dataset with joint hyperparameter adaptation. Here, we run the variational algorithm for each model order  $I = 1 \dots 100$  independently and evaluate the lower bound after optimising the

hyperparameters. The Gibbs sampler is not found practical and is omitted here. The lower bound behaves as is expected from marginal likelihood, reflecting the tradeoff between too many and too few templates. Higher resolution implies more templates, consistent with our intuition that detail requires more templates for accurate representation.

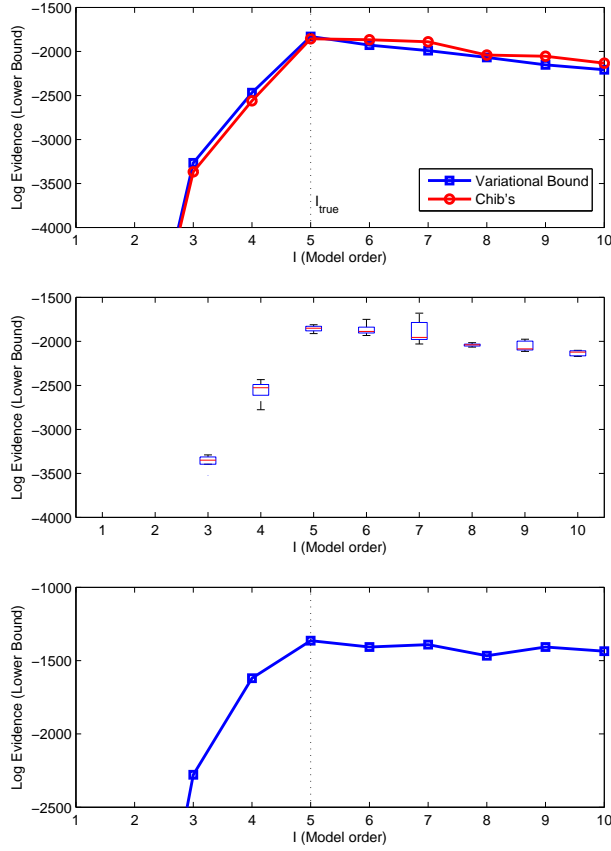


Figure 3: Model selection results. (top) Comparison of model selection by variational bound (squares) and marginal likelihood estimated by Chib's (circles) method. The hyperparameters are assumed to be known. (Middle) Box-plot of marginal likelihood estimates by Chib's method using 5000, 10000 and 10000 iterations for burn-in, free and clamped sampling. The boxes show the lower quartile, median, and upper quartile values. (Bottom) Model selection by variational bound when hyperparameters are unknown and jointly estimated.

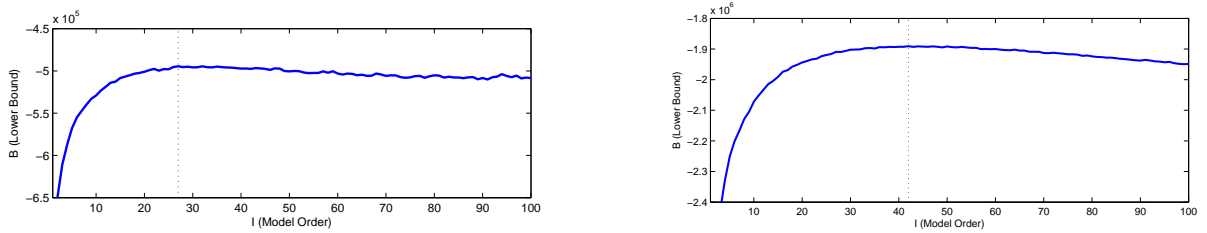


Figure 4: Model selection using variational bound with adapted hyperparameters on face data  $16 \times 16$  with  $I^* = 27$  (left) and  $32 \times 32$  with  $I^* = 42$  (right).

We also investigate the nature of the representations (see Figure 5). Here, for each independent run, we fix the values of shape parameters to  $(a^t, a^v) = [(10, 10), (0.1, 0.1), (10, 0.2), (10, 0.5)]$  and only estimate  $b^t$  and  $b^v$ . This corresponds to enforcing sparse or non-sparse  $t$  and  $v$ . Each column shows  $I = 36$  templates estimated from the dataset conditioned on hyperparameters. The middle image is the same template image above weighted with the excitations corresponding to the reconstruction (the expected value of the predictive distribution) below. Here, we clearly see the effect of the hyperparameters. In the first condition  $(a^t, a^v) = (10, 10)$ , the prior does not enforce sparsity to the templates and excitations. Hence, for the representation of a given image, there are many active templates. In the second condition, we try to force both matrices to be sparse with  $(a^t, a^v) = (0.1, 0.1)$ . Here, the result is not satisfactory as isolated components of the templates are zeroed, giving a representation that looks like one contaminated by “salt-and-pepper” noise. The third condition  $((a^t, a^v) = (10, 0.2))$  forces only the excitations to be sparse. Here, we observe that the templates correspond to

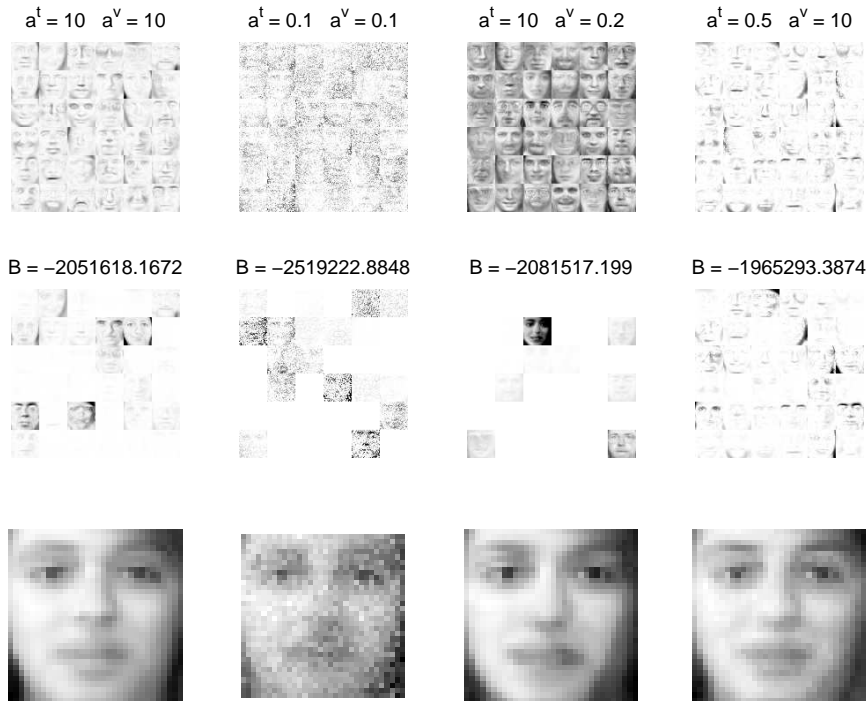


Figure 5: Templates, excitations for a particular example and the reconstructions obtained for different hyperparameter settings.  $B$  is the lower bound for the whole dataset.

some average face images. Qualitatively, each image is reconstructed using a superposition of a few of these templates. In the final representation, we enforce sparsity in the templates but not in the excitations. Here, our estimate finds templates that correspond to parts of individual face images (eyebrows, lips etc.). This solution, intuitively corresponding to a parsimonious representation, also is the best in terms of the marginal likelihood. With proper initialisation, our variational procedure is able to find such solutions.

**Prediction:** We now compare variational Bayesian NMF with the maximum likelihood NMF on a missing data prediction task. To illustrate the self regularisation effect, we set up an experiment in which we select a subset of the face data consisting of 50 images. From half of the images, we remove a the same patch (Fig. 6) and try to predict missing pixels. This is a rather small dataset for this task, as we have only 10 images for each of the 5 different persons, and half of these images have missing data at the same spot. We measure the quality of the prediction in terms of signal-to-noise ratio (SNR). The missing values are reconstructed using the mean of the predictive distribution  $\mathbf{X}_{\text{pred}} \equiv \langle X \rangle_{\mathcal{P}_{\mathcal{O}}(X; T^* V^*)} = T^* V^*$  where  $T^*$  and  $V^*$  are point estimates of the template and excitation matrix. We compare our variational algorithm with the classical NMF. For each algorithm, we test two different versions. The variational algorithms differ in how we estimate  $T^*$  and  $V^*$ . In the first variational algorithm, we use a crude estimate of  $T^*$  and  $V^*$  as the mean of the approximating  $q$  distribution. In the second condition, after convergence of hyperparameters via VB, we reinitialise  $T$  and  $V$  randomly and switch to an ICM algorithm (see Eq.21). This strategy finds a local mode  $(T^*, V^*)$  of the exact posterior distribution. In NMF, we test two initialisation strategies: in the first condition, we initialise the templates randomly. In the second we set them equal to the images in the dataset with random perturbations.

In Fig. 6, we show the reconstruction results for a typical run, for a model with  $I = 100$  templates. Note that this an overcomplete model, with twice as many templates as there are images. To characterise the nature of the estimated template and excitation matrices, we use the sparseness criteria [11] of an  $m \times n$  matrix  $X$ , defined as  $\text{Sparseness}(X) = (\sqrt{mn} - (\sum_{i,j} |X_{i,j}|) / (\sum_{i,j} X_{i,j}^2)^{1/2}) / (\sqrt{mn} - 1)$ . This measure is 1 when the matrix  $X$  has only a single non-zero entry and 0 when all entries are equal. We see that the variational algorithms are superior in this case in terms of SNR as well as the visual quality of the reconstruction. This is perhaps not surprising, since with maximum likelihood estimation; if the model order is not carefully chosen, generalisation performance is poor: the “noise” in the observed data is fitted but the prediction quality drops on new data. An interesting observation is that highly sparse solutions (either in templates or excitations) do not give the best result, the solution that balances both seems to be the best in this setting. This example illustrates that sparseness in itself may not be necessarily a good criteria to optimise; for model selection, the marginal likelihood should be used as the natural quantity.

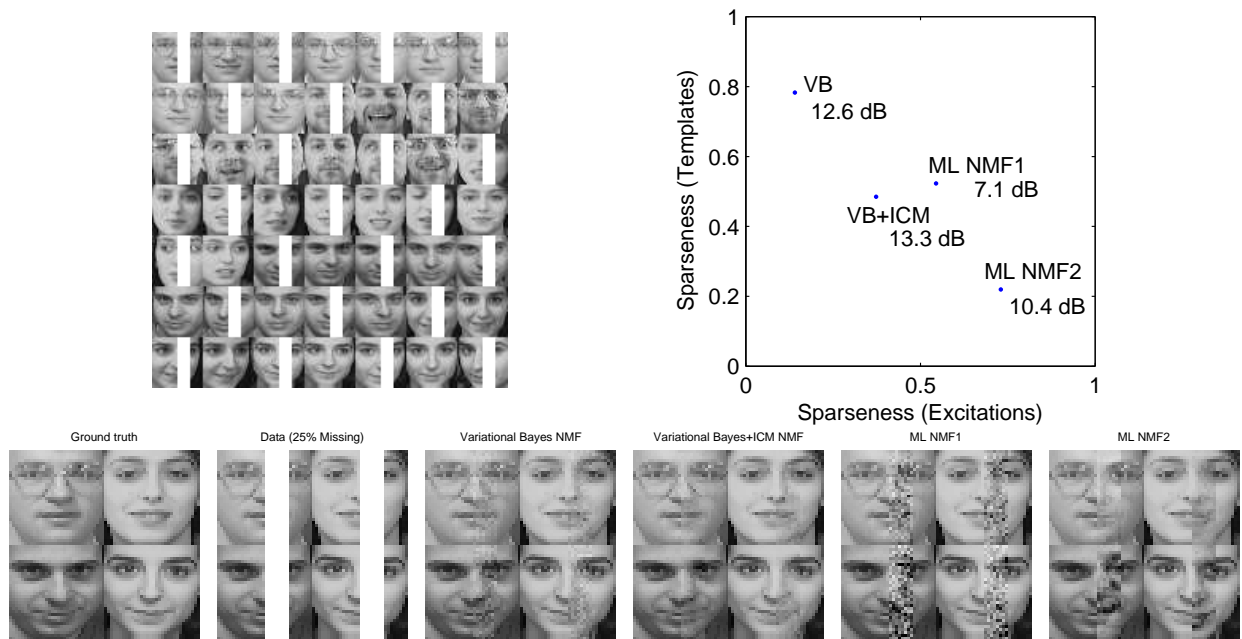


Figure 6: Results of a typical run. (Top) Example images from the data set. Left to right. The ground truth, data with missing pixels. The reconstructions of VB, VB+ICM, and ML NMF with two initialisation strategies (1=random, 2=to image). (Bottom) Comparison of the reconstruction accuracy of different methods in terms of SNR (in dB), organised according to the sparseness of the solution.

#### 4.1 Discussion and Conclusions

In this paper, we have investigated NMF from a statistical perspective. We have shown that KL minimisation formulation the original algorithm can be derived from a probabilistic model where the observations are superposition of  $I$  independent Poisson distributed latent sources. Here, the template and excitation matrices turn out to be latent intensity parameters. The interpretation of NMF as a maximum likelihood method in a Poisson model is mentioned in the original NMF paper [16] and discussed in more detail by [13, 22]. [13] focuses on the optimisation and gives an equivalent description using auxiliary function maximisation. In contrast, [22] illustrates that the auxiliary variables can be viewed as model variables (the sources  $s$ ) that are analytically integrated out. The novel observation in the current article is the exact characterisation of the approximating distribution  $q(S)$  or full conditionals  $p(S|T, V, X)$  as a product of multinomial distributions, leading to a richer approximation distribution than a naive mean field or single site Gibbs (which would freeze due to deterministic  $p(X|S)$ ). This conjugate form leads to significant simplifications in full Bayesian integration. Apart from the conditionally Gaussian case, NMF with KL objective seems to be unique in this respect. For several other distance metrics  $D(\cdot||\cdot)$ , we find full Bayesian inference not as practical, as  $p(S|T, V, X)$  is not standard.

We have also shown that the standard NMF algorithm with multiplicative update rules is in fact an EM algorithm with data augmentation. Extending upon this observation, we have developed an hierarchical model with conjugate Gamma priors. We have developed a variational Bayes algorithm and a Gibbs sampler for inference in this hierarchical model. We have also developed methods for estimating the marginal likelihood for model selection.

Our simulations suggest that the variational bound seems to be a reasonable approximation to the marginal likelihood and can guide model selection for NMF. The computational requirements are comparable to the ML NMF. A potentially time consuming step in the implementation of the variational algorithm is the evaluation of the  $\Psi$  function, but this step can also be replaced by a simple piecewise polynomial approximation since  $\exp(\Psi(x)) \approx x - 0.5$  for  $x > 5$ .

We first compare the variational inference with a Gibbs sampler. In our simulations, we observe that both algorithms give qualitatively very similar results, both for inference of templates and excitations as well as model order selection. We find the variational approach somewhat more practical as it can be expressed as simple matrix operations, where both the fixed point equations as well as the bound can be compactly and efficiently implemented using matrix computation software. In contrast, our Gibbs sampler is computationally more demanding and the calculation of marginal likelihood is somewhat more tricky. With our implementation of both algorithms the variational method is faster by a factor of around 13. Reference implementations of both algorithms in matlab are available from the following url: <http://www-sigproc.eng.cam.ac.uk/~atc27/bnmf/>.

In terms of computational requirements, the variational procedure has several advantages. First, we circumvent sampling from multinomial variables, which is the main computational bottleneck with the Gibbs sampler. Whilst efficient algorithms are developed for multinomial sampling [7], the procedure is time consuming when the number of latent sources  $I$  is large. In contrast, the variational method estimates the expected sufficient statistics directly by elementary

matrix operations. Another advantage is hyperparameter estimation. In principle, it is possible to maximise the marginal likelihood via a Monte Carlo EM procedure [21, 19], yet this potentially requires many more iterations. In contrast, the evaluation of the derivatives of the lower bound is straightforward and can be implemented without much additional computational cost.

The efficiency of the Gibbs sampler could be improved by working out the distribution of the sufficient statistics of sources directly (namely quantities  $\sum_{\tau} s_{\nu,i,\tau}$  or  $\sum_{\nu} s_{\nu,i,\tau}$ ) to circumvent multinomial sampling. Unfortunately, for the sum of binomial random variables with different cell probability parameters, the sum does not have a simple form but various approximations are possible [4].

From a modelling perspective, our hierarchical model provides some attractive properties. It is easy to incorporate prior knowledge about individual latent sources via hyper parameters and one can easily capture variability in the templates and excitations that is potentially useful for developing robust techniques. The prior structure here is qualitatively similar to an entropic prior [3, 20] and we find qualitatively similar representations to the ones found by NMF reported earlier by [16, 11]. However, none of the above mentioned methods provide an estimate of the marginal likelihood, which is useful for model selection. Our generative model formulation can be extended in various ways to suit the specific needs of particular applications. For example, one can enforce more structured prior models such as chains or fields [22]. As a second possibility, the Poisson observation model can be replaced with other models such as clipped Gaussian, Gamma or Gaussians which lead to alternative source separation algorithms. For example, the case of Gaussian sources where the excitations and templates correspond to the variances is discussed in [5].

Our main contribution here is the development of a principled and practical way to estimate both the optimal sparsity criteria and model order, in terms of marginal likelihood. By maximising the bound on marginal likelihood, we have a method where all the hyperparameters can be estimated from data, and the appropriate sparseness criteria is found automatically. We believe that our approach provides a practical improvement to the highly popular NMF algorithm without incurring much additional computational cost.

## Acknowledgements

We like to thank to Nick Whiteley, Tuomas Virtanen and Paul Peeling for fruitful discussion and for their comments on earlier drafts of this paper.

## A Standard Distributions in Exponential Form, their sufficient statistics and entropies

- Gamma

$$\begin{aligned} \mathcal{G}(\lambda; a, b) &\equiv \exp\left(+ (a-1) \log \lambda - \frac{1}{b} \lambda - \log \Gamma(a) - a \log b\right) \\ \langle \lambda \rangle_{\mathcal{G}} &= ab \quad \langle \log \lambda \rangle_{\mathcal{G}} = \Psi(a) + \log(b) \\ H_{\mathcal{G}}[\lambda] &\equiv -\langle \log \mathcal{G} \rangle_{\mathcal{G}} = -(a-1)\Psi(a) + \log b + a + \log \Gamma(a) \end{aligned}$$

Here,  $\Psi$  denotes the digamma function defined as  $\Psi(a) \equiv d \log \Gamma(a) / da$ .

- Poisson

$$\begin{aligned} \mathcal{PO}(s; \lambda) &= \exp(s \log \lambda - \lambda - \log \Gamma(s+1)) \\ \langle s \rangle_{\mathcal{PO}} &= \lambda \end{aligned}$$

- Multinomial

$$\begin{aligned} \mathcal{M}(\mathbf{s}; x, \mathbf{p}) &= \delta\left(x - \sum_i s_i\right) \exp\left(\log \Gamma(x+1) + \sum_{i=1}^I (s_i \log p_i - \log \Gamma(s_i+1))\right) \\ \langle s_i \rangle_{\mathcal{M}} &= x p_i \end{aligned}$$

Here,  $\mathbf{s} = \{s_1, s_2, \dots, s_I\}$  and  $\mathbf{p} = \{p_1, p_2, \dots, p_I\}$  and  $p_1 + p_2 + \dots + p_I = 1$ . Here,  $p_i, i = 1 \dots I$  are the cell probabilities and  $x$  is the index parameter where  $s_1 + s_2 + \dots + s_I = x$ . The entropy is given as:

$$\begin{aligned} H_{\mathcal{M}}[s_{\nu,1:I,\tau}] &= -\log \Gamma(x_{\nu,\tau} + 1) - \sum_{i=1}^I \langle s_{\nu,i,\tau} \rangle \log p_{\nu,i,\tau} \\ &\quad + \sum_{i=1}^I \langle \log \Gamma(s_{\nu,i,\tau} + 1) \rangle - \left\langle \log \delta\left(x_{\nu,\tau} - \sum_i s_{\nu,i,\tau}\right) \right\rangle \end{aligned}$$

A closed form expression for the entropy is not known due to  $\langle \log \Gamma(s+1) \rangle$  terms, but asymptotic expansions exist [8, 12]. Computationally efficient sampling from a multinomial distribution is not trivial, see [7] for a comparison of various methods and detailed discussion of tradeoffs.

## A.1 Summary of the Generative Model

Indices

$i = 1 \dots I$  Source index

$\nu = 1 \dots W$  Row (frequency bin) index

$\tau = 1 \dots K$  Column (time frame) index

$t_{\nu,i}$  Template variable at  $\nu$ 'th row of the  $i$ 'th source

$$t_{\nu,i} \sim \mathcal{G}(t_{\nu,i}; a_{\nu,i}^t, b_{\nu,i}^t/a_{\nu,i}^t)$$

$v_{i,\tau}$  Excitation variable of the  $i$ 'th source at  $\tau$ 'th column

$$v_{i,\tau} \sim \mathcal{G}(v_{i,\tau}; a_{i,\tau}^v, b_{i,\tau}^v/a_{i,\tau}^v)$$

$s_{\nu,i,\tau}$  Source variable of  $i$ 'th source at  $\nu$ 'th row (frequency bin) and  $\tau$ 'th column (time frame)

$$s_{\nu,i,\tau} \sim \mathcal{PO}(s_{\nu,i,\tau}; t_{\nu,i}v_{i,\tau})$$

$x_{\nu,\tau}$  Observation at  $\nu$ 'th row (frequency bin) and  $\tau$ 'th column (time frame)

$$x_{\nu,\tau} \sim \sum_i s_{\nu,i,\tau}$$

## A.2 Expression of the full joint distribution

Here,  $\phi \equiv p(X, S, T, V|\Theta) = p(X|S)p(S|T, V)p(V|\Theta^v)p(V|\Theta^t)$

$$\begin{aligned} \log \phi &= \sum_{\nu} \sum_i \sum_{\tau} (-t_{\nu,i}v_{i,\tau} + s_{\nu,i,\tau} \log(t_{\nu,i}v_{i,\tau}) - \log \Gamma(s_{\nu,i,\tau} + 1)) \\ &+ \sum_{\nu} \sum_{\tau} \log \delta(x_{\nu,\tau} - \sum_i s_{\nu,i,\tau}) \\ &+ \sum_{\nu} \sum_i (a_{\nu,i}^t - 1) \log t_{\nu,i} - \frac{a_{\nu,i}^t}{b_{\nu,i}^t} t_{\nu,i} - \log \Gamma(a_{\nu,i}^t) - a_{\nu,i}^t \log(b_{\nu,i}^t/a_{\nu,i}^t) \\ &+ \sum_{\tau} \sum_i (a_{i,\tau}^v - 1) \log v_{i,\tau} - \frac{a_{i,\tau}^v}{b_{i,\tau}^v} v_{i,\tau} - \log \Gamma(a_{i,\tau}^v) - a_{i,\tau}^v \log(b_{i,\tau}^v/a_{i,\tau}^v) \end{aligned}$$

## A.3 The Variational Bound

The variational bound in (13) can be written as

$$\mathcal{L}_X(\Theta) \equiv \log p(X|\Theta) \geq \langle \log \phi \rangle_q + H[q] = \mathcal{B}_{VB}$$

where the energy term is given by the expectation of the expression in section A.2 and  $H[q]$  denotes the entropy of the variational approximation distribution  $q$  where the individual entropies are defined in section A.

$$H[q] = -\langle \log q \rangle = \sum_{\nu} \sum_{\tau} H_{\mathcal{M}}[s_{\nu,1:I,\tau}] + \sum_{\nu} \sum_i H_{\mathcal{G}}[t_{\nu,i}] + \sum_i \sum_{\tau} H_{\mathcal{G}}[v_{i,\tau}]$$

One potential problem is that this expression requires the entropy of a multinomial distribution for which there is no known simple expression. This is due to terms of form  $\langle \log \Gamma(s+1) \rangle$  where only asymptotic expansions are known.

Fortunately, the difficult terms in the energy term can be canceled by the corresponding terms in the entropy term and one obtains the following expression that only depends on known sufficient statistics.

$$\begin{aligned}
\mathcal{B} = & -\sum_{\nu} \sum_{\tau} \sum_i \langle t_{\nu,i} \rangle \langle v_{i,\tau} \rangle \\
& + \sum_{\nu} \sum_i \langle \log t_{\nu,i} \rangle \left( a_{\nu,i}^t - 1 + \sum_{\tau} \langle s_{\nu,i,\tau} \rangle \right) + \sum_{\tau} \sum_i \langle \log v_{i,\tau} \rangle \left( a_{i,\tau}^v - 1 + \sum_{\nu} \langle s_{\nu,i,\tau} \rangle \right) \\
& + \sum_{\nu} \sum_i -\frac{a_{\nu,i}^t}{b_{\nu,i}} \langle t_{\nu,i} \rangle - \log \Gamma(a_{\nu,i}^t) - a_{\nu,i}^t \log(b_{\nu,i}/a_{\nu,i}^t) \\
& + \sum_{\tau} \sum_i -\frac{a_{i,\tau}^v}{b_{i,\tau}^v} \langle v_{i,\tau} \rangle - \log \Gamma(a_{i,\tau}^v) - a_{i,\tau}^v \log(b_{i,\tau}^v/a_{i,\tau}^v) \\
& + \sum_{\nu} \sum_{\tau} \left( -\log \Gamma(x_{\nu,\tau} + 1) - \sum_{i=1}^I \langle s_{\nu,i,\tau} \rangle \log p_{\nu,i,\tau} \right) \\
& + \sum_{\nu} \sum_i (-\langle \alpha_{\nu,i}^t \rangle - 1) \Psi(\langle \alpha_{\nu,i}^t \rangle) + \log \beta_{\nu,i}^t + \langle \alpha_{\nu,i}^t \rangle + \log \Gamma(\langle \alpha_{\nu,i}^t \rangle) \\
& + \sum_i \sum_{\tau} (-\langle \alpha_{i,\tau}^v \rangle - 1) \Psi(\langle \alpha_{i,\tau}^v \rangle) + \log \beta_{i,\tau}^v + \langle \alpha_{i,\tau}^v \rangle + \log \Gamma(\langle \alpha_{i,\tau}^v \rangle)
\end{aligned}$$

After some careful manipulations, the following expression is obtained where  $\log L$  denotes here element wise logarithm of matrix  $L$ .

$$\begin{aligned}
\mathcal{B} = & \sum_{\nu} \sum_{\tau} (-E_t E_v - \log \Gamma(X + 1)) \\
& + \sum_{\nu} \sum_{\tau} -X .* (((L_t .* \log(L_t)) L_v + L_t (L_v .* \log(L_v))) ./ (L_t L_v) - \log(L_t L_v)) \\
& + \sum_{\nu} \sum_i -(A_t ./ B_t) .* E_t - \log \Gamma(A_t) + A_t .* \log(A_t ./ B_t) \\
& + \sum_{\nu} \sum_i \alpha_t .* (\log \beta_t + 1) + \log \Gamma(\alpha_t) \\
& + \sum_i \sum_{\tau} -(A_v ./ B_v) .* E_v - \log \Gamma(A_v) + A_v .* \log(A_v ./ B_v) \\
& + \sum_i \sum_{\tau} \alpha_v .* (\log \beta_v + 1) + \log \Gamma(\alpha_v)
\end{aligned} \tag{28}$$

#### A.4 Handling Missing Data and MAP estimation

When there is missing data, i.e., when some of the  $x_{\nu,\tau}$  are not observed, computation is still straightforward in our framework and can be accomplished by a simple modification to the original algorithm. We first define a *mask* matrix  $\mathbf{M} = \{m_{\nu,\tau}\}$ , same size as  $X$  where

$$m_{\nu,\tau} = \begin{cases} 0 & x_{\nu,\tau} \text{ is missing} \\ 1 & \text{otherwise} \end{cases} \tag{29}$$

Using the mask variables, the observation model with missing data can be written as

$$p(X|S)p(S|T, V) = \prod_{\nu,\tau} (p(x_{\nu,\tau} | s_{\nu,1:I,\tau}) p(s_{\nu,1:I,\tau} | t_{\nu,1:I}, v_{1:I,\tau}))^{m_{\nu,\tau}}$$

The prior is not affected. Hence, we merely replace the first two lines of the expression for the full joint distribution (given in the appendix A.2) as

$$\begin{aligned}
\log \phi = & \sum_{\nu} \sum_{\tau} m_{\nu,\tau} \sum_i (-t_{\nu,i} v_{i,\tau} + s_{\nu,i,\tau} \log(t_{\nu,i} v_{i,\tau}) - \log \Gamma(s_{\nu,i,\tau} + 1)) \\
& + \sum_{\nu} \sum_{\tau} m_{\nu,\tau} \log \delta(x_{\nu,\tau} - \sum_i s_{\nu,i,\tau}) + \dots
\end{aligned}$$

Consequently, it is easy to see that

$$\begin{aligned}
q(v_{i,\tau}) & \propto \exp \left( (a_{i,\tau}^v + \sum_{\nu} m_{\nu,\tau} \langle s_{\nu,i,\tau} \rangle - 1) \log v_{i,\tau} - \left( \frac{a_{i,\tau}^v}{b_{i,\tau}^v} + \sum_{\nu} m_{\nu,\tau} \langle t_{\nu,i} \rangle \right) v_{i,\tau} \right) \\
& \propto \mathcal{G}(v_{i,\tau}; \alpha_{i,\tau}^v, \beta_{i,\tau}^v)
\end{aligned}$$



$$\alpha_{i,\tau}^v \equiv a_{i,\tau}^v + \sum_{\nu} m_{\nu,\tau} \langle s_{\nu,i,\tau} \rangle \quad \beta_{i,\tau}^v \equiv \left( \frac{a_{i,\tau}^v}{b_{i,\tau}^v} + \sum_{\nu} m_{\nu,\tau} \langle t_{\nu,i} \rangle \right)^{-1}$$

By a derivation analogous to one detailed in section 3.3, we see that the excitation update equations in Algorithm 1, line 4, can be written using matrix notation as

$$\begin{aligned} \Sigma_v &= L_v .* (L_t^\top ((\mathbf{M} .* X) ./ (L_t L_v))) \\ \boldsymbol{\alpha}_v &= A_v + \Sigma_v \quad \boldsymbol{\beta}_v = 1 ./ (A_v ./ B_v + E_t^\top \mathbf{M}) \end{aligned}$$

The update rules for the templates are similar. Note that when there is no missing data, we have  $\mathbf{M} = \mathbf{1}_W \mathbf{1}_K^\top$  which gives the original algorithm. The bound in (28) can also be easily modified for handling missing data. We merely replace  $\mathbf{X} \leftarrow \mathbf{M} .* \mathbf{X}$  and the first term  $E_t E_v \leftarrow \mathbf{M} .* E_t E_v$ .

We conclude this subsection by noting that the standard NMF update equations, given in Equations (10) and (11) can be also rewritten to handle missing data:

$$\begin{aligned} V^{(n+1)} &= V^{(n)} .* (T^\top ((\mathbf{M} .* X) ./ (TV))) ./ (T^\top \mathbf{M}) \\ T^{(n+1)} &= T^{(n)} .* (((\mathbf{M} .* X) ./ (TV)) V^\top) ./ (\mathbf{M} V^\top) \end{aligned}$$

Here, the denominator has to be nonzero. Similarly, an iterative conditional modes (ICM) algorithm can be derived to compute the maximum a-posteriori (MAP) solution

$$V^{(n+1)} = (A_v + V^{(n)} .* (T^\top ((\mathbf{M} .* X) ./ (TV)))) ./ (A_v ./ B_v + T^\top \mathbf{M}) \quad (30)$$

$$T^{(n+1)} = (A_t + T^{(n)} .* (((\mathbf{M} .* X) ./ (TV)) V^\top)) ./ (A_t ./ B_t + \mathbf{M} V^\top) \quad (31)$$

Note that when the shape parameters go to zero, i.e.  $A_t, A_v \rightarrow \mathbf{0}$ , we obtain the maximum likelihood NMF algorithm.

## A.5 Hyperparameter optimisation

The hyperparameters  $\Theta = (\Theta^t, \Theta^v)$  can be estimated by maximising the bound in 13. Below, we will derive the results for the excitations; the results for templates are similar. The solution for shape parameters involves finding the zero of a function  $f(a) - c$  where

$$\begin{aligned} f(a) &= \log a - \Psi(a) + 1 \\ a^* &= f^{-1}(c) \end{aligned}$$

The solution can be found by Newton's method by iteration of the following fixed point equation

$$a^{(n+1)} = a^{(n)} - \frac{f(a^{(n)}) - c}{f'(a^{(n)})} = a^{(n)} - \frac{\log(a^{(n)}) - \Psi(a^{(n)}) + 1 - c}{1/a^{(n)} - \Psi'(a^{(n)})} = a^{(n)} - \Delta^{(n)}$$

It is well known that Newton iterations can diverge if started away from the root. Occasionally, we observe that  $a$  can become negative. If this is the case, we set  $\Delta^{(n)} \leftarrow \Delta^{(n)}/2$ , and try again. The digamma function  $\Psi$  function and its derivative  $\Psi'$  are available in numeric computation libraries (e.g. in Matlab as `psi(0, a)` and `psi(1, a)`, respectively).

The derivation of the hyperparameter update equations is straightforward:

$$\begin{aligned} \frac{\partial \mathcal{B}}{\partial a_{i,\tau}^v} &= \langle \log v_{i,\tau} \rangle - \frac{1}{b_{i,\tau}^v} \langle v_{i,\tau} \rangle - \Psi(a_{i,\tau}^v) - \log b_{i,\tau}^v + \log a_{i,\tau}^v + 1 = 0 \\ c_{i,\tau} &= \log a_{i,\tau}^v - \Psi(a_{i,\tau}^v) + 1 \\ c_{i,\tau} &\equiv \frac{\langle v_{i,\tau} \rangle}{b_{i,\tau}^v} - (\langle \log v_{i,\tau} \rangle - \log b_{i,\tau}^v) \end{aligned}$$

$$\begin{aligned} \frac{\partial \mathcal{B}}{\partial b_{i,\tau}^v} &= \frac{a_{i,\tau}^v}{(b_{i,\tau}^v)^2} \langle v_{i,\tau} \rangle - a_{i,\tau}^v \frac{1}{b_{i,\tau}^v} = 0 \\ b_{i,\tau}^v &= \langle v_{i,\tau} \rangle \end{aligned}$$

Tying parameters across  $\tau$  as  $a_i^v = a_{i,\tau}^v$  and  $b_i^v = b_{i,\tau}^v$  yields

$$\begin{aligned} \frac{\partial \mathcal{B}}{\partial a_i^v} &= \sum_{\tau} \langle \log v_{i,\tau} \rangle - \sum_{\tau} \frac{1}{b_{i,\tau}^v} \langle v_{i,\tau} \rangle - K \Psi(a_i^v) - \sum_{\tau} \log b_{i,\tau}^v + K = 0 \\ c_i &= \log a_i^v - \Psi(a_i^v) + 1 \\ c_i &= \frac{1}{K} \sum_{\tau} \left( \frac{\langle v_{i,\tau} \rangle}{b_{i,\tau}^v} - (\langle \log v_{i,\tau} \rangle - \log b_{i,\tau}^v) \right) \end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{B}}{\partial b_i^v} &= \sum_{\tau} \frac{a_{i,\tau}^v}{(b_i^v)^2} \langle v_{i,\tau} \rangle - \frac{1}{b_i^v} \sum_{\tau} a_{i,\tau}^v \\ b_i^v &= \frac{\sum_{\tau} a_{i,\tau}^v \langle v_{i,\tau} \rangle}{\sum_{\tau} a_{i,\tau}^v}\end{aligned}$$

Tying parameters across  $\tau$  and  $i$ ,  $a^v = a_{i,\tau}^v$  and  $b^v = b_{i,\tau}^v$  yields

$$\begin{aligned}c &= \log a^v - \Psi(a^v) + 1 \\ c &= \frac{1}{KI} \sum_{\tau} \sum_i \left( \frac{\langle v_{i,\tau} \rangle}{b_{i,\tau}^v} - (\langle \log v_{i,\tau} \rangle - \log b_{i,\tau}^v) \right)\end{aligned}$$

$$\begin{aligned}\frac{\partial \mathcal{B}}{\partial b^v} &= \sum_i \sum_{\tau} \frac{a_{i,\tau}^v}{(b^v)^2} \langle v_{i,\tau} \rangle - \frac{1}{b^v} \sum_i \sum_{\tau} a_{i,\tau}^v \\ b^v &= \frac{\sum_i \sum_{\tau} a_{i,\tau}^v \langle v_{i,\tau} \rangle}{\sum_i \sum_{\tau} a_{i,\tau}^v}\end{aligned}$$

The derivation of the template parameters is exactly analogous. We can express the update equations once again in compact matrix notation

$$\begin{aligned}Z &\leftarrow E_v^{(n)} ./ B_v^{(n)} - \log(L_v^{(n)} ./ B_v^{(n)}) \\ C &\leftarrow \begin{cases} Z & \text{Not tied} \\ (Z \mathbf{1}_K) / K & \text{Tie columns (over } \tau) \\ (\mathbf{1}_I^\top Z) / I & \text{Tie rows (over } i) \\ (\mathbf{1}_I^\top Z \mathbf{1}_K) / (KI) & \text{Tie all (over } \tau \text{ and } i) \end{cases} \\ A_v^{(n+1)} &\leftarrow \text{SolveByNewton}(A_v^{(n)}, C) \\ B_v^{(n+1)} &\leftarrow \begin{cases} E_v^{(n)} & \text{Not tied} \\ (((A_v^{(n)} .* E_v^{(n)}) \mathbf{1}_K) ./ (A_v^{(n)} \mathbf{1}_K)) \mathbf{1}_K^\top & \text{Tie columns (over } \tau) \\ \mathbf{1}_I ((\mathbf{1}_I^\top (A_v^{(n)} .* E_v^{(n)})) ./ (\mathbf{1}_I^\top A_v^{(n)})) & \text{Tie rows (over } i) \\ \mathbf{1}_I ((\mathbf{1}_I^\top (A_v^{(n)} .* E_v^{(n)}) \mathbf{1}_K) ./ (\mathbf{1}_I^\top A_v^{(n)} \mathbf{1}_K)) \mathbf{1}_K^\top & \text{Tie all (over } \tau \text{ and } i) \end{cases}\end{aligned}$$

Here, we assume  $\text{SolveByNewton}(A_0, C)$  is a matrix valued function that finds root  $C_{i,j} = f(A_{i,j})$  for each element of  $A$ , starting from the initial matrix  $A_0$ . If  $C$  is a scalar or vector, it is repeated over the missing index to implement parameter tying. For example, if  $C$  is a  $I \times 1$  vector and  $A_0$  is  $I \times K$ , we assume  $C_i = c_{i,\tau}$  for all  $\tau = 1 \dots K$  and the output is the same size as  $A_0$ . This is only a notational convenience, an actual implementation can be achieved more efficiently. Again, the implementation of the template parameters is exactly analogous; merely replace above the subscripts as  $v \leftarrow t$ ,  $(i, \tau) \leftarrow (\nu, i)$  and  $(I, K) \leftarrow (W, I)$ .

## References

- [1] P. Smaragdis W. Wang A. Cichocki, M. Mrup and R. Zdunek. Advances in nonnegative matrix and tensor factorization. *Computational Intelligence and Neuroscience*, 2008, 2008.
- [2] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [3] W. Buntine. Variational extensions to EM and multinomial PCA. In *ECML*, LNAI 2430, pages 23–34. Springer, 2002.
- [4] K. Butler and M. Stephens. The distribution of a sum of binomial random variables. Prepared for the Office of Naval Research 467, Stanford University, April 1993.
- [5] A. T. Cemgil, P. Peeling, O. Dikmen, and S. J. Godsill. Prior structures for time-frequency energy distributions. In *Proc. of IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, October 2007.
- [6] S. Chib. Marginal likelihood from the gibbs output. *JASA*, 90(432):1313–1321, Dec. 1995.
- [7] C. S. Davis. The computer generation of multinomial random variates. *Computational Statistics and Data Analysis*, 16:205–217, 1993.
- [8] P. Flajolet. Singularity analysis and asymptotics of bernoulli sums. *Theor. Comput. Sci.*, 275:371–387, 1999.
- [9] Z. Ghahramani and M. Beal. Propagation algorithms for variational Bayesian learning. In *Neural Information Processing Systems 13*, 2000.

- [10] W. R. Gilks, S. Richardson, and D. J Spiegelhalter, editors. *Markov Chain Monte Carlo in Practice*. CRC Press, London, 1996.
- [11] P. O. Hoyer. Non-negative matrix factorization with sparseness constraints. *J. Mach. Learn. Res.*, 5:1457–1469, 2004.
- [12] P. Jacquet and W. Szpankowski. Entropy computations via analytic depoissonization. *IEEE Transactions on Information theory*, 45(4), May 1999.
- [13] H. Kameoka. *Statistical Approach to Multipitch Analysis*. PhD thesis, University of Tokyo, 2007.
- [14] J. F. C. Kingman. *Poisson Processes*. Oxford Science Publications, 1993.
- [15] J. Le Roux, H. Kameoka, N. Ono, and S. Sagayama. On the relation between divergence-based minimization and maximum-likelihood estimation for the i-divergence. *Submitted to IEEE Trans. Signal Processing*, 2008.
- [16] D. D. Lee and H. S. Seung. Learning the parts of objects with nonnegative matrix factorization. *Nature*, 401:788–791, 1999.
- [17] D. D. Lee and H. S. Seung. Algorithms for non-negative matrix factorization. In *NIPS*, pages 556–562, 2000.
- [18] J. S. Liu. *Monte Carlo strategies in scientific computing*. Springer, 2004.
- [19] F. A. Quintana, J. S. Liu, and G. E. del Pino. Monte Carlo EM with importance reweighting and its applications in random effects models. *Computational Statistics and Data Analysis*, 29:429–444, 1999.
- [20] M. V. Shashanka, B. Raj, and P. Smaragdis. Sparse overcomplete latent variable decomposition of counts data. *Neural Information Processing Systems (NIPS)*, December 2007.
- [21] M. A. Tanner. *Tools for Statistical Inference: Methods for the Exploration of Posterior Distributions and Likelihood Functions*. Springer, New York, 3rd edition, 1996.
- [22] T. O. Virtanen, A. T. Cemgil, and S. J. Godsill. Bayesian extensions to nonnegative matrix factorisation for audio signal modelling. In *Proc. of IEEE ICASSP 08*, Las Vegas, 2008. IEEE.
- [23] Tuomas Virtanen. *Sound Source Separation in Monaural Music Signals*. PhD thesis, Tampere University of Technology, November 2006.