

Supervised Dictionary Learning by a Variational Bayesian Group Sparse Nonnegative Matrix Factorization

Ivan Ivek

Abstract— Nonnegative matrix factorization (NMF) with group sparsity constraints is formulated as a probabilistic graphical model and, assuming some observed data have been generated by the model, a feasible variational Bayesian algorithm is derived for learning model parameters. When used in a supervised learning scenario, NMF is most often utilized as an unsupervised feature extractor followed by classification in the obtained feature subspace. Having mapped the class labels to a more general concept of groups which underlie sparsity of the coefficients, what the proposed group sparse NMF model allows is incorporating class label information to find low dimensional label-driven dictionaries which not only aim to represent the data faithfully, but are also suitable for class discrimination. Experiments performed in face recognition and facial expression recognition domains point to advantages of classification in such label-driven feature subspaces over classification in feature subspaces obtained in an unsupervised manner.

Index Terms— Face and gesture recognition, Markov random fields, Pattern analysis

1 INTRODUCTION

SINCE the appearance of the seminal paper [1], NMF has become a popular data decomposition technique due to successful applications in a still growing number of fields where data are nonnegative, such as pixel intensities in computer vision, amplitude spectra in audio signal analysis and EEG signal analysis, term counts in document clustering problems, and item ratings in collaborative filtering.

NMF aims at decompositions $\mathbf{X} \approx \mathbf{T}\mathbf{V}$, where \mathbf{X} , \mathbf{T} and \mathbf{V} are all nonnegative matrices. Throughout this paper \mathbf{X} will be regarded as a collection of data samples organized columnwise, \mathbf{T} as a dictionary of features organized columnwise, and \mathbf{V} as matrix of coefficients when \mathbf{X} is projected onto the dictionary \mathbf{T} . Under assumptions of linearity and nonnegativity, when underlying dimensionality is lower than dimensionality of the original space of the data \mathbf{X} , dimensionality reduction of the data can effectively be achieved this way.

Although the decomposition is nonunique in general, NMF is able to produce strictly additive decompositions perceived as part-based by adding additional bias in the model [1], [2]. To this end, different sparsity promoting regularizers have been proposed for divergence-based NMF [3]. Also, to include higher order data descriptions, many other variants have been developed, e.g. Local NMF [4] with locality constraints, Non-smooth NMF [5] with regularization for sparse and localized features, NMF with smoothness constraints [6], graph regularized NMF [7], [8] manifold regularized NMF [9]. More recently, alternative formulations of NMF in probabilistic framework have been developed [10], [11], allowing for

explicit modeling of richer structural constraints as graphical models [12], [13], [14], [15].

1.1 Context and Contributions

Closely related to work in this paper are NMF formulations and extensions with group sparsity constraints. Addressing EEG classification where problem is to classify different tasks being performed by different subjects, in [16] a divergence based NMF with mixed-norm regularization imposed on dictionary elements to get task-related (common) features that are as close as possible and, separately, features which reflect task-independent (individual) characteristics, required to be as far as possible, has been proposed. Another divergence based algorithm but with group sparsity penalizations on the coefficient matrix has been proposed in [17]. In [18] a generative model with the same aim, enforcing two groups of features with the previously mentioned properties, is trained using a variational Bayesian approach. In [19] a NMF variant with Itakura-Saito divergence with a direct group-sparsity enforcing penalization of coefficient matrix has been successfully applied to blind audio source separation. Very recently, in [20] a generative model trained with Markov chain Monte Carlo is proposed to separate features into groups of common bases and individual bases with Laplacian scale mixture distributions as priors for the two groups, and applied to blind music source separation.

Work presented in this paper had come out of the probabilistic NMF modeling track of research, with group sparsity constraints as exponential scale mixture distributions imposed on coefficient matrix directly rather than seeking group sparsity by constraining the two groups of features as in [16], [18], [20]. Comparatively, it may be

• Ivek, I. is with the Division of Electronics, Rudjer Boskovic Institute, Croatia. E-mail: ivan.ivek@irb.hr

noted that algorithms [16], [17], [19] are not of Bayesian type. Although it aims to impose group sparsity through common and individual features, most resemblance to the here presented model bears [20] because of Laplacian scale mixture as a prior for the two groups, but it uses Metropolis-Hastings algorithm for parameter estimation due to choices of distributions in their hierarchical model. Group sparse NMF presented in this paper however, is engineered in a way that it allows efficient learning, here performed using mean-field variational Bayesian methodology, which is deterministic and features explicit variational bound calculation based upon which model comparison and selection can be made [11].

Most commonly, NMF is utilized in two scenarios: unsupervised learning, to get decompositions suitable for clustering, and supervised learning, where NMF is used as an unsupervised feature extractor followed by classification in the obtained feature space. In both those cases NMF does not incorporate label information. Applications presented in this paper are focused on the latter case, but, instead of ignoring the data labeling, the proposed group sparse NMF model is put in a setup where it is label-driven in an attempt to bring out potential information in the labeling to find feature subspaces which aid classification. One remaining scenario is NMF in semi-supervised learning applications, which divergence based algorithms [21], [22], [23], [24] have been developed for, also possessing the ability to include and use label information.

1.2 Probabilistic Formulations of NMF

It has been shown in [10] that maximum-likelihood (ML) estimation of factor matrices \mathbf{T} and \mathbf{V} under different noise distributions are equivalent to NMF algorithms with different corresponding divergences. More precisely, ML estimation of probabilistic NMF under Gaussian, Poissonian, and Gamma noise correspond to minimization of Euclidean, Kullback-Leibler (KL) and Itakura-Saito divergences, respectively. In case of maximum-a-posteriori (MAP) estimation, exponential prior on a factor corresponds to sparsity promoting l1 regularization. Another connection between probabilistic modeling and NMF worth noting is that Probabilistic Latent Semantic Analysis, which is an expectation-maximization algorithm, is equivalent to KL-NMF algorithm using multiplicative updates [25]. Apart from ML, MAP and EM estimators, Bayesian methods (Monte Carlo, variational Bayes) have successfully been employed for efficient NMF parameter learning [11], [14], [15].

1.3 Variational Bayesian Learning

Consider the objective of minimizing dissimilarity between conditional distribution $P(H|D, \theta)$, where H denotes hidden variables in the model, D observed variables and θ model parameters, and its instrumental variational approximation $q(H)$, quantified by Kullback-Leibler divergence

$$\text{KL}(q||p) = \int_H q(H) \log \left(\frac{q(H)}{p(H|D, \theta)} \right) dH. \quad (1)$$

Equation (1) can be rewritten as

$$\begin{aligned} \text{KL}(q||p) &= \int_H q(H) \log \left(\frac{q(H)}{p(H|D, \theta)} \right) dH \\ &= \int_H q(H) \log \left(\frac{q(H)p(D|\theta)}{p(H, D|\theta)} \right) dH \\ &= \int_H q(H) \log \left(\frac{q(H)}{p(H, D|\theta)} \right) dH + \int_H q(H) p(D|\theta) dH \\ &= \int_H q(H) \log \left(\frac{q(H)}{p(H, D)} \right) dH + \log p(D|\theta) \\ &= -\mathcal{L}(q) + \log p(D|\theta). \end{aligned}$$

Because KL divergence is nonnegative, it turns out that $\mathcal{L}(q)$ is lower bound on the marginal loglikelihood of the observed data. On the other hand, because $\log p(D|\theta)$ is constant, the objective of minimizing KL divergence can be reformulated as maximization of $\mathcal{L}(q)$.

By expanding the expression for the variational bound

$$\begin{aligned} \mathcal{L}(q) &= \int_H q(H) \log p(H, D) dH - \int_H q(H) \log q(H) dH \\ &= \langle \log p(H, D|\theta) \rangle_{q(H)} + \mathcal{H}(q), \end{aligned} \quad (2)$$

where $\mathcal{H}(q)$ denotes entropy of the approximation $q(H)$, and supposing that the approximative variational distribution takes a factorized form $q(H) = \prod_{\alpha \in C} q_{\alpha}$, it can be shown that iterative local updates alternating over C of form

$$(q_{\alpha})^{(n+1)} \propto \exp(\langle \log p(H, D|\theta) \rangle_{q_{-\alpha}^{(n)}}), \quad (3)$$

improve lower bound on the marginal loglikelihood monotonically, with

$$q_{-\alpha} = \frac{\prod_{\alpha' \neq \alpha} q_{\alpha'}}{q_{\alpha}}.$$

Should the mode be conjugate-exponential, all the expectations in (3) necessarily assume analytical forms [26].

2 METHODOLOGY

2.1 Exponential Scale Mixture Distribution

Let random variable V be a product of reciprocal of some positive random variable λ and exponentially distributed random variable U with scale 1,

$$V = \lambda^{-1}U.$$

Supposing independency of λ and u , conditioned on λ , v is exponentially distributed with scale λ^{-1} ,

$$p(v|\lambda) = \text{Exponential}(v|\lambda^{-1}). \quad (4)$$

and its marginal distribution assumes form of a continuous mixture with mixing variable λ ,

$$p(v) = \int_0^{\infty} p(v|\lambda)p(\lambda) d\lambda.$$

Note that, should the distribution of the mixing variable be discrete, the above expression collapses to a discrete mixture of exponentials with specific shape parameters.

Because exponential distribution can be obtained by truncation only of Laplacian distribution, it follows that exponential scale mixture is a special case of Laplacian scale mixture distribution [20], [27].

2.2 Group Sparsity Model

Placing exponential scale mixture in setting of

probabilistic inference and learning, feasibility of related algorithms depends on the choice of prior $p(\lambda)$ and, as feasibility is prioritized, this choice of suitable priors gets narrowed down. As (3) suggests, conjugacy is desirable when computational complexity is taken into consideration. For examples of how conjugate-exponential models are used for variational Bayesian NMF, the interested reader is referred to [11], [14] and also [15], where a variational family more general than one which conjugate model would suggest to get analytical expressions appropriate for optimization.

In this paper prior $p(\lambda)$ is engineered as a hierarchical graphical model with a discrete mixture

$$p(\lambda|\lambda_c, z) = \delta\left(\lambda - \sum_{c=1}^C \delta(z-c)\lambda_c\right) \quad (5)$$

of C gamma distributed variables λ_c ,

$$p(\lambda_c|a_c^\lambda, b_c^\lambda) = \text{Gamma}(\lambda_c|a_c^\lambda, b_c^\lambda), c \in \{1, \dots, C\} \quad (6)$$

with z as a categorical mixture selector variable.

According to (4), (5) and (6), v is exponentially distributed with different scale parameters λ_c^{-1} for different selections of z .

To clarify on relationship to group sparsity, suppose first that multiple variables v_τ exist organized in groups $\{1, \dots, C\}$, their affiliation indicated by variables $z_\tau \in \{1, \dots, C\}$. Because reciprocals of the variables λ_c , representing mean values of the exponential distributions in the mixture and being interpreted as continuous indicators of how large the averaged outcomes in a group are, have inverse gamma probability density functions, the masses of these density functions can be tuned to be concentrated on small values. Such priors act as constraints in a way that only the minority of the groups are expected to have exponential distributions with significantly large mean values, which is a suitable way to describe group sparse processes consistently in a probabilistic setup.

2.3 Probabilistic NMF with Group Sparsity Prior

The proposed generative NMF model consist of:

- 1) mixing stage with dimensionality reduction effect under Poissonian noise

$$p(x_{v\tau}|s_{v:\tau}) = \delta(x_{v\tau} - \sum_i s_{vi})$$

$$p(s_{v\tau}|t_{vi}, v_{i\tau}) = \text{Poisson}(s_{v\tau}; t_{vi}v_{i\tau}),$$

- 2) gamma priors for the left matrix \mathbf{T}

$$p(t_{vi}|a_{vi}^t, b_{vi}^t) = \text{Gamma}(t_{vi}|a_{vi}^t, b_{vi}^t)$$

- 3) group sparsity structural constraints over the right matrix \mathbf{V}

$$p(v_{i\tau}|z_\tau, \lambda_{ic}) = \text{Exponential}(v_{i\tau} | (\sum_c \delta(z_\tau - c)\lambda_{ic})^{-1})$$

$$p(\lambda_{ic}|a_{ic}^\lambda, b_{ic}^\lambda) = \text{Gamma}(\lambda_{ic}|a_{ic}^\lambda, b_{ic}^\lambda).$$

Both the mixing stage and the gamma priors over \mathbf{T} are designed as in [11], while the structural constraints over \mathbf{V} are introduced as a novelty.

Joint distribution of the proposed model is

$$p(H, D|\theta) = p(\mathbf{X}, \mathbf{S}, \mathbf{T}, \mathbf{V}, \mathbf{A} | \mathbf{A}^t, \mathbf{B}^t, \mathbf{A}^\lambda, \mathbf{B}^\lambda, \vec{z})$$

$$= p(\mathbf{X}|\mathbf{S})p(\mathbf{S}|\mathbf{T}, \mathbf{V})p(\mathbf{T}|\mathbf{A}^t, \mathbf{B}^t)$$

$$* p(\mathbf{V}|\mathbf{A}, \vec{z})p(\mathbf{A}|\mathbf{A}^\lambda, \mathbf{B}^\lambda),$$

with

$$p(\mathbf{X}|\mathbf{S}) = \prod_{v,\tau} p(x_{v\tau}|s_{v:\tau})$$

$$p(\mathbf{S}|\mathbf{T}, \mathbf{V}) = \prod_{v,\tau} p(s_{v:\tau}|t_{vi}, v_{i\tau})$$

$$p(\mathbf{T}|\mathbf{A}^t, \mathbf{B}^t) = \prod_{v,i} p(t_{vi}|a_{vi}^t, b_{vi}^t)$$

$$p(\mathbf{V}|\mathbf{A}, \vec{z}) = \prod_{i,\tau} p(v_{i\tau}|\lambda_{ic}, z_\tau)$$

$$p(\mathbf{A}|\mathbf{A}^\lambda, \mathbf{B}^\lambda) = \prod_{i,c} p(\lambda_{ic}|a_{ic}^\lambda, b_{ic}^\lambda).$$

2.4 Variational Learning Algorithm

In order to obtain convenient analytical forms of iterative updates, the variational distribution is chosen to be factorized as

$$q(H) = q(\mathbf{S}, \mathbf{T}, \mathbf{V}, \mathbf{A}) = q(\mathbf{S})q(\mathbf{T})q(\mathbf{V})q(\mathbf{A}),$$

with

$$q(\mathbf{S}) = \prod_{v,\tau} q(s_{v:\tau})$$

$$q(\mathbf{T}) = \prod_{v,i} q(t_{vi})$$

$$q(\mathbf{V}) = \prod_{i,\tau} q(v_{i\tau})$$

$$q(\mathbf{A}) = \prod_{i,c} q(\lambda_{ic}).$$

As the proposed learning algorithm will be laid out in matrix form with computationally efficient matrix operations, hyperparameters and variational parameters of the model are organized as matrices according to Table 1 and Table 2, respectively.

$[\mathbf{X}]_{v\tau} = x_{v\tau}$	$[\mathbf{A}]_{\tau c} = \delta(z_\tau - c)$
$[\mathbf{A}^t]_{vi} = a_{vi}^t$	$[\mathbf{A}^\lambda]_{ic} = a_{ic}^\lambda$
$[\mathbf{B}^t]_{vi} = b_{vi}^t$	$[\mathbf{B}^\lambda]_{ic} = b_{ic}^\lambda$

Table 1. Parameters and hyperparameters of the proposed model

$[\mathbf{E}_t^{(n)}]_{vi} = \langle t_{vi} \rangle^{(n)}$	$[\mathbf{E}_v^{(n)}]_{i\tau} = \langle v_{i\tau} \rangle^{(n)}$
$[\mathbf{L}_t^{(n)}]_{vi} = \langle \log t_{vi} \rangle^{(n)}$	$[\mathbf{L}_v^{(n)}]_{i\tau} = \langle \log v_{i\tau} \rangle^{(n)}$
$[\mathbf{S}_t^{(n)}]_{vi} = \sum_\tau \langle s_{v\tau} \rangle^{(n)}$	$[\mathbf{E}_\lambda^{(n)}]_{ic} = \langle \lambda_{ic} \rangle^{(n)}$
$[\mathbf{S}_v^{(n)}]_{i\tau} = \sum_v \langle s_{v\tau} \rangle^{(n)}$	$[\mathbf{L}_\lambda^{(n)}]_{ic} = \langle \log \lambda_{ic} \rangle^{(n)}$

Table 2. Variational parameters of the proposed model

To derive iterative alternating updates for $q(s_{v:\tau})^{(n+1)}$, applying (5) yields

$$q(s_{v:\tau})^{(n+1)} \propto \exp(\log p(H, D|\theta)) \frac{q(H)^{(n)}}{q(s_{v:\tau})^{(n)}}$$

$$= \text{Multinomial}(s_{v:\tau} | x_{v\tau}, p_{v\tau}^{(n)})$$

with natural parameters

$$p_{v\tau}^{(n)} = \frac{\exp((\log t_{vi})^{(n)} + (\log v_{i\tau})^{(n)})}{\sum_i \exp((\log t_{vi})^{(n)} + (\log v_{i\tau})^{(n)})}.$$

Analytically, variational factor $q(s_{v:\tau})^{(n+1)}$ assumed form of a multinomial distribution. Using analytical form of expectation of sufficient statistics of a multinomial distribution, they get updated according to

$$\langle s_{v\tau} \rangle^{(n+1)} = x_{v\tau} p_{v\tau}^{(n)}. \quad (7)$$

Specifically, as will be seen later, alternating updates for other hidden variables in the model will require those expectations in forms $\sum_v \langle s_{v_i\tau} \rangle^{(n+1)}$ and $\sum_\tau \langle s_{v_i\tau} \rangle^{(n+1)}$, which, by putting summation over (7) and by simple algebraic manipulation, compactly become

$$\sum_v \langle s_{v_i\tau} \rangle^{(n+1)} = \exp(\log v_{i\tau})^{(n)} \sum_v \exp(\log t_{vi})^{(n)} \xi_{v\tau}^{(n)} \quad (8)$$

$$\sum_\tau \langle s_{v_i\tau} \rangle^{(n+1)} = \exp(\log t_{vi})^{(n)} \sum_\tau \xi_{v\tau}^{(n)} \exp(\log v_{i\tau})^{(n)}, \quad (9)$$

with repeating term substituted as

$$\xi_{v\tau}^{(n)} = \frac{x_{v\tau}}{\sum_i \exp(\log t_{vi})^{(n)} \exp(\log v_{i\tau})^{(n)}}. \quad (10)$$

Rewritten in matrix notation, expressions (8), (9) and (10) now become

$$\begin{aligned} \mathbf{\Sigma}_v^{(n+1)} &= \exp \mathbf{L}_v^{(n)} * (\exp \mathbf{L}_t^{(n)})^T * \boldsymbol{\xi}^{(n)} \\ \mathbf{\Sigma}_t^{(n+1)} &= \exp \mathbf{L}_t^{(n)} * (\boldsymbol{\xi}^{(n)} * (\exp \mathbf{L}_v^{(n)})^T) \\ \boldsymbol{\xi}^{(n+1)} &= \mathbf{X} ./ ((\exp \mathbf{L}_t^{(n)}) * (\exp \mathbf{L}_v^{(n)})). \end{aligned}$$

Approached in the same manner, $q(t_{vi})^{(n+1)}$ analytically assumes form of a gamma probability density function:

$$\begin{aligned} q(t_{vi})^{(n+1)} &\propto \exp(\log p(H, D | \theta)) \frac{q(H)^{(n)}}{q(t_{vi})^{(n)}} \\ &= \text{Gamma}(t_{vi} | \alpha_{vi}^t{}^{(n)}, \beta_{vi}^t{}^{(n)}) \end{aligned}$$

with shape and scale parameters

$$\begin{aligned} \alpha_{vi}^t{}^{(n)} &= a_{vi}^t + \sum_\tau \langle s_{v_i\tau} \rangle^{(n)} \\ \beta_{vi}^t{}^{(n)} &= ((b_{vi}^t)^{-1} + \sum_\tau \langle v_{i\tau} \rangle^{(n)})^{-1}, \end{aligned}$$

respectively. Again, it suffices to update and store the expectations of sufficient statistics of $q(t_{vi})^{(n+1)}$,

$$\begin{aligned} \langle t_{vi} \rangle^{(n+1)} &= \alpha_{vi}^t{}^{(n)} \beta_{vi}^t{}^{(n)} \\ \langle \log t_{vi} \rangle^{(n+1)} &= \Psi(\alpha_{vi}^t{}^{(n)}) + \log \beta_{vi}^t{}^{(n)}. \end{aligned}$$

In matrix notation, these updates assume forms

$$\begin{aligned} \mathbf{A}_t^{(n)} &= \mathbf{1} * \mathbf{A}^t + \mathbf{\Sigma}_t^{(n)} \\ \mathbf{B}_t^{(n)} &= \mathbf{1} ./ (\mathbf{1} ./ \mathbf{B}^t + \mathbf{1} * \mathbf{E}_v^{(n)T}) \\ \mathbf{E}_t^{(n+1)} &= \mathbf{A}_t^{(n)} * \mathbf{B}_t^{(n)} \\ \mathbf{L}_t^{(n+1)} &= \Psi(\mathbf{A}_t^{(n)}) + \log \mathbf{B}_t^{(n)}, \end{aligned}$$

where by $\mathbf{1}$ a unity matrix of appropriate size is denoted and $\Psi(\cdot)$ is elementwise digamma function.

Likewise, iterative update equations for $q(v_{i\tau})^{(n+1)}$ are derived as follows:

$$\begin{aligned} q(v_{i\tau})^{(n+1)} &\propto \exp(\log p(H, D | \theta)) \frac{q(H)^{(n)}}{q(v_{i\tau})^{(n)}} \\ &= \text{Gamma}(v_{i\tau} | \alpha_{i\tau}^v{}^{(n)}, \beta_{i\tau}^v{}^{(n)}), \\ \alpha_{i\tau}^v{}^{(n)} &= \mathbf{1} + \sum_v \langle s_{v_i\tau} \rangle^{(n)} \\ \beta_{i\tau}^v{}^{(n)} &= (\sum_c \delta(z_\tau - c)^{(n)} \langle \lambda_{ic} \rangle^{(n)} + \sum_v \langle t_{vi} \rangle^{(n)})^{-1}, \end{aligned}$$

and the corresponding updates are to be done as

$$\langle v_{i\tau} \rangle^{(n+1)} = \alpha_{i\tau}^v{}^{(n)} \beta_{i\tau}^v{}^{(n)}$$

$$\langle \log v_{i\tau} \rangle^{(n+1)} = \Psi(\alpha_{i\tau}^v{}^{(n)}) + \log \beta_{i\tau}^v{}^{(n)},$$

or, in matrix form,

$$\begin{aligned} \mathbf{A}_v^{(n)} &= \mathbf{1} + \mathbf{\Sigma}_v^{(n)} \\ \mathbf{B}_v^{(n)} &= \mathbf{1} ./ (\mathbf{E}_\lambda^{(n)} * \mathbf{\Delta}^{(n)T} + \mathbf{E}_t^{(n)T} * \mathbf{1}) \\ \mathbf{E}_v^{(n+1)} &= \mathbf{A}_v^{(n)} * \mathbf{B}_v^{(n)} \\ \mathbf{L}_v^{(n+1)} &= \Psi(\mathbf{A}_v^{(n)}) + \log \mathbf{B}_v^{(n)}. \end{aligned}$$

Iterative update equations for λ are derived as follows:

$$\begin{aligned} q(\lambda_{ic})^{(n+1)} &\propto \exp(\log p(H, D | \theta)) \frac{q(H)^{(n)}}{q(\lambda_{ic})^{(n)}} \\ &\propto \sum_\tau \langle \log p(v_{i\tau} | \lambda_{ic}, z_\tau) \rangle + \langle \log p(\lambda_{ic} | a_{ic}^\lambda, b_{ic}^\lambda) \rangle, \end{aligned}$$

where all the expectations are taken with respect to

$$q(v_{i\tau})^{(n)} \frac{\prod_c q(\lambda_{ic})^{(n)}}{q(\lambda_{ic})^{(n)}} q(\bar{z})^{(n)}, \text{ not explicitly noted for clarity.}$$

The first term,

$$\begin{aligned} \sum_\tau \langle \log p(v_{i\tau} | \lambda_{ic}, z_\tau) \rangle &\propto \sum_\tau (-\langle v_{i\tau} \rangle \sum_c \delta(z_\tau - c) \lambda_{ic}) \\ &\quad + \sum_\tau \langle \log (\sum_c \delta(z_\tau - c) \lambda_{ic}) \rangle, \end{aligned}$$

is seemingly more difficult than what has been encountered so far due to expectation operator over the logarithmic function. Luckily, by observing that logarithmic function is concave, the lower bound can be relaxed using Jensen's inequality,

$$\langle \log (\sum_c \delta(z_\tau - c) \lambda_{ic}) \rangle \geq \sum_c \delta(z_\tau - c) \langle \log \lambda_{ic} \rangle.$$

In this relaxed lower bound, coordinate ascent now admits a closed form

$$\begin{aligned} q(\lambda_{ic})^{(n+1)} &= \text{Gamma}(\lambda_{ic} | \alpha_{ic}^\lambda{}^{(n)}, \beta_{ic}^\lambda{}^{(n)}), \\ \alpha_{ic}^\lambda{}^{(n)} &= a_{ic}^\lambda + \sum_\tau \delta(z_\tau - c) \\ \beta_{ic}^\lambda{}^{(n)} &= (b_{ic}^\lambda{}^{-1} + \sum_\tau \langle v_{i\tau} \rangle \delta(z_\tau - c))^{-1}. \end{aligned}$$

The corresponding expectations of natural parameters are updated according to

$$\begin{aligned} \langle \lambda_{ic} \rangle^{(n+1)} &= \alpha_{ic}^\lambda{}^{(n)} \beta_{ic}^\lambda{}^{(n)} \\ \langle \log \lambda_{ic} \rangle^{(n+1)} &= \Psi(\alpha_{ic}^\lambda{}^{(n)}) + \log \beta_{ic}^\lambda{}^{(n)}, \end{aligned}$$

which can compactly be rewritten in matrix form as

$$\begin{aligned} \mathbf{A}_\lambda^{(n)} &= \mathbf{A}^\lambda + \mathbf{\Delta}^{(n)} * \mathbf{1} \\ \mathbf{B}_\lambda^{(n)} &= \mathbf{1} ./ (\mathbf{1} ./ \mathbf{B}^\lambda + \mathbf{E}_v^{(n)} * \mathbf{\Delta}^{(n)}) \\ \mathbf{E}_\lambda^{(n+1)} &= \mathbf{A}_\lambda^{(n)} * \mathbf{B}_\lambda^{(n)} \\ \mathbf{L}_\lambda^{(n+1)} &= \Psi(\mathbf{A}_\lambda^{(n)}) + \log \mathbf{B}_\lambda^{(n)}. \end{aligned}$$

The learning algorithm in matrix form is recapitulated in Appendix C of the paper, with the iterative updates in the same order as presented above.

In the presented algorithm, group affiliation variables z_τ are assumed to be fully observed, i.e. groups are to be explicitly defined beforehand depending on the application, as the model does not attempt to learn the group affiliations.

2.5 Variational Bound

Expression for the variational lower bound on the marginal loglikelihood of the observed data is obtained by expanding (2), as presented in Appendix B.

3 EXPERIMENTS

In this subsection performance of a simple classification method based on the proposed group sparse NMF algorithm is investigated on three publicly available benchmark datasets, Yale [28] and ORL [29], [30] for face recognition and JAFFE [31] for facial expression recognition. Yale dataset consists of 11 grayscale images per subject of 15 subjects, one image per different facial expression of configuration: center-light, w/glasses, happy, left-light, w/no glasses, normal, right-light, sad, sleepy, surprised and wink. ORL dataset consists of 10 different images of 40 distinct subjects, taken at different times, varying the lightning, facial expressions (open or close eyes, smiling or not smiling) and facial details (glasses or no glasses). JAFFE dataset is a collection of 213 images of 7 facial expressions - angry, disgust, fear, happy, neutral, sad surprise - posed by 10 Japanese female models.

Results are compared to classification methods based on PCA and related NMF algorithms, and also, where available, to relevant published results on the same datasets with different approaches.

MATLAB/Octave implementation of the algorithm as well as scripts used to generate the results are available from the author's homepage or will have been received upon request.

3.1 Data Preprocessing

Images in Yale dataset used in the experiments have been prepared by the MIT media laboratory [32] - aligned by rotation and centering of the manually determined locations of the eyes and then cropped. Additionally, specifically for this paper, after downsampling the images by factor 0.5 to alleviate computational load, pixelwise masking has been applied to remove most of the background, torso and the hair. Finally, histograms of masked images have been equalized.

In case of ORL dataset, because faces have been taken at different angles, centering has not been attempted. Images have only been downsampled by factor of 0.5, with histogram equalization.

Images from JAFFE dataset have first been roughly aligned by congealing [33], [34], then resized by factor 0.75 and masked leaving only pixels which roughly correspond to locations of faces, followed by histogram equalization.

For all datasets, images have been vectorized in a way that each image is a column vector of input matrix which is to be factorized. Examples of preprocessed images are shown in Fig. 1.

3.2 Experimental Setting

The classification method consists of three consecutive stages: image preprocessing stage, dictionary learning stage and classification stage. Preprocessing is dependent on the dataset and has been described in detail in the preceding subsection. In the dictionary learning stage a dictionary is obtained by the proposed group sparse NMF algorithm, or PCA or the standard sparse NMF algorithms; only in case of the probabilistic group sparsity NMF algorithm are the class labels taken into account,

while in other cases the algorithms cannot include this information straightforwardly and is therefore done in a fully unsupervised manner, which is exactly where the comparative advantage of the proposed algorithm for classification problems lies.

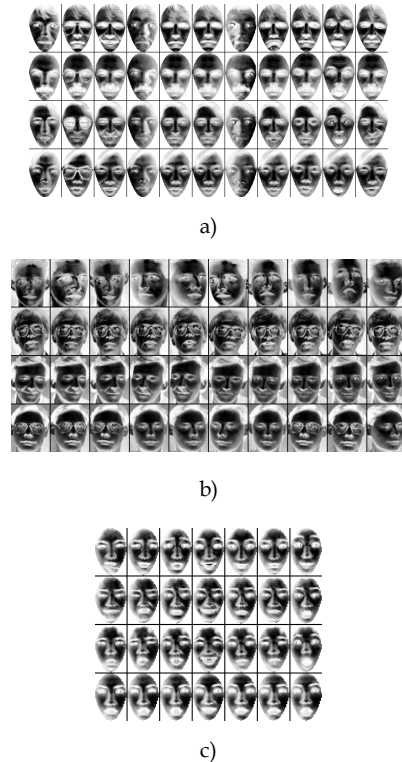


Fig. 1. Examples of preprocessed face samples (in negative): a) Yale, b) ORL, c) JAFFE.

To put more weight on role of the quality of the decomposition obtained in the dictionary learning stage, in a sense of class-to-class discriminative information it is able to bring out by itself rather than on the role of the classifier, the simple classifier of choice in the classification stage is 1-nearest neighbor with cosine distance metric.

Performance quantifiers have been obtained by 5 runs of 10-fold crossvalidation for 10 different random restarts of NMF algorithms; at each pass dictionary learning by one of the algorithms had been performed on the training set only, followed by obtaining the representation of the test set in the feature space (in which the classifier had been built) by linear least squares with nonnegativity constraints [35]. NMF parameter optimization has been done using parameter sweeps using the previously mentioned crossvalidation scheme to obtain performance measures; the criterion for parameter selection is chosen to be highest crossvalidated estimate of maximal accuracy out of 10 NMF restarts with random initializations. Having found the best set of parameters for each of the NMF algorithms, reported accuracies are

- 1) crossvalidated estimates of maximal accuracy (the criterion itself) out of 10 random restarts
- 2) mean and variance of accuracies of the entire crossvalidation scheme.

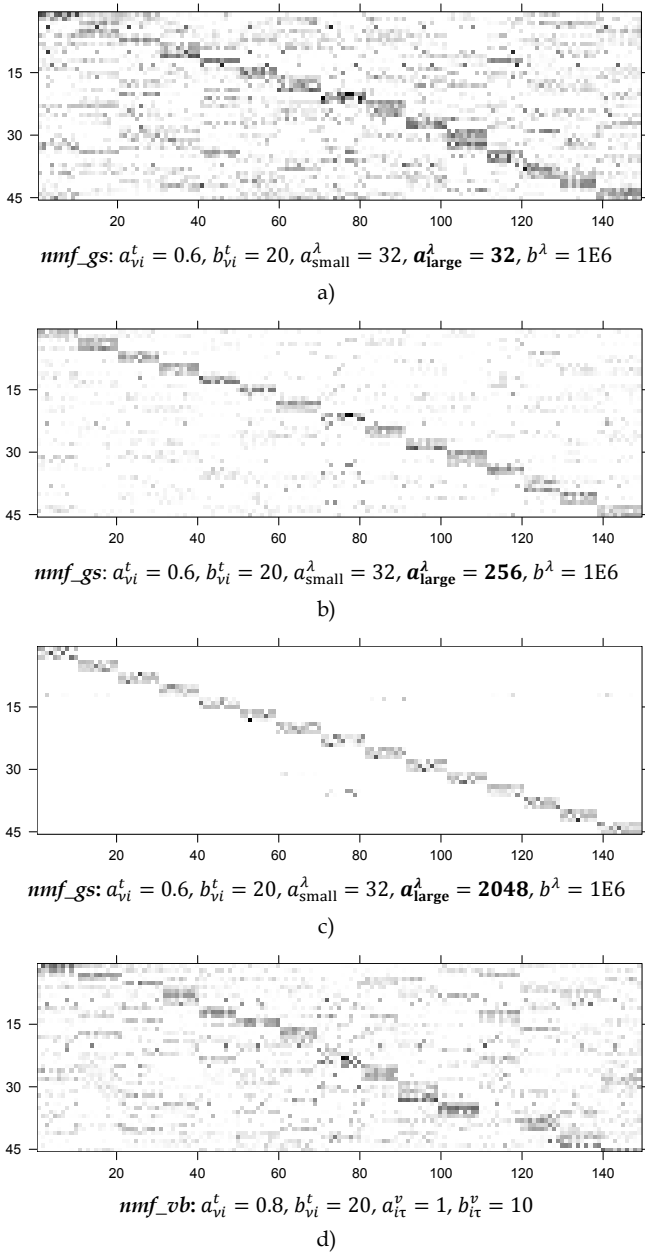


Fig. 2. Yale dataset, samples in feature space, i.e. the coefficient matrix E_v , appropriately sorted: a), b) and c) *nmf_gs* with increasing parameter a_{large}^λ ; d) *nmf_vb*, with parameters listed below their corresponding subimages. The darker the shade of grey, the higher the magnitude.

Number of iterations for family of NMF algorithms has been set to 300. In case of PCA, data is projected to space of most significant principal components, and the crossvalidated estimate of the accuracy is reported. Decomposition algorithms used in the experiments for comparison will be referred to in short as

- 1) *pca*, principal component analysis [36]
- 2) *nmf_kl*, NMF with KL-divergence that includes a weighted penalty term to encourage sparsity in the right matrix [6]

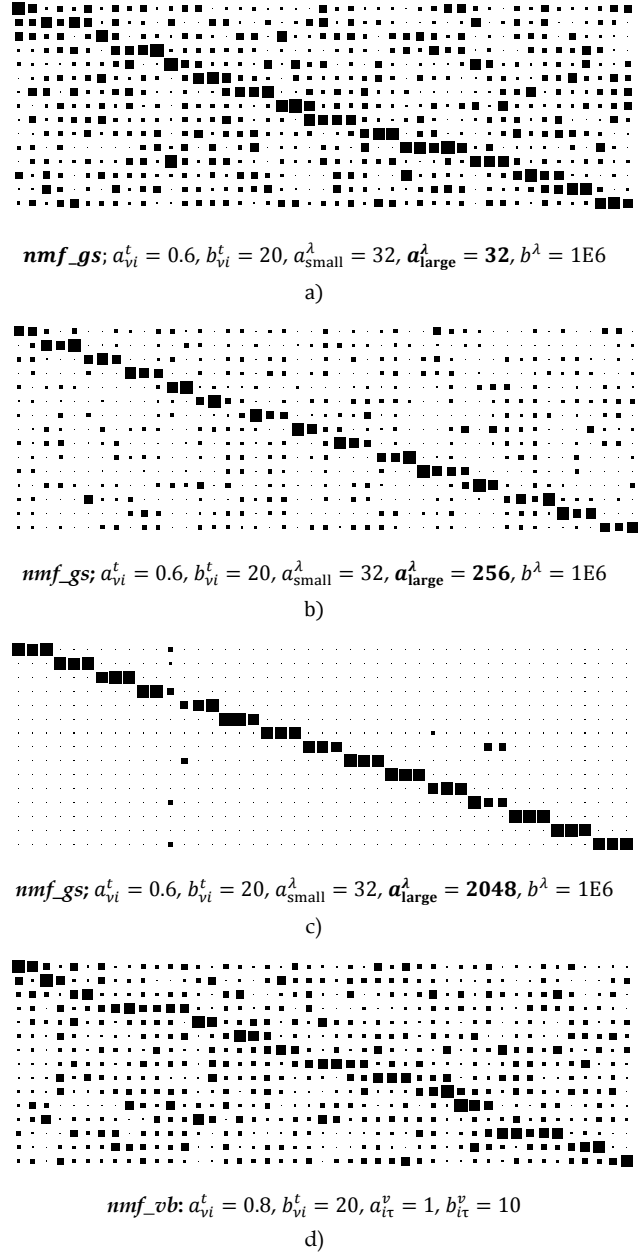


Fig. 3. Controlling prevalences of features in labels on Yale dataset - Hinton diagrams of normalized l1 norm of the coefficient matrix E_v accumulated across labels: a), b) and c) *nmf_gs* with increasing parameter a_{large}^λ ; d) *nmf_vb*, with parameters listed below their corresponding subimages. Rows correspond to labels and columns to features.

- 3) *nmf_vb*, a variational Bayesian NMF which can model sparse decompositions [11]
- 4) *nmf_gs*, the proposed variational Bayesian NMF with group sparsity constraints.

The implementation *nmf_vb* is available from [37] and *nmf_kl* is short for *nmf_kl_sparse_v*, a part of NMFlib v0.1.3 library for MATLAB [38].

3.3 Results and Discussion

First, effect of hyperpriors α_{ic}^λ and β_{ic}^λ on continuous indicators of presence of a feature across group labels, λ_{ic} are

going to the examined. Specifically, it will be shown how features representative of specific group labels, i.e. which are prevalent in specific group labels, can be extracted from data this way.

Consider the hyperpriors A^λ and B^λ of forms

$$\begin{aligned} A^\lambda &= \begin{bmatrix} \vec{a}_1^\lambda & \dots & \vec{a}_1^\lambda & \vec{a}_2^\lambda & \dots & \vec{a}_2^\lambda & \dots & \vec{a}_c^\lambda & \dots & \vec{a}_c^\lambda \end{bmatrix}^T, \\ \vec{a}_1^\lambda &= [a_{\text{small}}^\lambda \quad a_{\text{large}}^\lambda \quad \dots \quad a_{\text{large}}^\lambda \quad a_{\text{large}}^\lambda]^T, \\ \vec{a}_2^\lambda &= [a_{\text{large}}^\lambda \quad a_{\text{small}}^\lambda \quad \dots \quad a_{\text{large}}^\lambda \quad a_{\text{large}}^\lambda]^T, \\ &\vdots \\ \vec{a}_c^\lambda &= [a_{\text{large}}^\lambda \quad a_{\text{large}}^\lambda \quad \dots \quad a_{\text{large}}^\lambda \quad a_{\text{small}}^\lambda]^T, \\ [B^\lambda]_{ic} &= b^\lambda, \end{aligned} \quad (11)$$

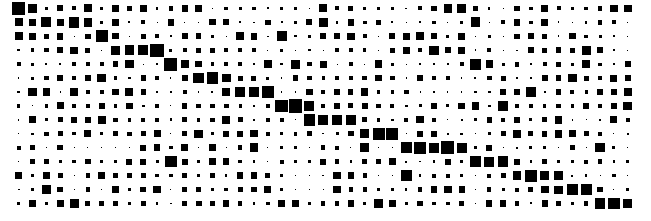
where C denotes number of groups. If a_{small}^λ is smaller than a_{large}^λ , the hyperprior is such that rows of A^λ which contain a_c^λ at some set of indices S will bias the corresponding λ_{ic}^{-1} , $i \in S$ towards larger values, which hierarchically propagates to the related hidden factors v_i , $i \in S$, giving their elements $v_{i\tau}$ a sparse prior with large expected value if $z_\tau = c$ and with small expected value otherwise. Thus, such a hyperprior describes the tendency of the coefficients $v_{i\tau}$ to be significantly large in samples belonging to a single group only. In the presented experiments number of such representative coefficients is chosen to be equally distributed among groups, i.e. each group will have the same number of representative features bound.

Behavior of prior (11) on the Yale dataset can be eyeballed from Fig. 2, Fig. 3 and Fig. 4, which all correspond to and visualize qualitatively typical mappings in the feature space, representative of several chosen parameter setups. Specifically, enforced by the prior (11) here are 3 representative features per each label; note that the dataset has 15 labels, which totals to 45 features. The number of samples in the training set is 149.

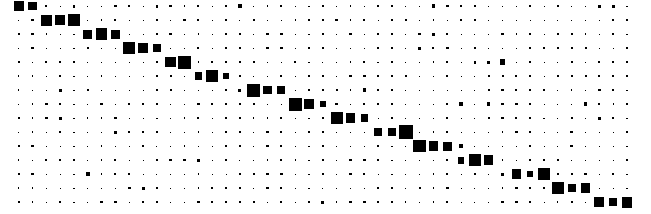
Fig. 2 presents samples in the feature space as heatmaps, having the samples sorted according to their labels and the coefficients according to their cumulative l1 norm per label averaged by number of samples per label. As a baseline, an example of *nmf_vb* decomposition is presented in Fig. 2d).

What is observed is that via magnitude of difference between a_{large}^λ and a_{small}^λ degree of mixing between specified group-prevalent features can be controlled. When $a_{\text{large}}^\lambda = a_{\text{small}}^\lambda$, group sparse decompositions are obtained with no prior which would bias distinct features to be prevalent across specific groups, as depicted by Fig. 2a). The case when $a_{\text{large}}^\lambda \gg a_{\text{small}}^\lambda$ resembles the result of concatenating the NMF decompositions obtained for each label separately, i.e. each label has its group of features which are groupwise very strictly separated in terms of mixing, represented by Fig. 2c). Between those two extremes, a sweet spot for obtaining representation spaces with good discriminative properties may be found, a case which relates to Fig. 2b). Indicative of the justifiedness of this line of reasoning is the fact that Fig. 2b) has been obtained using parameters which yielded the best

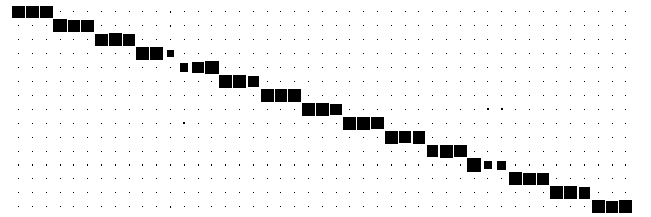
crossvalidated classification performance, as reported in Table 3.



nmf_gs; $a_{vi}^t = 0.6$, $b_{vi}^t = 20$, $a_{\text{small}}^\lambda = 32$, $a_{\text{large}}^\lambda = 32$, $b^\lambda = 1E6$
a)



nmf_gs; $a_{vi}^t = 0.6$, $b_{vi}^t = 20$, $a_{\text{small}}^\lambda = 32$, $a_{\text{large}}^\lambda = 256$, $b^\lambda = 1E6$
b)



nmf_gs; $a_{vi}^t = 0.6$, $b_{vi}^t = 20$, $a_{\text{small}}^\lambda = 32$, $a_{\text{large}}^\lambda = 2048$, $b^\lambda = 1E6$
c)

Fig. 4. Yale dataset, Hinton diagrams of l1 norm of the reciprocal of group indicator variables E_λ^{-1} accumulated across labels for *nmf_gs* with increasing parameter a_{large}^λ : a) $a_{\text{large}}^\lambda = 32$, b) $a_{\text{large}}^\lambda = 256$, c) $a_{\text{large}}^\lambda = 2048$. Rows correspond to labels and columns to features.

From a different perspective, the same structural pattern over labels can be recognized using quantifiers as in Fig. 3 where, rather than presented directly across samples, l1 norm of the coefficients has been accumulated across samples having same label then normalized by cardinalities of the corresponding labels. Furthermore, in Fig. 4 diagrams of the same type but having the reciprocal of the group indicator matrix, E_λ^{-1} , as the target variable bear resemblance of a high degree to Fig. 3 a) b) and c), reason for which is that E_λ directly specifies the prevalences of the features in each of the groups and propagates them hierarchically to variables $v_{i\tau}$, according to (5).

Impact of choice of a_{small}^λ and a_{large}^λ on classification performance is illustrated on Fig. 6, Fig. 7 and Fig. 10 for Yale, ORL and JAFFE datasets, respectively. On Yale dataset, for a fixed a_{small}^λ as a_{large}^λ increases the accuracy improves, hitting a peak after which it begins to deteriorate, but not below the case when $a_{\text{small}}^\lambda = a_{\text{large}}^\lambda$. Qualitatively similar is the behavior on ORL dataset, but the

accuracy improvement is more modest. On JAFFE dataset, improvement in the accuracy is notable but, more significantly, a pronounced droop when $a_{\text{large}}^\lambda \gg a_{\text{small}}^\lambda$ is present.

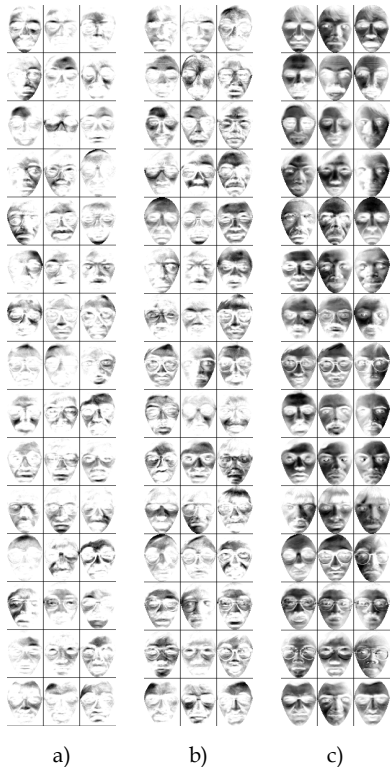


Fig. 5. Yale dataset, extracted features (in negative) for nmf_gs with increasing parameter a_{large}^λ : a) $a_{\text{large}}^\lambda = 32$, b) $a_{\text{large}}^\lambda = 256$, c) $a_{\text{large}}^\lambda = 2048$. Other parameters are $a_{vi}^t = 0.6$, $b_{vi}^t = 20$, $a_{\text{small}}^\lambda = 32$, $b^\lambda = 1E6$.

Explanation for this effect is that, on JAFFE dataset, the droop is caused by a too restrictive mixing (Fig. 8c)), resulting in decompositions where same subjects with different expressions exhibit hardly any common features, i.e. expression-independent and subject-specific information is not shared between groups, and, consequently, the features are forced to be too holistic to extract expressions exclusively, (Fig. 9c)). The same effect is behind the behavior with the Yale dataset, however, Yale dataset is such that distances between same subjects having different expressions or configuration is smaller than distances between different subjects with same expressions or under same configuration, allowing satisfactory class discrimination even with such extremely holistic features (Fig. 5). To remind the reader, the goal on Yale dataset is subject recognition regardless of different expressions and configuration and on JAFFE face expression recognition regardless of the subject making it.

Experimental results on the Yale dataset are summarized in Table 3. Compared to the classification using nmf_kl and nmf_vb , improvement in the performance turned out to be significant when classification is performed in feature subspaces obtained by nmf_gs . It showed significantly higher average peak performance and higher average performance, with smaller variance also.

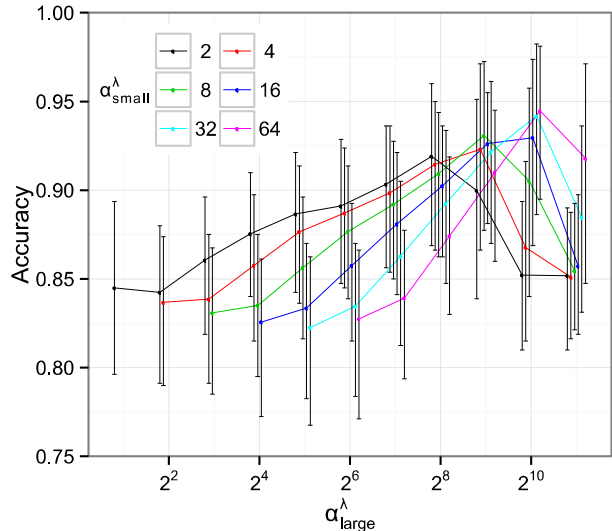


Fig. 6. Yale dataset, dependency of mean classification accuracy using nmf_gs on a_{small}^λ and a_{large}^λ . Error bars represent crossvalidated estimates of maximal and minimal accuracies of a number of nmf_gs runs. Other parameters are fixed as $a_{vi}^t = 0.6$, $b_{vi}^t = 20$, $b^\lambda = 1E6$.

Algorithm	pca	nmf_kl	nmf_vb	nmf_gs
Accuracy				
maximum	0.7987	0.9267	0.8875	0.9825
mean \pm variance	0.7987 \pm 0.0000	0.8690 \pm 0.0054	0.8246 \pm 0.0077	0.9417 \pm 0.0036
Subspace dimension	73	120	60	45

Table 3. Classification results on the Yale dataset

On ORL dataset, as shown in Table 4, improvements are still observed, but to a far lesser extent. In the community, classification problem on the ORL dataset is known to lay on the easier side, as pca alone gives high accuracy. Regarding NMF algorithms, it can be concluded that sparsity constraints only are sufficient to give performance of high quality, leaving little space for improvement by group sparsity constraints.

Similar are the results on JAFFE dataset, presented in Table 5, but the improvement when using $gsNMF$ is less marginal - compared to the results in Table 6 it can be seen that in this case nmf_gs in conjunction with 1-NN classifier can output classification results on the level of [39], where discrete wavelet transform with 2D linear discriminant analysis (LDA) is used to find features followed by classification using support vector machines with different kernel choices. Somewhat lower accuracies have been reported in [40] and [41]: experimental setup used in [40] consists of processing the samples by Gabor filtering, then sampling at fiducial points followed by PCA to get the features, finalized by LDA as classifier, and [41] uses Gabor filtering to get the features and two-layered perceptron for label discrimination. In [42] the authors use classifier based on Gaussian processes in the original pixel space.

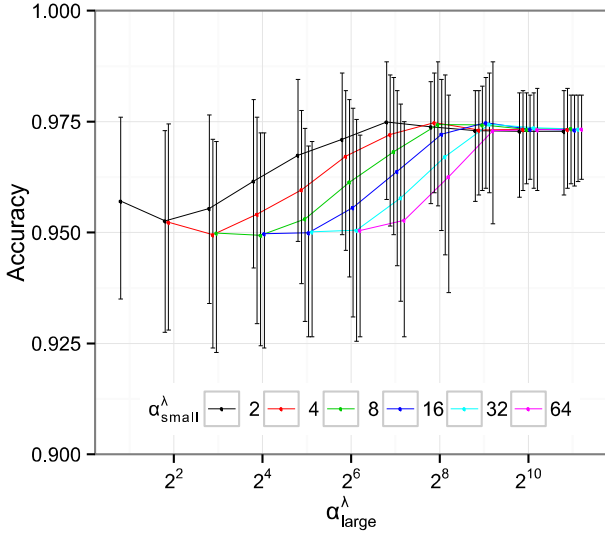


Fig. 7. ORL dataset, dependency of mean classification accuracy using *nmf_gs* on $\alpha_{\text{small}}^\lambda$ and $\alpha_{\text{large}}^\lambda$. Error bars represent crossvalidated estimates of maximal and minimal accuracies of a number of *nmf_gs* runs. Other parameters are fixed as $\alpha_{vi}^t = 0.5$, $b_{vi}^t = 10$, $b^\lambda = 1E6$.

Algorithm	<i>pca</i>	<i>nmf_kl</i>	<i>nmf_vb</i>	<i>nmf_gs</i>
Accuracy				
maximum	0.9605	0.9775	0.9780	0.9885
mean \pm variance	0.9605 \pm 0.0000	0.9529 \pm 0.0013	0.9511 \pm 0.0013	0.9750 \pm 0.0006
Subspace dimension	67	140	60	160

Table 4. Classification results on the ORL dataset

Several other reported results on the same datasets are known to the author but are unfortunately of little use as the results have been evaluated nonuniformly across publications.

From the practical point of view, however, even though peak performance of classification with *nmf_gs* is admirable, the problem of *a priori* selection of a *nmf_gs* decomposition which is bound to produce this peak remains open. This problem is not characteristic only of *nmf_gs*, but also of other NMF methods due to dependency of decompositions on initial values. Even though variational Bayes methodology allows calculation of variational bound which model comparisons can be made based upon, in the presented experiments the variational bound has been found to be uncorrelated with the classification accuracy, which is attributed to the fact that the classifier stands outside the Bayesian framework, i.e. that no objective directly connected with the classification had been embedded in the probabilistic model. Still, a solution always remains, which is to evaluate classification performance on a separate validation set and use it as an optimality indicator to determine which dictionary to select.

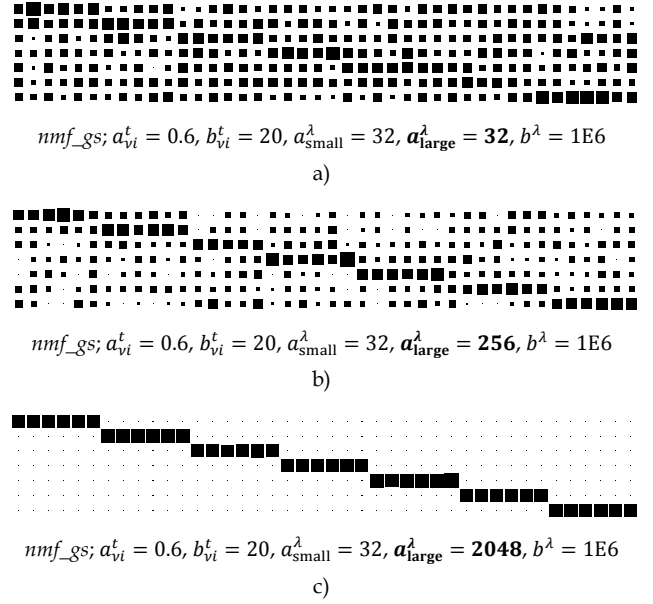


Fig. 8. Controlling prevalences of features in labels on JAFFE dataset - Hinton diagrams of normalized l1 norm of the coefficient matrix E_v accumulated across labels for *nmf_gs* with increasing parameter $\alpha_{\text{large}}^\lambda$: a) $\alpha_{\text{large}}^\lambda = 32$, b) $\alpha_{\text{large}}^\lambda = 256$, c) $\alpha_{\text{large}}^\lambda = 2048$. Rows correspond to labels and columns to features.

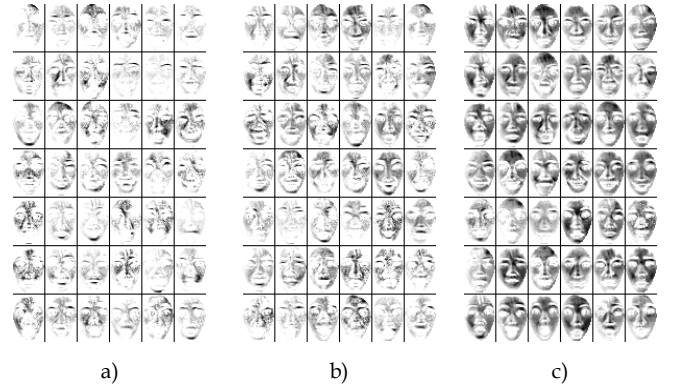


Fig. 9. JAFFE dataset, extracted features (in negative) for *nmf_gs* with increasing parameter $\alpha_{\text{large}}^\lambda$: a) $\alpha_{\text{large}}^\lambda = 32$, b) $\alpha_{\text{large}}^\lambda = 256$, c) $\alpha_{\text{large}}^\lambda = 2048$. Other parameters are $\alpha_{vi}^t = 0.6$, $b_{vi}^t = 20$, $\alpha_{\text{small}}^\lambda = 32$, $b^\lambda = 1E6$.

4 CONCLUSION

A probabilistic formulation of NMF with group sparsity constraints has been laid out with an efficient variational Bayesian algorithm for approximate learning of the model parameters. It has been shown how prevalence of specific features across groups and the degree of their mixing between groups can be controlled. Having identity-mapped the class labels to a more general notion of groups, the presented model has been utilized as a supervised feature extractor in face recognition and facial expression recognition applications and beneficial effects of such decomposition subspaces on classification performance have been observed.

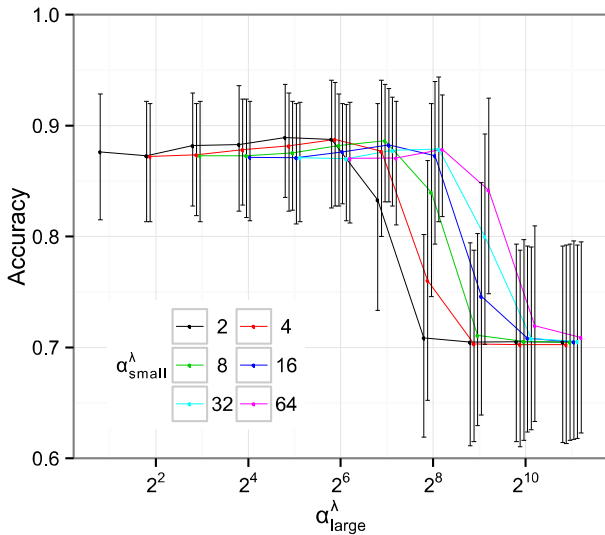


Fig. 10. JAFFE dataset, dependency of mean classification accuracy using *nmf_gs* on a_{small}^λ and a_{large}^λ . Error bars represent crossvalidated estimates of maximal and minimal accuracies of a number of *nmf_gs* runs. Other parameters are fixed as $a_{v_i}^t = 0.5$, $b_{v_i}^t = 20$, $b^\lambda = 1E6$.

To contribute to community where easily reproducible research is appreciated, the implementation used in the experiments is made publicly available (which unfortunately is not the case still too often).

Algorithm	<i>Pca</i>	<i>nmf_kl</i>	<i>nmf_vb</i>	<i>nmf_gs</i>
Accuracy				
maximum	0.9067	0.9267	0.9295	0.9438
mean \pm variance	0.9067 \pm 0.0000	0.8690 \pm 0.0054	0.8597 \pm 0.0063	0.8789 \pm 0.0051
Subspace dimension	18	28	14	42

Table 5. Classification results on the JAFFE dataset

Method	Shih et al. [39]	Lyons et al. [40]	Zhang et al. [41]	Cheng et al. [42]
Accuracy	0.9413	0.92	0.901	0.8689

Table 6. Relevant classification results on the JAFFE dataset reported in the literature where accuracy has been obtained by 10-fold crossvalidation

Future work on the presented subject includes pursuing modifications of the group sparse formulation which would work in semi-supervised settings. The model should ideally allow efficient inference and learning of the labels of unlabeled data, avoiding sampling techniques for execution speed if possible.

ACKNOWLEDGMENTS

This work was supported by the Croatian Ministry of Science, Education and Sports through the project "Computational Intelligence Methods in Measurement Systems", No. 098-0982560-2565

BIBLIOGRAPHY

- [1] D. D. Lee and H. S. Seung, "Learning the parts of objects by non-negative matrix factorization.," *Nature*, vol. 401, pp. 788–791, 1999.
- [2] D. Donoho and V. Stodden, "When does non-negative matrix factorization give a correct decomposition into parts?," *NIPS*, vol. 16, pp. 1141–1148, 2003.
- [3] P. O. Hoyer, "Non-negative matrix factorization with sparseness constraints," *Journal of Machine Learning Research*, vol. 5, pp. 1457–1469, 2004.
- [4] S. Z. Li, X. Hou, H. Zhang, and Q. Cheng, "Learning spatially localized, parts-based representation," *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition CVPR 2001*, vol. 1, pp. 207–212, 2001.
- [5] A. Pascual-Montano, J. M. Carazo, K. Kochi, D. Lehmann, and R. D. Pascual-Marqui, "Nonsmooth nonnegative matrix factorization (nsNMF).," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 28, no. 3, pp. 403–15, Mar. 2006.
- [6] T. V. T. Virtanen, "Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria," *IEEE Transactions On Audio Speech And Language Processing*, vol. 15, pp. 1066–1074, 2007.
- [7] R. Zhi, M. Flierl, Q. Ruan, and W. B. Kleijn, "Graph-preserving sparse nonnegative matrix factorization with application to facial expression recognition," *IEEE Transactions on Systems, Man and Cybernetics*, vol. 41, pp. 38–52, 2011.
- [8] D. Cai, X. He, J. Han, and T. S. Huang, "Graph regularized non-negative matrix factorization for data representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, pp. 1548–1560, 2010.
- [9] Z. Zhang and K. Zhao, "Low-rank matrix approximation with manifold regularization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, pp. 1717–29, 2013.
- [10] C. Fevotte and T. Cemgil, "Nonnegative matrix factorization as probabilistic inference in composite models," *Proc. 17th European Signal Processing Conference (EUSIPCO'09)*, pp. 1–5, 2009.
- [11] A. T. Cemgil, "Bayesian inference for nonnegative matrix factorisation models," *Computational Intelligence and Neuroscience*, vol. 2009, Jan. 2009.
- [12] O. Dikmen and C. Fevotte, "Maximum marginal likelihood estimation for nonnegative dictionary learning in the gamma-Poisson model," *IEEE Transactions on Signal Processing*, vol. 60, no. 10, pp. 5163–5175, Oct. 2012.
- [13] M. N. Schmidt and H. Laurberg, "Nonnegative matrix factorization with Gaussian process priors," *Computational intelligence and neuroscience*, 2008.

- [14] T. Virtanen, A. T. Cemgil, and S. Godsill, "Bayesian extensions to non-negative matrix factorisation for audio signal modelling," *2008 IEEE International Conference on Acoustics Speech and Signal Processing*, pp. 1825–1828, 2008.
- [15] D. Blei, P. Cook, and M. Hoffman, "Bayesian nonparametric matrix factorization for recorded music," *In Proceedings of the International Conference on Machine Learning (ICML)*, pp. 439–446.
- [16] H. Lee and S. Choi, "Group nonnegative matrix factorization for EEG classification," *Journal of Machine Learning Research - Proceedings*, vol. 5, pp. 320–327, 2009.
- [17] J. Kim, R. Monteiro, and H. Park, "Group sparsity in nonnegative matrix factorization," *In proc. of the 2012 SIAM International Conference on Data Mining (SDM)*, pp. 851–862, 2012.
- [18] B. Shin and A. Oh, "Bayesian group nonnegative matrix factorization for EEG analysis," *arXiv:1212.4347*, 2012.
- [19] A. Lefevre, F. Bach, and C. Févotte, "Itakura-Saito nonnegative matrix factorization with group sparsity," *IEEE International Conference on Acoustics Speech and Signal Processing*, vol. 31, pp. 1–4, 2011.
- [20] J.-T. Chien and H.-L. Hsieh, "Bayesian group sparse learning for music source separation," *EURASIP Journal on Audio, Speech, and Music Processing*, vol. 18, p. 15, 2013.
- [21] S. Zafeiriou, A. Tefas, I. Buciu, and I. Pitas, "Exploiting discriminant information in nonnegative matrix factorization with application to frontal face verification," *IEEE Transactions on Neural Networks*, vol. 17, pp. 683–695, 2006.
- [22] Y. Zhu, L. Jing, and J. Yu, "Text clustering via constrained nonnegative matrix factorization," *2011 IEEE 11th International Conference on Data Mining*, vol. 0, pp. 1278–1283, 2011.
- [23] R. He, W.-S. Zheng, B.-G. Hu, and X.-W. Kong, "Nonnegative sparse coding for discriminative semi-supervised learning," *CVPR 2011*, pp. 2849–2856, 2011.
- [24] H. Liu, Z. Wu, S. Member, and X. Li, "Constrained nonnegative matrix factorization for image representation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, pp. 1299–1311, 2012.
- [25] C. Goutte and E. Gaussier, "Relation between PLSA and NMF and implications," *Proceedings of the 28th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 601–602, 2005.
- [26] J. M. Winn, "Variational message passing and its applications," *Ph.D. thesis, Department of Physics, University of Cambridge*, 2003.
- [27] P. J. Garrigues and B. A. Olshausen, "Group sparse coding with a Laplacian scale mixture prior," *Advances in Neural Information Processing Systems*, vol. 23, pp. 1–9, 2010.
- [28] "Yale face database." [Online]. Available: <http://cvc.yale.edu/projects/yalefaces/yalefaces.html>. [Accessed: 01-Jul-2013].
- [29] F. S. Samaria and A. C. Harter, "Parameterisation of a stochastic model for human face identification," *Proceedings of 1994 IEEE Workshop on Applications of Computer Vision*, vol. 13, pp. 138–142, 1994.
- [30] AT&T Laboratories Cambridge, "The database of faces," 2002. [Online]. Available: <http://www.cl.cam.ac.uk/research/dtg/attarchive/facedatabase.html>. [Accessed: 01-Jul-2013].
- [31] M. J. Lyons, M. Kamachi, and J. Gyoba, "Japanese female facial expressions (JAFFE), database of digital images," 1997. [Online]. Available: <http://www.kasrl.org/jaffe.html>. [Accessed: 01-Jul-2013].
- [32] Vision and Modeling Group MIT Media Laboratory, "The normalized Yale face database." [Online]. Available: <http://vismod.media.mit.edu/vismod/classes/mas622-00/datasets/>. [Accessed: 01-Jul-2013].
- [33] G. B. Huang, V. Jain, and E. Learned-Miller, "Unsupervised joint alignment of complex images," *IEEE 11th International Conference on Computer Vision (2007)*, pp. 1–8, 2007.
- [34] University of Massachusetts Vision Laboratory, "Methods and theory," *Congeaing - source code*. [Online]. Available: <http://vis-www.cs.umass.edu/congeal.html>. [Accessed: 01-Jul-2013].
- [35] C. L. Lawson and R. J. Hanson, "Solving least squares problems," vol. 15, Prentice-Hall, 1974, p. 337.
- [36] H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of Educational Psychology*, vol. 24, pp. 417–441, 1933.
- [37] A. T. Cemgil, "Variational Bayesian nonnegative matrix factorization." [Online]. Available: <http://www.cmpe.boun.edu.tr/~cemgil/bnmf/>. [Accessed: 01-Jul-2013].
- [38] G. Grindlay, "NMFLib - Efficient Matlab library implementing a number of common NMF variants." [Online]. Available: <https://code.google.com/p/nmflib/>. [Accessed: 01-Jul-2013].
- [39] F. Y. Shih, C. F. Chuang, and P. S. P. Wang, "Performance comparisons of facial expression recognition in JAFFE database," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 22, pp. 445–459, 2008.
- [40] M. J. Lyons, J. Budynek, and S. Akamatsu, "Automatic classification of single facial images," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 21, pp. 1357–1362, 1999.
- [41] Z. Zhang, M. Lyons, M. Schuster, and S. Akamatsu, "Comparison between geometry-based and Gabor-wavelets-based facial expression recognition using multi-layer perceptron," *Proceedings Third IEEE International Conference on Automatic Face and Gesture Recognition*, pp. 454–459, 1998.
- [42] F. C. F. Cheng, J. Y. J. Yu, and H. X. H. Xiong, "Facial expression recognition in JAFFE dataset based on Gaussian process classification," *IEEE Transactions on Neural Networks*, vol. 21, pp. 1685–1690, 2010.

Ivan Ivek received his B.S. in electrical engineering from the Faculty of Electrical Engineering and Computing, University of Zagreb, Zagreb, Croatia in 2007. He joined the Ruđer Bošković Institute in 2008, where his current position is Research Assistant for the Computational Intelligence Methods in Measurement Systems project. Currently, he is pursuing his Ph.D. in electronics at the Faculty of Electrical Engineering and Computing, University of Zagreb, Croatia.

Appendix A - Probability Density Functions

By $\mathbf{1}$ (column) vector of ones and by $\Psi(\cdot)$ elementwise digamma function, $\Psi(x) = \frac{d}{dx} \log \Gamma(x)$.

<p><i>Poisson</i>($x \lambda$) = $\exp(-\lambda + x \log(\lambda) - \log \Gamma(x + 1))$,</p> <p>with $\lambda > 0$</p> <p>$\langle x \rangle = \lambda$;</p>
<p><i>Gamma</i>($x a, b$) = $\exp\left(-\frac{1}{b}x + (a - 1) \log x - a \log b - \log \Gamma(a)\right)$,</p> <p>with $a > 0, b > 0$;</p> <p>$\langle x \rangle = ab$,</p> <p>$\langle \log x \rangle = \Psi(a) + \log b$;</p> <p>$\mathcal{H}(x) = -(a - 1)\psi(a) + \log b + a + \log \Gamma(a)$;</p>
<p><i>Exponential</i>($x b$) = <i>Gamma</i>($x 1, b$);</p>
<p><i>Multinomial</i> $\left(\begin{bmatrix} x_1 \\ \vdots \\ x_C \end{bmatrix} \middle s, \begin{bmatrix} p_1 \\ \vdots \\ p_C \end{bmatrix}\right) = \delta(s - \sum_i x_i) \exp(\log \Gamma(s + 1) + \sum_i (x_i \log p_i - \log \Gamma(x_i + 1)))$,</p> <p>with $\sum_i p_i = 1, \sum_i x_i = s$;</p> <p>$\begin{bmatrix} \langle \log x_1 \rangle \\ \vdots \\ \langle \log x_C \rangle \end{bmatrix} = s \begin{bmatrix} p_1 \\ \vdots \\ p_C \end{bmatrix}$;</p> <p>$\mathcal{H}(\vec{x}) = -\log \Gamma(s + 1) - \sum_i \langle x_i \rangle \log p_i + \sum_i \langle \log \Gamma(x_i + 1) \rangle - \langle \log \delta(s - \sum_i x_i) \rangle$;</p>
<p><i>Dirichlet</i> $\left(\begin{bmatrix} x_1 \\ \vdots \\ x_C \end{bmatrix} \middle \begin{bmatrix} u_1 \\ \vdots \\ u_C \end{bmatrix}\right) = \exp\left(\begin{bmatrix} u_1 - 1 \\ \vdots \\ u_C - 1 \end{bmatrix}^T \begin{bmatrix} \log x_1 \\ \vdots \\ \log x_C \end{bmatrix} + \log \Gamma(\sum_c^C u_c) - \sum_c^C \log \Gamma(u_c)\right)$,</p> <p>with $\sum_c x_c = 1, u_c > 0$;</p> <p>$\begin{bmatrix} \langle \log x_1 \rangle \\ \vdots \\ \langle \log x_C \rangle \end{bmatrix} = \begin{bmatrix} \Psi(u_1) \\ \vdots \\ \Psi(u_C) \end{bmatrix} - \Psi(\sum_c^C u_c)$;</p> <p>$\mathcal{H}\left(x \begin{bmatrix} x_1 \\ \vdots \\ x_C \end{bmatrix}\right) = -\log \Gamma(\sum_{c=1}^C u_c) + \Gamma(\sum_{c=1}^C u_c) + (\sum_{c=1}^C u_c - C)\psi(\sum_{c=1}^C u_c) - \sum_{c=1}^C (u_c - 1)\psi(u_c)$;</p>
<p><i>Discrete</i>($x \middle \begin{bmatrix} \log p_1 \\ \vdots \\ \log p_C \end{bmatrix}$) = $\exp\left(\begin{bmatrix} \log p_1 \\ \vdots \\ \log p_C \end{bmatrix}^T \begin{bmatrix} \delta(x - 1) \\ \vdots \\ \delta(x - C) \end{bmatrix}\right)$,</p> <p>with $\sum_i p_i = 1$;</p> <p>$\begin{bmatrix} \langle \delta(x - 1) \rangle \\ \vdots \\ \langle \delta(x - C) \rangle \end{bmatrix} = \begin{bmatrix} p_1 \\ \vdots \\ p_C \end{bmatrix}^T$;</p> <p>$\mathcal{H}(x) = -\sum_{c=1}^C p_c \log p_c$;</p>

Appendix B - Variational Bound

Applying (2) on the presented model, the bound is expanded as

$$\begin{aligned}
\mathcal{L}(q)^{(n)} &= \sum_{v,\tau} \langle \log p(x_{v\tau} | s_{v:\tau}) \rangle_{q(s_{v:\tau})}^{(n)} \\
&+ \sum_{v,i,\tau} \langle \log p(s_{v\tau} | t_{vi}, v_{i\tau}) \rangle_{q(s_{v\tau})}^{(n)} q(t_{vi})^{(n)} q(v_{i\tau})^{(n)} + \sum_{v,i,\tau} \mathcal{H}(q(s_{v\tau})^{(n)}) \\
&+ \sum_{v,i} \langle \log p(t_{vi} | a_{vi}^t, b_{vi}^t) \rangle_{q(t_{vi})}^{(n)} + \sum_{v,i} \mathcal{H}(q(t_{vi})^{(n)}) \\
&+ \sum_{i,\tau} \langle \log p(v_{i\tau} | z_\tau, \lambda_{i\cdot}) \rangle_{q(v_{i\tau})}^{(n)} q(z_\tau)^{(n)} q(\lambda_{i\cdot})^{(n)} + \sum_{i,\tau} \mathcal{H}(q(v_{i\tau})^{(n)}) \\
&+ \sum_{i,c} \langle \log p(\lambda_{ic} | a_{ic}^\lambda, b_{ic}^\lambda) \rangle_{q(\lambda_{ic})}^{(n)} + \sum_{i,c} \mathcal{H}(q(\lambda_{ic})^{(n)}) \\
&+ \sum_{\tau,c} \langle \log p(z_\tau | \pi_\tau) \rangle_{q(z_\tau)^{(n)} q(\pi_\tau)^{(n)}} + \sum_{\tau} \mathcal{H}(q(z_\tau)^{(n)}) \\
&+ \sum_{\tau} \langle p([\pi_{\tau 1}, \dots, \pi_{\tau c}]^T | [u_{\tau 1}, \dots, u_{\tau c}]^T) \rangle_{q(\pi_\tau)^{(n)}} + \sum_{\tau} \mathcal{H}(q([\pi_{\tau 1}, \dots, \pi_{\tau c}]^T)^{(n)}).
\end{aligned}$$

With update expressions ordered as in Appendix C and bound calculated where marked, substitutions

$$\begin{aligned}
\Psi(\alpha_{vi}^t)^{(n)} &= \langle \log t_{vi} \rangle^{(n+1)} - \log \beta_{vi}^t)^{(n)} \\
\Psi(\alpha_{i\tau}^v)^{(n)} &= \langle \log v_{i\tau} \rangle^{(n+1)} - \log \beta_{i\tau}^v)^{(n)} \\
\alpha_{vi}^t)^{(n)} &= a_{vi}^t + \sum_{\tau} \langle s_{v\tau} \rangle^{(n)} \\
\alpha_{i\tau}^v)^{(n)} &= 1 + \sum_v \langle s_{v\tau} \rangle^{(n)}
\end{aligned}$$

can be used to eliminate more costly evaluations of digamma function. Then, the bound adopts the form

$$\begin{aligned}
\mathcal{L}(q)^{(n)} &= - \sum_{v,\tau} \left(\mathbf{E}_t^{(n)} \mathbf{E}_v^{(n)} + \log \Gamma(\mathbf{X} + 1) + \log(\exp \mathbf{L}_t^{(n)} * \exp \mathbf{L}_v^{(n)}) \right) \\
&+ \sum_{v,\tau} \left(-\mathbf{X} * \left((\exp \mathbf{L}_t^{(n)} * \mathbf{L}_t^{(n)}) * \exp \mathbf{L}_v^{(n)} + \exp \mathbf{L}_t^{(n)} * (\exp \mathbf{L}_v^{(n)} * \mathbf{L}_v^{(n)}) / (\exp \mathbf{L}_t^{(n)} * \exp \mathbf{L}_v^{(n)}) \right) \right) \\
&+ \sum_{v,i} \left(-1./\mathbf{B}^t * \mathbf{E}_t^{(n)} - \log \Gamma(\mathbf{A}^t) - \mathbf{A}^t * \log \mathbf{B}^t + \mathbf{A}_t^{(n)} * (\log \mathbf{B}_t^{(n)} + 1) + \log \Gamma(\mathbf{A}_t^{(n)}) \right) \\
&+ \sum_{i,\tau} \left(-\mathbf{E}_\lambda^{(n)} * \mathbf{\Delta}^{(n)T} * \mathbf{E}_v^{(n)} + \log(\mathbf{E}_\lambda^{(n)} * \mathbf{\Delta}^{(n)T}) - \mathbf{A}^v * \log \mathbf{B}^v + \mathbf{A}_v^{(n)} * (\log \mathbf{B}_v^{(n)} + 1) + \log \Gamma(\mathbf{A}_v^{(n)}) \right) \\
&+ \sum_{i,c} \left(-1./\mathbf{B}^\lambda * \mathbf{E}_\lambda^{(n)} + (\mathbf{A}^\lambda - 1) * \mathbf{L}_\lambda^{(n)} - \mathbf{A}^\lambda * \log \mathbf{B}^\lambda - \log \Gamma(\mathbf{A}^\lambda) \right) \\
&+ \sum_{i,c} \left((1 - \mathbf{A}_\lambda^{(n)}) * \psi(\mathbf{A}_\lambda^{(n)}) + \log \mathbf{B}_\lambda^{(n)} + \mathbf{A}_\lambda^{(n)} + \log \Gamma(\mathbf{A}_\lambda^{(n)}) \right) \\
&+ \sum_{\tau,c} (\mathbf{U} - 1) * \mathbf{\Pi}^{(n)T} + \sum_{\tau} (\mathbf{U} * \vec{1} - \sum_c 1) * \psi(\mathbf{U} * \vec{1}) - \sum_{\tau,c} (\mathbf{U} - 1) * \psi(\mathbf{U}) \quad , \quad (14)
\end{aligned}$$

where by $\sum_c 1$ the number of columns of \mathbf{U} is denoted.

Appendix C - Summary of the Learning Algorithm

Inputs: $X, A^t, B^t, A^\lambda, B^\lambda, U$

Initialize (randomly):

$$\mathbf{E}_t^{(0)}, \mathbf{L}_t^{(0)}, \mathbf{E}_v^{(0)}, \mathbf{L}_v^{(0)}, \boldsymbol{\Sigma}_t^{(0)}, \boldsymbol{\Sigma}_v^{(0)}, \boldsymbol{\Delta}^{(0)}, \mathbf{E}_\lambda^{(0)}, \mathbf{L}_\lambda^{(0)}, \boldsymbol{\Pi}^{(0)}, \mathbf{Y}^{(0)}$$

$$j = [j_1, \dots, j_i, \dots]^T, j_i = i$$

$n = 0$;

Repeat:

$$\boldsymbol{\xi}^{(n+1)} = \mathbf{X} ./ ((\exp \mathbf{L}_t^{(n)}) * (\exp \mathbf{L}_v^{(n)}))$$

$$\boldsymbol{\Sigma}_v^{(n+1)} = \exp \mathbf{L}_v^{(n)} .* (\exp \mathbf{L}_t^{(n)})^T * \boldsymbol{\xi}^{(n)}$$

$$\boldsymbol{\Sigma}_t^{(n+1)} = \exp \mathbf{L}_t^{(n)} .* (\boldsymbol{\xi}^{(n)} * (\exp \mathbf{L}_v^{(n)})^T)$$

$$\mathbf{A}_t^{(n)} = \mathbf{1} .* \mathbf{A}^t + \boldsymbol{\Sigma}_t^{(n)}$$

$$\mathbf{B}_t^{(n)} = \mathbf{1} ./ (\mathbf{1} ./ \mathbf{B}^t + \mathbf{1} * \mathbf{E}_v^{(n)T})$$

$$\mathbf{E}_t^{(n+1)} = \mathbf{A}_t^{(n)} .* \mathbf{B}_t^{(n)}$$

$$\mathbf{A}_v^{(n)} = \mathbf{1} + \boldsymbol{\Sigma}_v^{(n)}$$

$$\mathbf{B}_v^{(n)} = \mathbf{1} ./ (\mathbf{E}_\lambda^{(n)} * \boldsymbol{\Delta}^{(n)T} + \mathbf{E}_t^{(n)T} * \mathbf{1})$$

$$\mathbf{E}_v^{(n+1)} = \mathbf{A}_v^{(n)} .* \mathbf{B}_v^{(n)}$$

Optional: calculate bound according to (14)

$$\mathbf{L}_t^{(n+1)} = \Psi(\mathbf{A}_t^{(n)}) + \log \mathbf{B}_t^{(n)}$$

$$\mathbf{L}_v^{(n+1)} = \Psi(\mathbf{A}_v^{(n)}) + \log \mathbf{B}_v^{(n)}$$

$$\mathbf{A}_\lambda^{(n)} = \mathbf{A}^\lambda + \boldsymbol{\Delta}^{(n)} * \mathbf{1}$$

$$\mathbf{B}_\lambda^{(n)} = \mathbf{1} ./ (\mathbf{1} ./ \mathbf{B}^\lambda + \mathbf{E}_v^{(n)} * \boldsymbol{\Delta}^{(n)})$$

$$\mathbf{E}_\lambda^{(n+1)} = \mathbf{A}_\lambda^{(n)} .* \mathbf{B}_\lambda^{(n)}$$

$$\mathbf{L}_\lambda^{(n+1)} = \Psi(\mathbf{A}_\lambda^{(n)}) + \log \mathbf{B}_\lambda^{(n)}$$

$$\mathbf{P}^{(n)} = \boldsymbol{\Pi}^{(n)} + \mathbf{E}_t^{(n)T} * \mathbf{E}_\lambda^{(n)} + \mathbf{1} * \mathbf{L}_\lambda^{(n)}$$

$$\boldsymbol{\Delta}^{(n+1)} = \mathbf{P}^{(n)} ./ (\mathbf{P}^{(n)} * \mathbf{1})$$

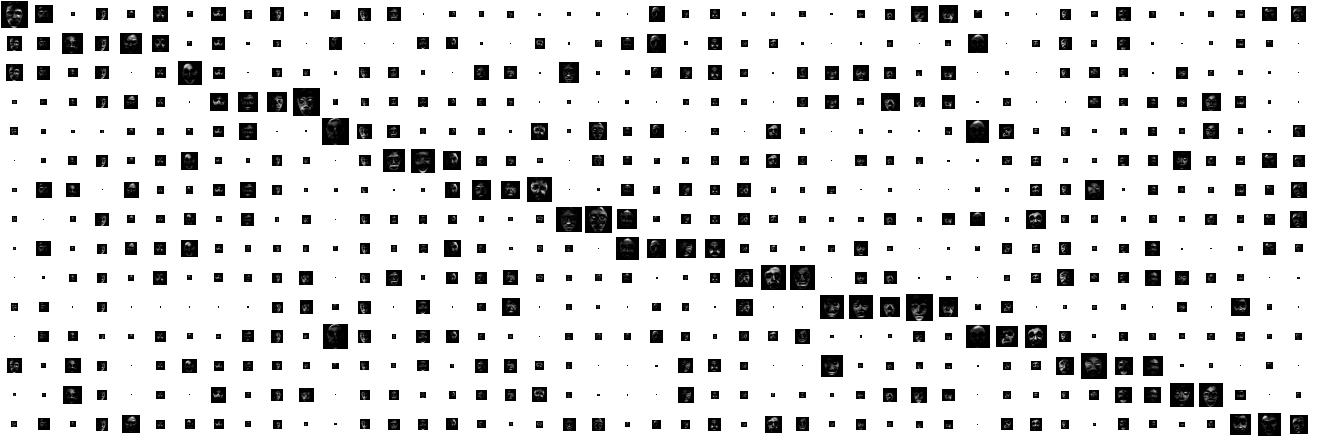
$$\mathbf{Y}^{(n)} = \mathbf{U}^{(n)} + \boldsymbol{\Delta}^{(n)}$$

$$\boldsymbol{\Pi}^{(n+1)} = \psi(\mathbf{Y}^{(n)}) - \psi(\mathbf{Y}^{(n)} * \mathbf{1})$$

$$n = n + 1$$

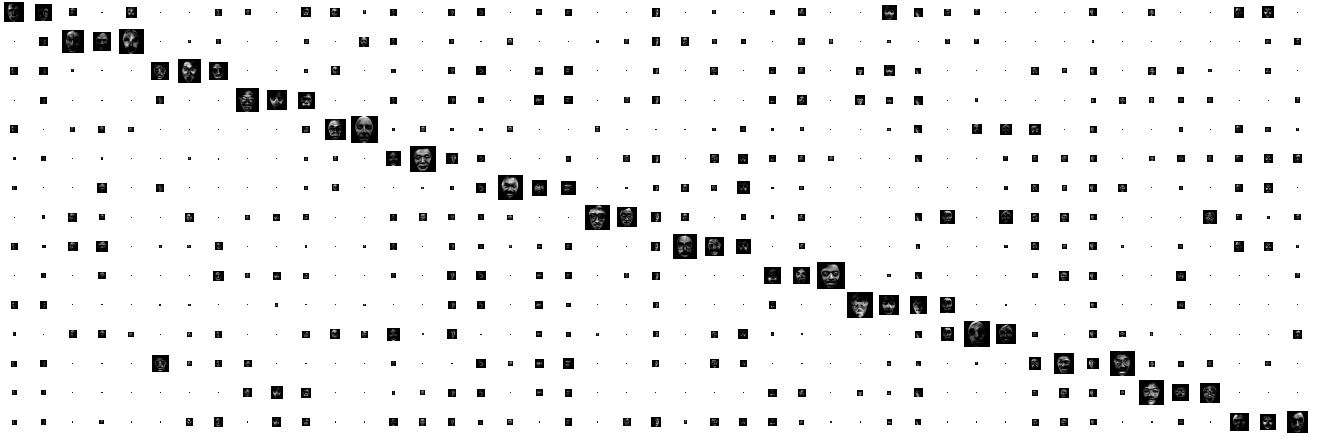
until termination criterion not satisfied.

Appendix D - Additional Figures



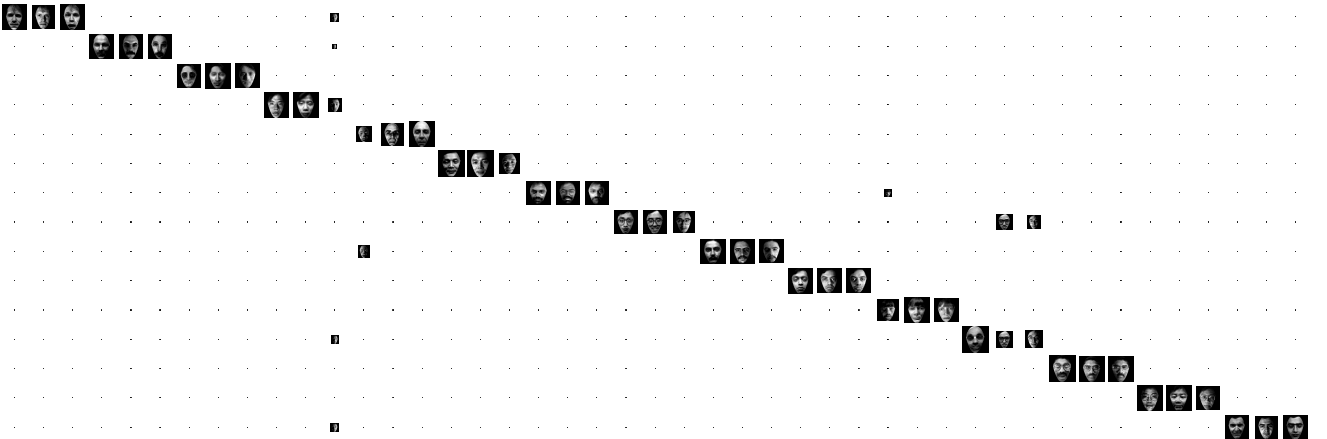
$$nmf_gs; a_{vi}^t = 0.6, b_{vi}^t = 20, a_{small}^\lambda = 32, a_{large}^\lambda = 32, b^\lambda = 1E6$$

a)



$$nmf_gs; a_{vi}^t = 0.6, b_{vi}^t = 20, a_{small}^\lambda = 32, a_{large}^\lambda = 256, b^\lambda = 1E6$$

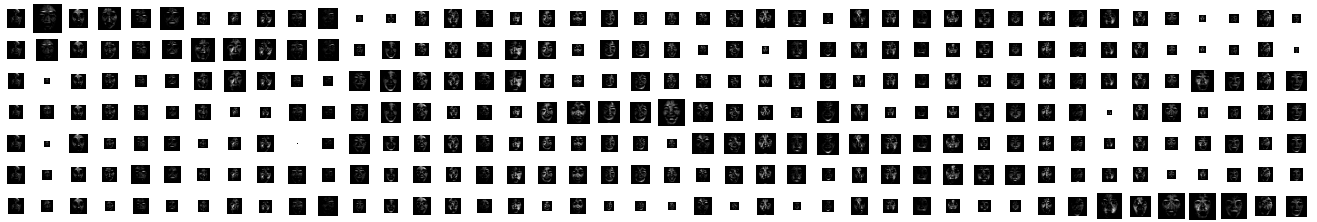
b)



$$nmf_gs; a_{vi}^t = 0.6, b_{vi}^t = 20, a_{small}^\lambda = 32, a_{large}^\lambda = 2048, b^\lambda = 1E6$$

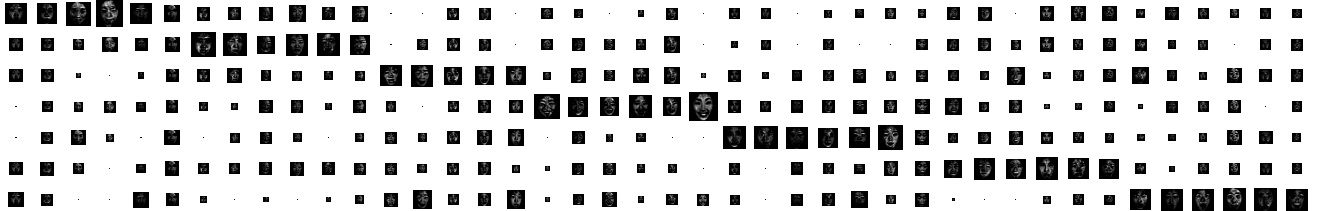
c)

Fig. 11. Controlling prevalences of features in labels on Yale dataset - Hinton diagrams of normalized l_1 norm of the coefficient matrix E_v accumulated across labels, with corresponding features overlaid, for nmf_gs with increasing parameter a_{large}^λ ; a) $a_{large}^\lambda = 32$, b) $a_{large}^\lambda = 256$, c) $a_{large}^\lambda = 2048$. Rows correspond to labels and columns to features.



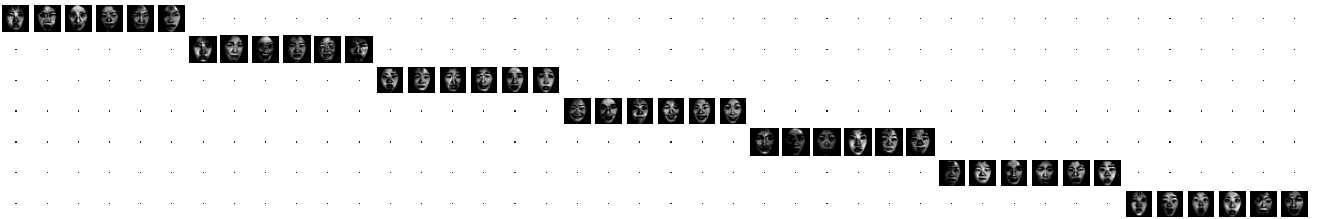
$$nmf_gs; a_{vi}^t = 0.6, b_{vi}^t = 20, a_{small}^\lambda = 32, a_{large}^\lambda = 32, b^\lambda = 1E6$$

a)



$$nmf_gs; a_{vi}^t = 0.6, b_{vi}^t = 20, a_{small}^\lambda = 32, a_{large}^\lambda = 256, b^\lambda = 1E6$$

b)



$$nmf_gs; a_{vi}^t = 0.6, b_{vi}^t = 20, a_{small}^\lambda = 32, a_{large}^\lambda = 2048, b^\lambda = 1E6$$

c)

Fig. 12. Controlling prevalences of features in labels on JAFFE dataset - Hinton diagrams of normalized l1 norm of the coefficient matrix E_v accumulated across labels, with corresponding features overlaid, for nmf_gs with increasing parameter a_{large}^λ : a) $a_{large}^\lambda = 32$, b) $a_{large}^\lambda = 256$, c) $a_{large}^\lambda = 2048$. Rows correspond to labels and columns to features.