

---

# A Deep Semi-NMF Model for Learning Hidden Representations

---

George Trigeorgis  
Konstantinos Bousmalis  
Stefanos Zafeiriou  
Björn W. Schuller

GEORGE.TRIGEORGIS08@IMPERIAL.AC.UK  
K.BOUSMALIS@IMPERIAL.AC.UK  
S.ZAFEIRIOU@IMPERIAL.AC.UK  
BJOERN.SCHULLER@IMPERIAL.AC.UK

Department of Computing, Imperial College London, United Kingdom

## Abstract

Semi-NMF is a matrix factorization technique that learns a low-dimensional representation of a dataset that lends itself to a clustering interpretation. It is possible that the mapping between this new representation and our original features contains rather complex hierarchical information with implicit lower-level hidden attributes, that classical one level clustering methodologies can not interpret. In this work we propose a novel model, Deep Semi-NMF, that is able to learn such hidden representations that allow themselves to an interpretation of clustering according to different, unknown attributes of a given dataset. We show that by doing so, our model is able to learn low-dimensional representations that are better suited for clustering, outperforming Semi-NMF, but also other NMF variants.

## 1. Introduction

Matrix factorization is a particularly useful family of techniques in data analysis. In recent years, there has been a significant amount of research on factorization methods that focus on particular characteristics of both the data matrix and the resulting factors. Non-negative matrix factorization (NMF), for example, focuses on the decomposition of non-negative multivariate data matrix  $X$  into factors  $Z$  and  $H$  that are also non-negative, such that  $X \approx ZH$ . The application area of the family of NMF algorithms has grown significantly during the past years. It has been shown that they can be a successful dimensionality reduction technique over a variety of areas including, but not limited to, environmetrics (Paatero & Tapper, 1994), microarray data analysis (Brunet et al., 2004; Devarajan, 2008), document clustering (Berry & Browne, 2005), face recog-

nition (Zafeiriou et al., 2006; Kotsia et al., 2007) and more. What makes NMF algorithms particularly attractive is the non-negativity constraints imposed on the factors they produce, allowing for better interpretability. Moreover, it has been shown that NMF variants (such as the Semi-NMF) are equivalent to  $k$ -means clustering, and that in fact, NMF variants are expected to perform better than  $k$ -means clustering particularly when the data is not distributed in a spherical manner (Ding et al., 2010; Cing et al., 2005). Nonlinear extensions of NMF have been also recently studied (Zafeiriou & Petrou, 2010).

In order to extend the applicability of NMF in cases where our data matrix  $H$  is not strictly non-negative, Ding et al. (2010) introduced the Semi-NMF, an NMF variant that imposes non-negativity constraints only on the second factor  $H$ , but allows mixed signs in both the data matrix  $X$  and the first factor  $Z$ . This was motivated from a clustering perspective, where  $Z$  represents cluster centroids, and  $H$  represents soft membership indicators for every data point, allowing Semi-NMF to learn new lower-dimensional features from the data that have a convenient clustering interpretation.

It is possible that the mapping  $Z$  between this new representation  $H$  and our original features  $X$  contains rather complex hierarchical and structural information. Consider for example the problem of mapping images of faces to their identities: a face image also contains information about attributes like pose and expression that can help identify the person depicted. One could argue that by further factorizing this mapping  $Z$ , in a way that each factor adds an extra layer of abstraction minimizing the dimensionality of the representation, one could automatically learn such latent attributes and the intermediate hidden representations that are implied, allowing for a better higher-level feature representation  $H$ , as demonstrated in Figure 1. In this work, we propose a novel Deep Semi-NMF approach, which applies the concept of Semi-NMF to a multi-layer structure that is able to learn hidden representations of the original data. As Semi-NMF has a close relation to  $k$ -

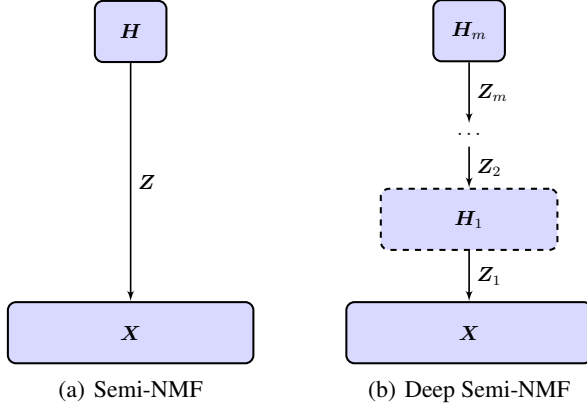


Figure 1. (a) A Semi-NMF model results in a linear transformation of the initial input space. (b) Deep Semi-NMF learns a hierarchy of hidden representations that aid in uncovering the final lower-dimensional representation of the data.

means clustering, by extending the model to a deep one we allow our model to learn new representations of our original data that continue to have a clustering interpretation according to the different latent attributes of our dataset, as demonstrated in Figure 2.

Closest to our proposal is recent work that has presented NMF-variants that factorize  $X$  into more than 2 factors. Specifically, Ahn et al. (2004) have demonstrated the concept of Multi-Layer NMF on a set of facial images and (Lyu & Wang, 2013; Cichocki & Zdunek, 2006; Song & Lee, 2013) have proposed similar NMF models that can be used for Blind Source Separation, classification of digit images (MNIST), and documents. The representations of the Multi-layer NMF however do not lend themselves to a clustering interpretation, as the representations learned from our model. Although the Multi-layer NMF is a promising technique for learning hierarchies of features from data, we show in this work that our proposed model, the Deep Semi-NMF outperforms the Multi-layer NMF and, in fact, all models we compared it with on the task of feature learning for clustering images of faces.

The novelty of this work can be summarized as follows: (1) we outline a novel deep framework for matrix factorization suitable for clustering of multimodally distributed objects such as faces, (2) we present a greedy algorithm to optimize the factors of the Semi-NMF problem, inspired by recent advances in deep learning (Hinton & Salakhutdinov, 2006), and (3) we evaluate the representations learned by different NMF-variants in terms of clustering performance.

## 2. Background

In this work, we assume that our data is provided in a matrix form  $X \in \mathbb{R}^{p \times n}$ , i.e.,  $X = [x_1, x_2, \dots, x_n]$  is a collection of  $n$  data vectors as columns, each with  $p$  features. Matrix factorization aims at finding factors of  $X$  that satisfy certain constraints. In Singular Value Decomposition (SVD) (Golub & Reinsch, 1970), the method that underlies Principal Component Analysis (PCA) (Wold et al., 1987), we factorize  $X$  into two factors: the loadings or bases  $Z \in \mathbb{R}^{p \times k}$  and the features or components  $H \in \mathbb{R}^{k \times n}$ , without imposing any sign restrictions on either our data or the resulting factors. In Non-negative Matrix Factorization (NMF) (Lee & Seung, 2001) we assume that all matrices involved contain only non-negative elements<sup>1</sup>, so we try to approximate a factorization  $X^+ \approx Z^+ H^+$ .

### 2.1. Semi-NMF

In turn, Semi-NMF (Ding et al., 2010) relaxes the non-negativity constraints of NMF and allows the data matrix  $X$  and the loadings matrix  $Z$  to have mixed signs, while restricting only the features matrix  $H$  to comprise of strictly non-negative components, thus approximating the following factorization:

$$X^\pm \approx Z^\pm H^+. \quad (1)$$

This is motivated from a clustering perspective. If we view  $Z = [z_1, z_2, \dots, z_k]$  as the cluster centroids, then  $H = [h_1, h_2, \dots, h_n]$  can be viewed as the cluster indicators for each datapoint. In fact, if we had a matrix  $H$  that was not only non-negative but also orthogonal, then every column vector would have only one positive element, making Semi-NMF equivalent to  $k$ -means. Thus Semi-NMF, which does not impose an orthogonality constraint on its features matrix, can be seen as a soft clustering method where the features matrix describes the compatibility of each component with a cluster centroid, a base in  $Z$ .

### 2.2. State-of-the-art for learning features for clustering based on NMF-variants

In this work, we compare our method with, among others, the state-of-the-art NMF techniques for learning features for the purpose of clustering. Cai et al. (2011) proposed a graph-regularized NMF (GNMF) which takes into account the intrinsic geometric and discriminating structure of the data space, which is essential to the real-world applications, especially in the area of clustering. To accomplish this, GNMF constructs a nearest neighbor graph to model the manifold structure. By preserving the graph structure,

<sup>1</sup>When not clear from the context we will use the notation  $A^+$  to state that a matrix  $A$  contains only non-negative elements. Similarly, when not clear, we will use the notation  $A^\pm$  to state that  $A$  may contain any real number.

it allows the learned features to have more discriminating power than the standard NMF algorithm, in cases that the data are sampled from a submanifold which lies in a higher dimensional ambient space. This combines two different notions, that of NMF and graph Laplacian regularization algorithms (Belkin & Niyogi, 2001).

Another state-of-the-art matrix factorization technique would be NeNMF (Guan et al., 2012). NeNMF makes use of Nesterov’s optimal gradient method to alternatively optimize one factor, with the other fixed. This allows for faster NMF optimization without time-consuming linesearch procedure or numerical instability problems that traditional NMF has. Guan et al. (2012) showed that it outperformed existing NMF solvers in terms of reconstruction error and document clustering performance.

### 3. Deep Semi-NMF

In Semi-NMF the goal is to construct a low-dimensional representation  $H^+$  of our original data  $X^\pm$ , with the bases matrix  $Z^\pm$  serving as the mapping between our original data and its lower-dimensional representation (see Equation 1). In many cases the data we wish to analyze is often rather complex and has a collection of distinct, often unknown, attributes. In this work for example, we deal with datasets of human faces, where the variability in the data does not only stem from the difference in the appearance of the subjects, but also from other attributes, such as the pose of the head in relation to the camera, or the facial expression of the subject. In addition faces compromise of mainly hierarchical features and thus face clustering problems can be better solved using our deep framework, as each subsequent layers can capture the hierarchical structure.

We propose here the Deep Semi-NMF model, which factorizes a given data matrix  $X$  into  $m + 1$  factors, as follows:

$$X^\pm \approx Z_1^\pm Z_2^\pm \dots Z_m^\pm H_m^+ \quad (2)$$

This formulation, as shown in with respect to Figures 2 and 1 allows for a hierarchy of  $m$  layers of implicit representations of our data that can be given by the following factorizations:

$$\begin{aligned} H_{m-1}^+ &\approx Z_m^\pm H_m^+ \\ &\vdots \\ H_2^+ &\approx Z_3^\pm \dots Z_m^\pm H_m^+ \\ H_1^+ &\approx Z_2^\pm \dots Z_m^\pm H_m^+ \end{aligned}$$

As one can see above, we further restrict these implicit representations ( $H_1^+, \dots, H_{m-1}^+$ ) to also be non-negative. By doing so, every layer of this hierarchy of representations also lends itself to a clustering interpretation. By examining Figure 2 one can acquire better intuition of how that

happens. In this case the input to the model, our data  $X$ , is a collection of face images from different subjects (identity), expressing a variety of emotions (expressions) taken from many angles (pose). A Semi-NMF model would find a representation  $H$  of  $X$ , which would be useful for performing clustering according to the identity of the subjects, and  $Z$  the mapping between these identities and the face images. A Deep Semi-NMF model also finds a representation of our data that has a similar interpretation at the top layer, its last factor  $H_m$ . However, the mapping from identities to face images is now further analyzed as a product of three factors  $Z = Z_1 Z_2 Z_3$ , with  $Z_3$  corresponding to the mapping of identities to emotions,  $Z_2 Z_3$  corresponding to the mapping of identities to poses, and finally  $Z_1 Z_2 Z_3$  corresponding to the mapping of identities to the face images. That means that, as shown in Figure 2 we are able to decompose our data in 3 different ways according to our 3 different attributes:

$$\begin{aligned} X^\pm &\approx Z_1^\pm H_1^+ \\ X^\pm &\approx Z_1^\pm Z_2^\pm H_2^+ \\ X^\pm &\approx Z_1^\pm Z_2^\pm Z_3^\pm H_3^+ \end{aligned}$$

Our hypothesis is that by further factorizing  $Z$  we are able to construct a deep model that is able to (1) automatically learn what this latent hierarchy of attributes is; (2) find representations of the data that are most suitable for clustering according to the attribute that corresponds to each layer in the model; and (3) find a better high-level, final-layer representation for clustering according to the attribute with the lowest variability, in our case the identity of the face depicted. In our example in Figure 2 we would expect to find better features for clustering according to identities  $H_3$  by learning the hidden representations at each layer most suitable for each of the attributes in our data, in this example:  $H_1 \approx Z_2 Z_3 H_3$  for clustering our original images in terms of poses and  $H_2 \approx Z_3 H_3$  for clustering the face images in terms of expressions.

In order to expedite the approximation of the factors in our model, we pre-train each of the layers to have an initial approximation of the matrices  $Z_i, H_i$  as this greatly improves the training time of the model. This is a tactic that has been employed successfully before (Hinton & Salakhutdinov, 2006) on deep autoencoder networks. To perform the pre-training, we first decompose the initial data matrix  $X \approx Z_1 H_1$ , where  $Z_1 \in \mathbb{R}^{p \times k_1}$  and  $H_1 \in \mathbb{R}_0^{+k_1 \times n}$ . Following this, we decompose the features matrix  $H_1 \approx Z_2 H_2$ , where  $Z_2 \in \mathbb{R}^{k_1 \times k_2}$  and  $H_1 \in \mathbb{R}_0^{+k_2 \times n}$ , continuing to do so until we have pre-trained all of the layers. Afterwards, we can fine-tune the weights of each layer, by employing alternating minimization (with respect to the objective function in Equation 3) of the two factors in each layer, in order to reduce the total reconstruction error of the model, according to the cost

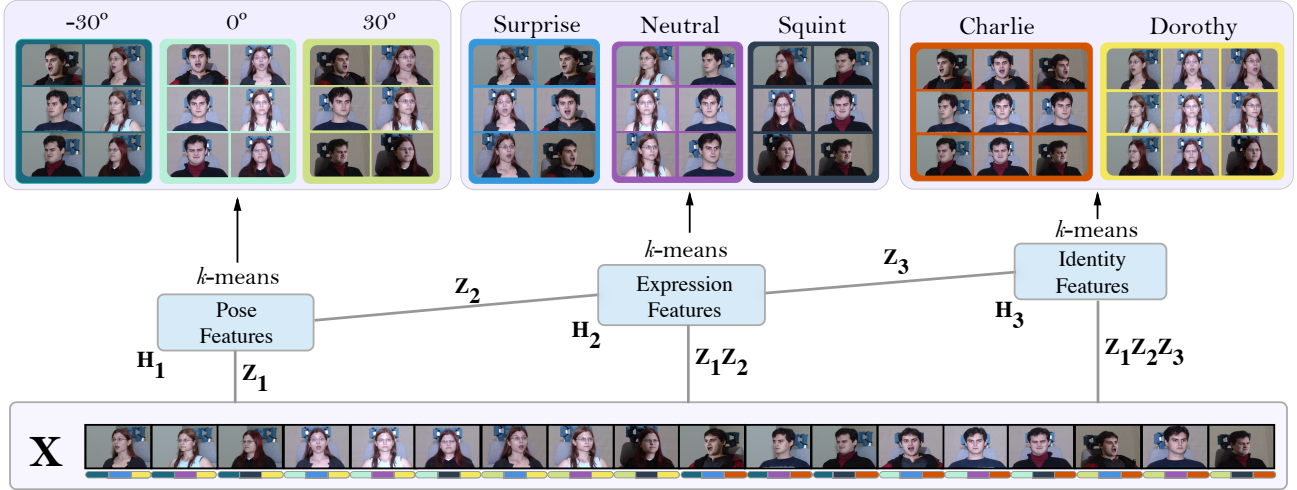


Figure 2. A Deep Semi-NMF model learns a hierarchical structure of features, with each layer learning a representation suitable for clustering according to the different attributes of our data. In this simplified, for demonstration purposes, example from the CMU Multi-PIE database, a Deep Semi-NMF model is able to simultaneously learn features for pose clustering ( $H_1$ ), for expression clustering ( $H_2$ ), and for identity clustering ( $H_3$ ). Each of the images in  $X$  has an associated color coding that indicates its memberships according to each of these attributes (pose/expression/identity).

function in Equation 3.

$$\begin{aligned} C_{\text{deep}} &= \frac{1}{2} \|X - Z_1 Z_2 \cdots Z_m H_m\|_F^2 \\ &= \text{tr}[X^\top X - 2X^\top Z_1 Z_2 \cdots Z_m H_m \\ &\quad + H_m^\top Z_m^\top Z_{m-1}^\top \cdots Z_1^\top Z_1 Z_2 \cdots Z_m H_m] \quad (3) \end{aligned}$$

**Update rule for the weights matrix  $Z$**  We fix the rest of the weights for the  $i^{\text{th}}$  layer and we minimize the cost function with respect to  $Z_i$ . That is, we set  $\partial C_{\text{deep}} / \partial Z_i = 0$ , which gives us the updates:

$$\begin{aligned} Z_i &= (\Psi^\top \Psi)^{-1} \Psi^\top X \tilde{H}_i^\top (\tilde{H}_i \tilde{H}_i^\top)^{-1} \\ Z_i &= \Psi^\dagger X \tilde{H}_i^\dagger \quad (4) \end{aligned}$$

where  $\Psi = Z_1 \cdots Z_{i-1}$ ,  $\dagger$  denotes the Moore-Penrose pseudo-inverse and  $\tilde{H}_i$  is the reconstruction of the  $i^{\text{th}}$  layer's feature matrix.

**Update rule for features matrix  $H$**  Utilizing a similar proof to Ding et al. (2010), we can formulate the update rule for  $H_i$  as follows:

$$H_i = H_i \odot \sqrt{\frac{[\Psi^\top X]_{\text{pos}} + [\Psi^\top \Psi]_{\text{neg}} H_i}{[\Psi^\top X]_{\text{neg}} + [\Psi^\top \Psi]_{\text{pos}} H_i}} \quad (5)$$

Supplementary material including the implementation of Algorithm 1 and the proof of its convergence can be found at <http://trigeorgis.com/deepseminmf>.

## Complexity

The computational complexity for the pre-training stage of Deep Semi-NMF is of order  $\mathcal{O}(mt(pnk + nk^2 + kp^2 + kn^2))$ , where  $m$  is the number of layers,  $t$  the number of iterations until convergence and  $k$  is the maximum number of components out of all the layers. The complexity for the fine-tuning stage is  $\mathcal{O}(mt_f(pnk + (p+n)k^2))$  where  $t_f$  is the number of additional iterations needed.

## Non-linear Update Rules

One can use a non-linear function  $g(\cdot)$ , between each of the implicit representations  $(H_1^+, \dots, H_{m-1}^+)$ , in order to better approximate the non-linear manifolds which the given data matrix  $X$  originally lies on. Thus, one can represent the  $i^{\text{th}}$  feature matrix  $H_i$ , by

$$g(H_i) \approx Z_{i+1} H_{i+1}. \quad (6)$$

This, creates new derivatives update rules as the derivatives with respect to the objective function become:

$$\frac{\partial C_{\text{deep}}}{\partial Z_i} = (N_i - R_i) H_i^\top \quad (7)$$

$$\frac{\partial C_{\text{deep}}}{\partial H_i} = Z_i^\top (N_i - R_i) \quad (8)$$



**Algorithm 1** Suggested algorithm for training a Deep Semi-NMF model. Initially we approximate the factors greedily using the SEMI-NMF algorithm (Ding et al., 2010) and we fine-tune the factors until we reach the convergence criterion.

**function** DEEPSMINMF

**Input:**  $X \in \mathbb{R}^{p \times n}$ , list of layer sizes

**Output:** weight matrices  $Z_i$  and feature matrices  $H_i$  for each of the layers

Initialize Layers

**for all** layers **do**

$Z_i, H_i \leftarrow \text{SEMINMF}(H_{i-1}, \text{layers}(i))$

**end for**

**repeat**

**for all** layers **do**

$\tilde{H}_i \leftarrow \begin{cases} H_i & \text{if } i = k \\ Z_{i+1} \tilde{H}_{i+1} & \text{otherwise} \end{cases}$

$\Psi \leftarrow \prod_{k=1}^{i-1} Z_k$

$Z_i \leftarrow \Psi^\dagger X \tilde{H}_i^\dagger$

$H_i \leftarrow H_i \odot \sqrt{\frac{[\Psi^\top X]^\text{pos} + [\Psi^\top \Psi]^\text{pos} H_i}{[\Psi^\top X]^\text{neg} + [\Psi^\top \Psi]^\text{pos} H_i + \epsilon}}$

**end for**

**until** Stopping criterion is reached

**end function**

which gives the following update rules:

$$Z_i = \begin{cases} X \tilde{H}_i^\dagger; & \text{if } i = 1 \\ g(H_{i-1}) \tilde{H}_i^\dagger & \text{otherwise} \end{cases}, \quad (9)$$

$$H_i = H_i \odot \frac{[Z_i^\top R_i]^\text{pos} + [Z_i^\top N_i]^\text{neg}}{[Z_i^\top R_i]^\text{neg} + [Z_i^\top N_i]^\text{pos}} \quad (10)$$

where,  $\forall i. i \leq m$   $N_i, R_i$  are auxiliary matrices that facilitate the computation of the factors in each layer,  $A^\text{pos}$  is a matrix that has the negative elements of matrix  $A$  replaced with 0, and similarly  $A^\text{neg}$  is one that has the positive elements of  $A$  replaced with 0. The base case values  $i = 1$  for these matrices are:

$$N_1 = (Z_1 \tilde{H}_1) \odot \nabla g^{-1}(Z_1 H_1), \quad (11)$$

$$R_1 = X \odot \nabla g^{-1}(Z_1 H_1), \quad (12)$$

and for the cases where  $i > 1$  they are computed as such:

$$N_{i+1} = (Z_i^\top N_i) \odot \nabla g^{-1}(Z_i H_i), \quad (13)$$

$$R_{i+1} = (Z_i^\top R_i) \odot \nabla g^{-1}(Z_i H_i). \quad (14)$$

## 4. Experiments

Our main hypothesis is that a Deep Semi-NMF is able to learn better high-level representations of our original data

than an one-layer Semi-NMF for clustering according to the attribute with the lowest variability in the dataset. In order to evaluate this hypothesis, we have compared the performance of Deep Semi-NMF with that of other methods, on the task of clustering images of faces in 3 distinct datasets. These datasets are:

- **CMU MultiPIE:** The first dataset we examine is *The CMU Multi Pose, Illumination, and Expression* (MultiPIE) Database (Gross et al., 2010) and contains around 750,000 images of 337 subjects, captured under laboratory conditions in four different sessions. In this work, we used a subset of 13,230 images of 147 subjects in 5 different poses and 3 different illumination conditions, expressing 6 different emotions. Using the annotations from (Sagonas et al., 2013a;b), we aligned these images based on a common frame by using piece-wise affine warping. After that, we resized them to a smaller resolution of  $30 \times 30$ . The database comes with labels for each of the attributes mentioned above: identity, illumination, pose, expression.
- **CMU PIE:** We also used a freely available version of CMU Pie (Sim et al., 2003), which comprises of 2,856 grayscale  $32 \times 32$  face images of 68 subjects. Each person has 42 facial images under different light and illumination conditions. In this database we only know the identity of the face in each image and we could not use piece-wise affine warping, since we did not have facial point annotations for it.
- **XM2VTS:** *The Extended Multi Modal Verification for Teleservices and Security applications* (XM2VTS) (Messer et al., 1999) contains 2,360 frontal images of 295 different subjects. Each subject has two available images for each of the four different laboratory sessions, for a total of 8 images. The images were aligned based on the same piece-wise affine warping technique as the Multi-PIE dataset, after resizing the original images to  $42 \times 30$ .

In order to evaluate the performance of our model, we compared it against not only Semi-NMF (Ding et al., 2010), but also against other NMF variants that could be useful in learning such representations. More specifically, for each of our three datasets we performed the following experiments:

- **Pixel Intensities:** By using only the pixel intensities of the images in each of our datasets, which of course give us a strictly non-negative input data matrix  $X$ , we compare the reconstruction error and the clustering performance of our Deep Semi-NMF method against the Semi-NMF, NMF with multiplicative update rules (Lee & Seung, 2001), Multi-Layer NMF (Song & Lee,

Name	Number of components				
	30	40	50	60	70
CMU Multi-PIE – Pixel Intensities					
NMF	2.68	2.52	2.41	2.32	2.23
NeNMF	2.40	2.20	2.04	1.92	1.80
GNMF	3.05	3.03	2.98	3.02	3.02
Semi-NMF	2.37	2.16	2.01	1.89	1.78
Multi-layer NMF	2.71	2.54	2.42	2.31	2.22
Deep Semi-NMF	2.40	2.20	2.07	1.96	1.86
CMU Multi-PIE – Image Gradient Orientations					
Semi-NMF	0.14	0.13	0.13	0.12	0.12
Deep Semi-NMF	0.14	0.13	0.13	0.13	0.12

Table 1. The mean reconstruction error for each of the algorithms with a variable number of components. The error quantifies the deviation of the reconstruction with the original matrix  $\mathbf{X}$ , using the Euclidean objective function we used:  $\frac{1}{N} \|\mathbf{X} - \hat{\mathbf{X}}\|^2$ , with  $N$  amount of samples. Deep Semi-NMF has a comparable reconstruction error to that of Semi-NMF.

2013), GNMF (Cai et al., 2011), and NeNMF (Guan et al., 2012).

- **Image Gradient Orientations (IGO):** In general, the trend in Computer Vision is to use complicated engineered features like HoGs, SIFT, LBPs, etc. As a proof of concept, we choose to conduct experiments with simple gradient orientations (Zafeiriou et al., 2012) as features, instead of pixel intensities, which results into a data matrix of mixed signs, and expect that we can learn better data representations for clustering faces according to identities. In this case, we only compared our Deep Semi-NMF with its one-layer Semi-NMF equivalent, as the other techniques are not able to deal with mixed-sign matrices.

Finally, we evaluated our secondary hypotheses, i.e. that every hidden representation in each layer is in fact most suited for clustering according to the attributes that corresponds to the layer of interest. We performed clustering experiments by using the features learned in each layer of a two-layer Deep Semi-NMF on the case of CMU Multi-PIE, as this was the only dataset for which we had labels for attributes other than the identity.

#### 4.1. Implementation Details

In order to initiate the matrix factorization process, NMF and Semi-NMF algorithms start from some initial point  $(\mathbf{Z}^0, \mathbf{H}^0)$ , where usually  $\mathbf{Z}^0$  and  $\mathbf{H}^0$  are randomly initialized matrices. In order to speed up the convergence rate of NMF, Boutsidis & Gallopoulos (2008), suggested Non-negative Double Singular Value decomposition (NNDSVD), which is a new method based on two SVD

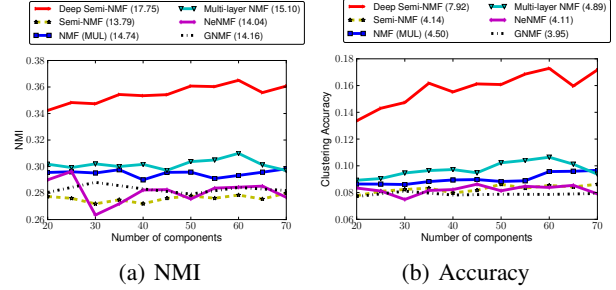


Figure 3. CMU MultiPIE–Pixel Intensities: NMI and Accuracy for clustering based on the representations learned by each model with respect to identities. The deep architectures are comprised of 2 representation layers 1988-625-a and the representations used were from the top layer. In parenthesis we show the AUC scores.

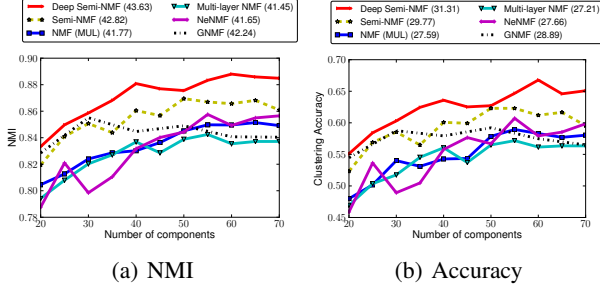
processes, one to approximate the initial data matrix  $\mathbf{X}$  and the other to approximate the positive sections of the resulting partial SVD factors. Although, the proposed initialization of Semi-NMF by its authors is by using the  $k$ -means algorithm (Ding et al., 2010), we found empirically that it was too computationally heavy when the number of components  $k$  was fairly high ( $k > 100$ ). As an alternative we used NNDSVD to gain an initial and deterministic approximation, after forcing the initial data to have non-negative values, by setting all negative values to zero.

For the GNMF experimental setup, we chose a suitable number of neighbours, in our case 5, using visualization of the datasets using Laplacian Eigenmaps (Belkin & Niyogi, 2001), such that we had visually distinct clusters.

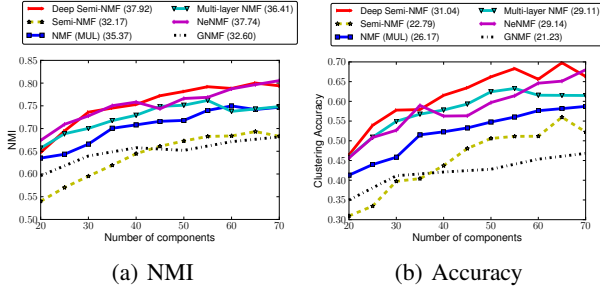
Important for the experimental setup is the selected structure of the multi-layered models. After careful preliminary experimentation, we focused on experiments that involve two hidden layer architectures for the Deep Semi-NMF and Multi-layer NMF. Although we experimented with a higher-number of layers, approximating the additional number of factors/layers did not seem to have a significant impact on our results for these datasets, and was significantly more computationally expensive. We specifically experimented with models that had a first hidden representation  $\mathbf{H}_1$  with 625 features, and a second representation  $\mathbf{H}_2$  with a number of features that ranged from 20 to 70. Again we chose these numbers after preliminary results showed us that the computational burden of additional features outweighed the performance increase obtained by increasing these numbers.

#### 4.2. Reconstruction Error Results

Our first experiment was to evaluate whether the extra layers, which naturally introduce more factors and are therefore more difficult to optimize, result in a lower quality lo-



**Figure 4. XM2VTS-Pixel Intensities:** NMI and Accuracy for clustering based on the representations learned by each model with respect to identities. The deep architectures are comprised of 2 representation layers 1260-625-a and the representations used were from the top layer. In parenthesis we show the AUC scores.



**Figure 5. CMU PIE-Pixel Intensities:** NMI and Accuracy for clustering based on the representations learned by each model with respect to identities. The deep architectures are comprised of 2 representation layers 1024-625-a and the representations used were from the top layer. In parenthesis we show the AUC scores.

cal optimum. We evaluated how well the matrix decomposition is performed by calculating the reconstruction error, the Frobenius norm of the difference between the original data and the reconstruction for all the methods we compared. Note that, in order to have comparable results, all of the methods have the same stopping criterion rules. We have set the maximum amount of iterations to 300 (usually  $\sim 100$  epochs are enough) and we use the convergence rule  $E_{i-1} - E_i \leq \kappa \max(1, E_{i-1})$  in order to stop the process when reconstruction error ( $E_i$ ) between the current and previous update is small enough. In our experiments we set  $\kappa = 10^{-6}$ . Table 1 shows the change in reconstruction error with respect to the selected number of features in  $H_2$  for all the methods we used on the Multi-PIE dataset. The results for the other datasets and for a larger variety of number of components were similar and are excluded due to lack of space.

The results show that Semi-NMF and Deep Semi-NMF manage to reach a much lower reconstruction error than

#	Pose Accuracy (%)		Identity Accuracy (%)	
	$H_1$	$H_2$	$H_1$	$H_2$
1	27.58	23.11	9.37	16.59
2	28.18	23.25	9.48	16.67

**Table 2.** Clustering accuracy according to pose and identity labels on the CMU Multi-PIE dataset using two Deep Semi-NMF models with two hidden layers each (1) 625-60 and (2) 625-70.

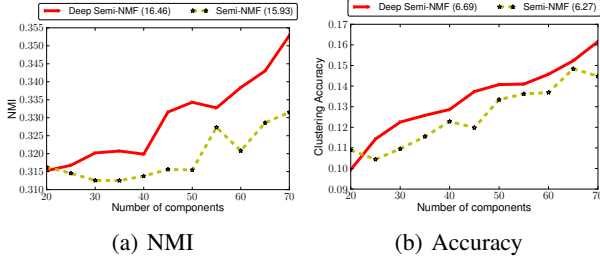
the rest consistently, which would match our expectations as they do not constrain the weights  $Z$  to be non-negative. What is important to note here is that the Deep Semi-NMF models do not have a significantly lower reconstruction error compared to the equivalent Semi-NMF models, even though the approximation involves more factors. This is in contrast to the multi-layer NMF and GNMF which sacrifice the reconstruction quality, in return for uncovering more meaningful features than their NMF counterpart.

### 4.3. Clustering Results

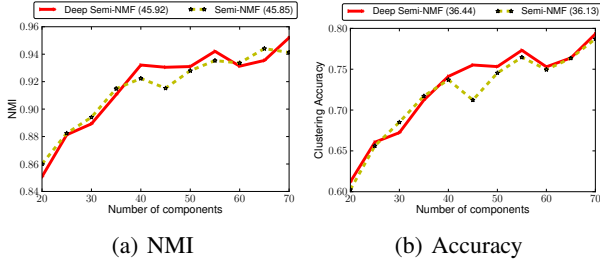
After achieving satisfactory reconstruction error for our method, we proceeded to evaluate the features learned at the final representation layer, by using  $k$ -means clustering, following the experimental protocol in (Cai et al., 2011). To assess the clustering quality of the representations produced by each of the algorithms we compared, we take advantage of the fact that the datasets are already labelled. The two metrics used were the accuracy (AC) and the normalized mutual information metric (NMI), as those are defined in (Xu et al., 2003).

Figures 3-5 show the comparison in clustering accuracy and NMI when using  $k$ -means on the feature representations produced by each of the techniques we compared, when our input matrix contained only the pixel intensities of each image. Our method significantly outperforms every method we compared it with on all three datasets. The difference is not as obvious in Figure 5, which could perhaps be a sign that our method is able to use the warping technique we used to its advantage more so than the other techniques we used.

By using IGOs, the Deep Semi-NMF was able to outperform the single-layer Semi-NMF as shown in Figures 6-8. Making use of these simple mixed-signed features improved the clustering accuracy considerably, especially for XM2VTS and CMU PIE. It should be noted that in all cases, with the exception of the XM2VTS experiment with IGOs, our Deep Semi-NMF outperformed all other methods with a difference in performance that is statistically significant (paired t-test,  $p \ll 0.01$ ).



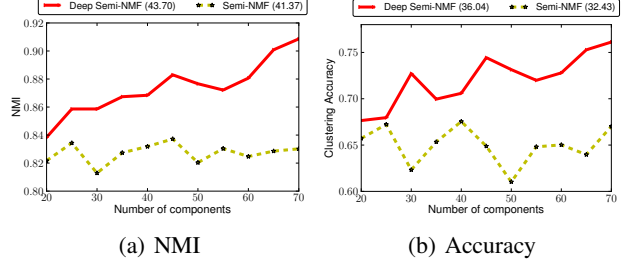
**Figure 6. CMU MultiPIE-IGO:** NMI and Accuracy scores on clustering based on the representations learned by each model with respect to identities. The deep architectures are comprised of 2 hidden layers 1988-625-a and the representations used were from the top layer. In parenthesis we show the AUC scores.



**Figure 7. XM2VTS-IGO:** NMI and Accuracy scores on clustering based on the representations learned by each model with respect to identities. The deep architectures are comprised of 2 hidden layers 2520-625-a and the representations used were from the top layer. In parenthesis we show the AUC scores.

#### 4.4. Clustering with Respect to Different Attributes

Finally, we conducted experiments by performing  $k$ -means clustering on each of the two representations learned by our Deep Semi-NMF models by using the raw intensities of our warped images of CMU Multi-PIE. We only used Multi-PIE since we only had identity labels for our other datasets. Since the warping technique we use gets rid of most of the variability shown in expressions, we evaluated how well we did on clustering according to pose and identities. As one can see in Table 2 our first layer indeed learns representations that are better suited for clustering according to poses. On the other hand, we do confirm that our final layer indeed learns representations that are better suited for clustering according to identities, thus confirming our secondary hypothesis, i.e. that every hidden representation in each layer is in fact most suited for clustering according to the attributes that corresponds to the layer of interest. As one can see in Figure 3, by learning the hidden representation that relates to poses, we are able to achieve significantly better results compared to the Semi-NMF, where we learn only one level of representation.



**Figure 8. CMU PIE-IGO:** NMI and Accuracy scores on clustering based on the representations learned by each model with respect to identities. The deep architectures are comprised of 2 hidden layers 2048-625-a and the representations used were from the top layer. In parenthesis we show the AUC scores.

## 5. Conclusion

We have introduced a novel deep architecture for semi-non-negative matrix factorization, the Deep Semi-NMF, that is able to automatically learn a hierarchy of attributes of a given dataset, as well as representations suited for clustering according to these attributes. We have also presented an algorithm for optimizing the factors of our Deep Semi-NMF, and we evaluate its performance compared to the single-layered Semi-NMF and other related work, on the problem of clustering faces with respect to their identities. We have shown that our technique is able to learn a high-level, final-layer representation for clustering with respect to the attribute with the lowest variability in the case of three popular datasets of face images, outperforming the other NMF-based techniques.

The next obvious step is to explore the suitability of such learned hidden representations for performing classification. Another line of work that will aid the learning process will be the initialization scheme for pre-training the Deep Semi-NMF model, as currently we used a rough approximation of the initial  $Z^0, H^0$  matrices using the NNDSVD algorithm. Finally, future avenues include experimenting with other applications, e.g. in the area of speech recognition, especially for multi-source speech recognition and we will investigate multilinear extensions of the proposed framework (Zafeiriou, 2009b;a).

## Acknowledgements

We would like to thank our anonymous reviewers for their useful comments. George Trigeorgis is a recipient of the fellowship of the Department of Computing, Imperial College London, and this work was partially funded by it. The work of Konstantinos Bousmalis was funded partially from the Google Europe Fellowship in Social Signal Processing. The work of Stefanos Zafeiriou was partially funded by the EPSRC project EP/J017787/1 (4D-FAB).



## References

- Ahn, Jong-Hoon, Choi, Seungjin, and Oh, Jong-Hoon. A multiplicative up-propagation algorithm. In *Proceedings of the twenty-first international conference on Machine learning*, pp. 3. ACM, 2004.
- Belkin, Mikhail and Niyogi, Partha. Laplacian eigenmaps and spectral techniques for embedding and clustering. In *NIPS*, volume 14, pp. 585–591, 2001.
- Berry, Michael W and Browne, Murray. Email surveillance using non-negative matrix factorization. *Computational & Mathematical Organization Theory*, 11(3):249–264, 2005.
- Boutsidis, Christos and Gallopoulos, Efstratios. Svd based initialization: A head start for nonnegative matrix factorization. *Pattern Recognition*, 41(4):1350–1362, 2008.
- Brunet, Jean-Philippe, Tamayo, Pablo, Golub, Todd R, and Mesirov, Jill P. Metagenes and molecular pattern discovery using matrix factorization. *Proceedings of the National Academy of Sciences*, 101(12):4164–4169, 2004.
- Cai, Deng, He, Xiaofei, Han, Jiawei, and Huang, Thomas S. Graph regularized nonnegative matrix factorization for data representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(8):1548–1560, 2011.
- Cichocki, Andrzej and Zdunek, Rafal. Multilayer nonnegative matrix factorization. *Electronics Letters*, 42:947–948, 2006.
- Cing, C., He, X., and Simon, H.D. On the equivalence of non-negative matrix factorization and spectral clustering. In *Proc. SIAM Data Mining*, 2005.
- Devarajan, Karthik. Nonnegative matrix factorization: an analytical and interpretive tool in computational biology. *PLoS computational biology*, 4(7):e1000029, 2008.
- Ding, Chris HQ, Li, Tao, and Jordan, Michael I. Convex and semi-nonnegative matrix factorizations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(1):45–55, 2010.
- Golub, Gene H and Reinsch, Christian. Singular value decomposition and least squares solutions. *Numerische Mathematik*, 14(5):403–420, 1970.
- Gross, Ralph, Matthews, Iain, Cohn, Jeffrey, Kanade, Takeo, and Baker, Simon. Multi-pie. *Image and Vision Computing*, 28(5): 807–813, 2010.
- Guan, Naiyang, Tao, Dacheng, Luo, Zhigang, and Yuan, Bo. Nnmf: an optimal gradient method for nonnegative matrix factorization. *Signal Processing, IEEE Transactions on*, 60(6): 2882–2898, 2012.
- Hinton, Geoffrey E and Salakhutdinov, Ruslan R. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- Kotsia, Irene, Zafeiriou, Stefanos, and Pitas, Ioannis. A novel discriminant non-negative matrix factorization algorithm with applications to facial image characterization problems. *IEEE Transactions on Information Forensics and Security*, 2(3-2): 588–595, 2007.
- Lee, Daniel D. and Seung, H. Sebastian. Algorithms for non-negative matrix factorization. *Advances in neural information processing systems*, 13:556–562, 2001.
- Lyu, Siwei and Wang, Xin. On algorithms for sparse multi-factor nmf. In *Advances in Neural Information Processing Systems*, pp. 602–610, 2013.
- Messer, Kieron, Matas, Jiri, Kittler, Josef, Luetttin, Juergen, and Maitre, Gilbert. Xm2vtsdb: The extended m2vts database. In *Second international conference on audio and video-based biometric person authentication*, volume 964, pp. 965–966. Citeseer, 1999.
- Paatero, Pentti and Tapper, Unto. Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2):111–126, 1994.
- Sagonas, C., Tzimiropoulos, G., Zafeiriou, S., and Pantic, M. 300 faces in-the-wild challenge: The first facial landmark localization challenge. In *Computer Vision Workshops (ICCVW), 2013 IEEE International Conference on*, pp. 397–403, Dec 2013a.
- Sagonas, C., Tzimiropoulos, G., Zafeiriou, S., and Pantic, M. A semi-automatic methodology for facial landmark annotation. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2013 IEEE Conference on*, pp. 896–903, June 2013b.
- Sim, Terence, Baker, Simon, and Bsat, Maan. "the cmu pose, illumination, and expression database." *IEEE Trans. Pattern Analysis and Machine Intelligence*, 25(12):1615–1618, 2003.
- Song, Hyun Ah and Lee, Soo-Young. Hierarchical data representation model - multi-layer nmf. *International Conference on Learning Representations*, abs/1301.6316, 2013.
- Wold, Svante, Esbensen, Kim, and Geladi, Paul. Principal component analysis. *Chemometrics and intelligent laboratory systems*, 2(1):37–52, 1987.
- Xu, Wei, Liu, Xin, and Gong, Yihong. Document clustering based on non-negative matrix factorization. In *Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 267–273. ACM, 2003.
- Zafeiriou, Stefanos. Algorithms for nonnegative tensor factorization. In *Tensors in Image Processing and Computer Vision*, pp. 105–124. Springer, 2009a.
- Zafeiriou, Stefanos. Discriminant nonnegative tensor factorization algorithms. *Neural Networks, IEEE Transactions on*, 20(2):217–235, 2009b.
- Zafeiriou, Stefanos and Petrou, Maria. Nonlinear non-negative component analysis algorithms. *Image Processing, IEEE Transactions on*, 19(4):1050–1066, 2010.
- Zafeiriou, Stefanos, Tefas, Anastasios, Buciu, Ioan, and Pitas, Ioannis. Exploiting discriminant information in nonnegative matrix factorization with application to frontal face verification. *Neural Networks, IEEE Transactions on*, 17(3):683–695, 2006.
- Zafeiriou, Stefanos, Tzimiropoulos, Georgios, Petrou, Maria, and Stathaki, Tania. Regularized kernel discriminant analysis with a robust kernel for face recognition and verification. *Neural Networks and Learning Systems, IEEE Transactions on*, 23(3): 526–534, 2012.