

---

# A Multiplicative Up-Propagation Algorithm

---

**Jong-Hoon Ahn**

Department of Physics, POSTECH, Korea

JONGHUN@POSTECH.AC.KR

**Seungjin Choi**

Department of Computer Science, POSTECH, Korea

SEUNGJIN@POSTECH.AC.KR

**Jong-Hoon Oh**

Department of Physics, POSTECH, Korea

JHOH@POSTECH.AC.KR

## Abstract

We present a generalization of the nonnegative matrix factorization (NMF), where a multilayer generative network with nonnegative weights is used to approximate the observed nonnegative data. The multilayer generative network with nonnegativity constraints, is learned by a multiplicative up-propagation algorithm, where the weights in each layer are updated in a multiplicative fashion while the mismatch ratio is propagated from the bottom to the top layer. The monotonic convergence of the multiplicative up-propagation algorithm is shown. In contrast to NMF, the multiplicative up-propagation is an algorithm that can learn hierarchical representations, where complex higher-level representations are defined in terms of less complex lower-level representations. The interesting behavior of our algorithm is demonstrated with face image data.

## 1. Introduction

Nonnegative matrix factorization (NMF) which was published in Nature a few years ago (Lee & Seung, 1999), drew extensive attraction in machine learning and pattern recognition communities. It is one of the efforts to realize a parts-based representation which is a way of understanding the perception in the brain and certain computational theories rely on such a representation. We also encounter into a variety of practical applications where parts-based representations of non-

negative data are useful. These include face images for computer vision (Li et al., 2001), medical image analysis (Lee et al., 2001), spectrograms of sound data (Cho et al., 2003), text data, and so on (Saul et al., 2003).

NMF is a simple linear coding algorithm for a single-layer generative network with nonnegativity constraints. Although it is an efficient linear coding method for nonnegative data, the linear single-layer generative network leads to several limitations in a certain general problem such as:

- when it learns data which lie on or near a nonlinear manifold;
- when it learns syntactic relationships of the given data;
- when it learns hierarchically generated data.

NMF is obviously not suitable for such cases. For instance, when NMF is applied to some sequences of images which contain facial expressions, it does not show a parts-based representation clearly. For images of objects viewed from extremely different viewpoints, it requires much more data and additional degree of freedom. NMF assumes that the hidden variables are nonnegative, but makes no further assumptions about their statistical dependency. Accordingly, NMF cannot learn anything about the syntactic relationships between parts. For instance, given two basis vectors of eyes and eyes+eyebrows computed by NMF, we cannot know the similarity or dependency between them. Let us consider a problem of finding a set of initial basis vectors, given a set of hierarchically generated data, which is described as follows: 1. Assume that we have a set of  $p_1$  nonnegative unit vectors; 2. Try to produce  $p_2$  composites by mixing and transforming them; 3. Try to repeat the step 2 and produce another

---

Appearing in *Proceedings of the 21<sup>st</sup> International Conference on Machine Learning*, Banff, Canada, 2004. Copyright 2004 by the authors.

$p_{i+1}$  composites from the  $i$ th composites; 4. When we see only the final composites, could we find the initial  $p_1$  nonnegative vectors;

Any modifications of the NMF algorithm in a single-layer network do not yield a solution to this problem. Therefore, we propose a generalization of NMF, where we employ a multilayer generative network with nonnegativity constraints, which still preserve a useful property of NMF such as parts-based representations. A generalization of NMF with employing a multilayer generative network, is a nonlinear extension since squashing functions should be incorporated with the multilayer network.

A conventional multilayer perceptron (MLP) that belong to a recognition model, can be efficiently trained by an error back-propagation algorithm. It was also shown that the multilayer generative network could be trained by an error up-propagation algorithm (Oh & Seung, 1997). A main difference between these two networks is that the former is a supervised network and the latter is an unsupervised network, whereas both of them are trained by gradient-based additive algorithms. The multilayer network that we consider in this paper, is a generative network with all of synaptic weights being nonnegative. Therefore, the error up-propagation algorithm (Oh & Seung, 1997) cannot be directly applied for learning our multilayer generative network. However, this suggests us to use an idea of error up-propagation. Motivated from an idea of multiplicative update used in NMF, we develop a new learning algorithm, *multiplicative up-propagation* for training the multilayer generative network with nonnegativity constraints. In our multiplicative up-propagation algorithm, synaptic weights for each layer are trained in a multiplicative fashion (in order to preserve nonnegativity constraints), while the mismatch (or the error) ratio is propagated from the bottom to the top layer.

Two exemplary NMF algorithms result from objective functions based on I-divergence and squared error. We show that NMF could be a general optimizing rule for various objective functions. As in gradient descent rules used in back-prop and up-prop algorithms for a given error function, our proposed multiplicative algorithm can also be used for general objective functions with nonnegativity constraints.

The rest of this paper is organized as follows. Next section briefly reviews the original up-propagation algorithm (Oh & Seung, 1997) and NMF (Lee & Seung, 1999). In Sec. 3, we illustrate a detailed derivation of the proposed multiplicative up-propagation algorithm. Its monotonic convergence is also shown in Sec. 3.

Then the algorithm is applied to face image data and its interesting behavior is shown in Sec. 4. Finally conclusion is drawn in Sec. 5.

## 2. Previous Work

### 2.1. Up-propagation

The up-propagation is an algorithm for inverting and learning a multi-layer generative model, which generalizes PCA and its variants such as conic coding (Oh & Seung, 1997). The generative model is a network of  $L$  layers of neurons with input layer at bottom, as depicted in fig. 2. The matrices  $\mathbf{H}^{(l)}$ ,  $l = 1, \dots, L$ , are the activations of the layers. The pattern  $\mathbf{H}^{(0)}$  is generated from the hidden variables  $\mathbf{H}^{(L)}$  by a top-down pass through the network,

$$\mathbf{H}^{(l-1)} = \mathbf{g}(\mathbf{W}^{(l)} \mathbf{H}^{(l)}), \quad l = L, \dots, 1.$$

The nonlinear function  $\mathbf{g}$  acts on matrices component by component. The matrix  $\mathbf{W}^{(l)}$  contains the synaptic connections from the neurons in layer  $l$  to the neurons in layer  $l - 1$ . A bias term  $\mathbf{B}^{(l)}$  can be added to the argument of  $\mathbf{g}$ , but is omitted here.

Given a sensory input  $\mathbf{V}$ , the top-down generative model can be inverted by finding hidden variables  $\mathbf{H}^{(L)}$  that generate a pattern  $\mathbf{H}^{(0)}$  matching  $\mathbf{V}$ . If some of the hidden variables represent the identity of the pattern, the inversion operation is called recognition. Alternatively, the hidden variables may just be a more compact representation of the pattern, in which case the operation is called analysis or encoding. The inversion is done iteratively, as described below.

In the following, the operator  $\odot$  denotes the Hadamard product (elementwise multiplication), that is,  $\mathbf{Z} = \mathbf{X} \odot \mathbf{Y}$  means  $z_{ij} = x_{ij}y_{ij}$ . The bottom-up pass starts with the mismatch between the sensory data  $\mathbf{V}$  and the generated pattern  $\mathbf{H}^{(0)}$ ,

$$\Delta^{(1)} = \mathbf{g}'(\mathbf{W}^{(1)} \mathbf{H}^{(1)}) \odot (\mathbf{V} - \mathbf{H}^{(0)}), \quad (1)$$

which is propagated upwards by

$$\Delta^{(l+1)} = \mathbf{g}'(\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)}) \odot ((\mathbf{W}^{(l)})^T \Delta^{(l)}). \quad (2)$$

When the error signal reaches the top of the network, it is used to update the hidden variables  $\mathbf{H}^{(L)}$ ,

$$\Delta \mathbf{H}^{(L)} = -\eta (\mathbf{W}^{(L)})^T \Delta^{(L)}. \quad (3)$$

This update closes the negative feedback loop. Then a new pattern  $\mathbf{H}^{(0)}$  is generated by a top-down pass, and the process starts over again.

This iterative inversion process performs gradient descent on the cost function  $\frac{1}{2} \|\mathbf{V} - \mathbf{H}^{(0)}\|^2$ , subject

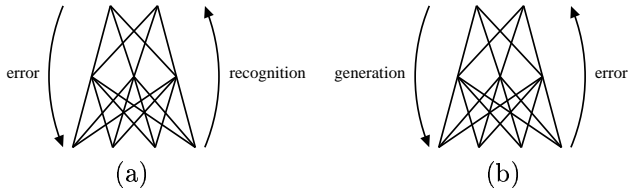


Figure 1. Bottom-up and top-down processing in neural network : (a) backprop network; (b) up-prop network

to the constraints. This can be proved using the chain rule, as in the traditional derivation of the backprop algorithm. Another method of proof is to add the equations as constraints, using Lagrange multipliers,

$$\frac{1}{2} \|\mathbf{V} - \mathbf{H}^{(0)}\|^2 + \sum_{l=1}^L (\Delta^{(l-1)})^T (\mathbf{H}^{(l-1)} - \mathbf{g}(\mathbf{W}^{(l)} \mathbf{H}^{(l)})). \quad (4)$$

This derivation has the advantage that the bottom-up activations  $\Delta^{(l)}$  have an interpretation as Lagrange multipliers.

Inverting the generative model by negative feedback can be interpreted as a process of sequential hypothesis testing. The top-down connections generate a hypothesis about the sensory data. The bottom-up connections propagate an error signal that is the disagreement between the hypothesis and data. When the error signal reaches the top, it is used to generate a revised hypothesis, and the generate-test-revise cycle starts all over again. Perception is the convergence of this feedback loop to the hypothesis that is most consistent with the data.

The synaptic weights  $\mathbf{W}^{(l)}$  determine the types of patterns that the network is able to generate. To learn from examples, the weights are adjusted to improve the network's generation ability. A suitable cost function for learning is the reconstruction error  $\frac{1}{2} \|\mathbf{V} - \mathbf{H}^{(0)}\|^2$  averaged over an ensemble of examples. Online gradient descent with respect to the synaptic weights is performed by a learning rule of the form

$$\Delta \mathbf{W}^{(l)} = \eta \Delta^{(l-1)} (\mathbf{H}^{(l)})^T. \quad (5)$$

The same error signal  $\Delta^{(l)}$  that was used to invert the generative model is also used to learn it.

## 2.2. Nonnegative Matrix Factorization

Let a set of data be given as an  $p \times N$  matrix  $\mathbf{V}$ , with each column consisting of the  $p$  non-negative pixel values of an image. Denote a set of basis images as a  $p \times q$  matrix  $\mathbf{W}$ . Each image can be approximated as a lin-

ear combination of the basis images using coefficients,

$$\mathbf{V} \approx \mathbf{W} \mathbf{H} \quad (6)$$

where  $\mathbf{H} \in \mathbf{R}^{q \times N}$  is the matrix which contains the coefficients.

PCA requires that the columns of  $\mathbf{W}$  be orthonormal and the rows of  $\mathbf{H}$  be mutually orthogonal. It imposes no other constraints than the orthogonality and hence allows the entries of  $\mathbf{W}$  and  $\mathbf{H}$  to be of arbitrary sign. Many basis images, for instance eigenfaces, lack intuitive meaning, and a linear combination of the basis images generally involves complex cancellations between positive and negative values. The NMF representations permit only nonnegative coefficients and thus non-subtractive combinations (Lee & Seung, 1999; Lee & Seung, 2000).

On the other hand NMF imposes the non-negativity constraints instead of the orthogonal basis and decorrelated hidden variables. The elements of  $\mathbf{W}$  and  $\mathbf{H}$  are all non-negative, and hence only non-subtractive combinations are allowed. This is believed to be compatible to the intuitive notion of combining parts to form a whole, and is how NMF learns the parts-based representation. It is also consistent with the physiological facts that the firing rates are non-negative and the signs of synapses do not change.

NMF uses the I-divergence of  $\mathbf{V}$  from  $\mathbf{W} \mathbf{H}$ , defined as

$$D(\mathbf{V} \parallel \mathbf{W} \mathbf{H}) = \sum_{i,j} \left( \mathbf{v}_{ij} \log \frac{\mathbf{v}_{ij}}{(\mathbf{W} \mathbf{H})_{ij}} - \mathbf{v}_{ij} + (\mathbf{W} \mathbf{H})_{ij} \right) \quad (7)$$

as the measure of fitness for factorizing  $\mathbf{V}$  into  $\mathbf{W} \mathbf{H}$ . The NMF factorization is defined as

$$\min_{\mathbf{W}, \mathbf{H}} D(\mathbf{V} \parallel \mathbf{W} \mathbf{H}), \quad (8) \\ s.t. \quad \mathbf{W}, \mathbf{H} \geq 0.$$

$D(\mathbf{V} \parallel \mathbf{W} \mathbf{H})$  reduces to Kullback-Leibler divergence when  $\sum_{i,j} \mathbf{v}_{ij} = \sum_{i,j} (\mathbf{W} \mathbf{H})_{ij} = 1$ . The above optimization can be done by using multiplicative update rules.

$$\mathbf{H}_{a\mu} \leftarrow \mathbf{H}_{a\mu} \frac{\sum_i \mathbf{W}_{ia} \mathbf{v}_{i\mu} / (\mathbf{W} \mathbf{H})_{i\mu}}{\sum_k \mathbf{W}_{ka}}, \\ \mathbf{W}_{ia} \leftarrow \mathbf{W}_{ia} \frac{\sum_\mu \mathbf{H}_{a\mu} \mathbf{v}_{i\mu} / (\mathbf{W} \mathbf{H})_{i\mu}}{\sum_\nu \mathbf{H}_{a\nu}}. \quad (9)$$

The algorithm performs both learning and inference simultaneously. That is, it learns a set of basis images and infers values for the hidden variables from

the visible variables. Although the generative model is linear, the inference computation is nonlinear due to the non-negativity constraints. The computation is a type of generalized EM algorithm.

Although NMF is successful in learning facial parts and semantic topics, this success does not imply that the method can learn parts from any database, such as images of objects viewed from extremely different viewpoints, or highly articulated objects. Learning parts for these complex cases is likely to require fully hierarchical models with multiple levels of hidden variables, instead of the single level in NMF (Lee & Seung, 1999).

### 3. The Multiplicative Up-Propagation Algorithm

There are a large number of neurons in inferior temporal cortex of monkeys which seem to encode an overall shape of biologically important objects - not specific features or parts. The finding agrees with hierarchical theories of object perception. According to these theories, cells in the cortical areas code elementary features such as line orientation and color. The outputs from these cells are then combined by detectors sensitive to higher-order features such as corners or intersections, an idea consistent with the findings of Hubel and Wiesel. The process is continued as each successive stage codes more complex combinations. At the top of the chain are IT neurons, selective for complex shapes like hands or faces. A huge number of hierarchical models for object recognition have been proposed over the years. Some of them were inspired by the desire to build intelligent machines, others by the desire to describe human recognition processes (Gazzaniga et al., 2001).

In this paper, we will try to invert the hierarchical recognition processes and view the visual perception as a hypothesis testing process. Helmholtz, in his doctrine of unconscious inference, argued that perceptions are formed by the interaction of bottom-up sensory data with top-down expectations. According to one interpretation of this doctrine, perception is a procedure of sequential hypothesis testing. We propose a new algorithm, called multiplicative up-propagation, that realizes this interpretation in layered networks. It uses top-down connections to generate hypotheses, and bottom-up connections to revise them.

How can we build such a hierarchical structure of special neurons which are responsible only for local sensory features by multi-layer generative model? We have already known that NMF can give local features

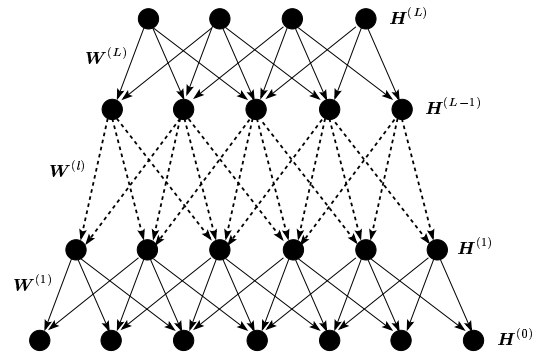


Figure 2. Multilayer generative network model:  $\mathbf{V} \approx \mathbf{H}^{(0)} = \mathbf{g}(\mathbf{W}^{(1)} \mathbf{g}(\mathbf{W}^{(2)} \mathbf{g}(\dots \mathbf{g}(\mathbf{W}^{(L)} \mathbf{H}^{(L)}))))$  and  $\mathbf{H}^{(l)} = \mathbf{g}(\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)})$ . Note that it will be fully connected.

and sparse codes from a given non-negative data set. It is just a single-layer generative network and gives local features and sparse codes. Then what about multi-layer network with nonnegative weights and hidden variables? It is what we will introduce in this section. The previous notations of  $\mathbf{W}$  and  $\mathbf{H}$  in a single-layer network are changed into  $\mathbf{W}^{(1)}$  and  $\mathbf{H}^{(1)}$ . We start by assuming further factorizations:

$$\mathbf{H}_{a\mu}^{(l)} = \mathbf{g} \left( (\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)})_{a\mu} \right), \quad (10)$$

where  $l = 0, 1, \dots, L - 1$ . The nonlinear function  $\mathbf{g}$ , which must output nonnegative values and increase monotonically, is applied elementwise to the matrix. Then we can construct a  $L$ -layer nonnegative networks as shown in Fig. 2. It is basically a generative model and constructs the nonnegative input data at its bottom layer. Before we show how it finds a new representation from the nonnegative data, we should think of the algorithm which we can find the optimized values of weights and hidden variables with. In order to extend NMF to multilayer case, let us introduce two matrices  $\mathbf{R}$  and  $\mathbf{N}$  to the update rules of equation (9):

$$\begin{aligned} \mathbf{H}_{a\mu} &\leftarrow \mathbf{H}_{a\mu} \frac{\sum_i \mathbf{W}_{ia} \frac{\mathbf{V}_{i\mu}}{(\mathbf{W}\mathbf{H})_{i\mu}}}{\sum_k \mathbf{W}_{ka}} \\ &= \mathbf{H}_{a\mu} \frac{\sum_i \mathbf{W}_{ia} \mathbf{R}_{i\mu}}{\sum_k \mathbf{W}_{ka} \mathbf{N}_{k\mu}} \\ &= \mathbf{H}_{a\mu} \frac{(\mathbf{W}^T \mathbf{R})_{a\mu}}{(\mathbf{W}^T \mathbf{N})_{a\mu}}, \end{aligned} \quad (11)$$

Table 1. Cost functions and initializations of  $\mathbf{N}^{(l)}$  and  $\mathbf{R}^{(l)}$  : If we use other cost functions, the initializations should be changed

F	$\mathbf{N}^{(1)}$	$\mathbf{R}^{(1)}$
$\ \mathbf{V} - \mathbf{H}^{(0)}\ ^R$	$\mathbf{N}_{i\mu}^{(1)} = (\mathbf{H}_{i\mu}^{(0)})^{R-1} \mathbf{g}'((\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu})$	$\mathbf{R}_{i\mu}^{(1)} = (\mathbf{V}_{i\mu}^{(0)})^{R-1} \mathbf{g}'((\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu})$
$\sum_{i\mu} \mathbf{V}_{i\mu} \log(\mathbf{H}_{i\mu}^{(0)}) - \mathbf{H}_{i\mu}^{(0)}$	$\mathbf{N}_{i\mu}^{(1)} = \mathbf{g}'((\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu})$	$\mathbf{R}_{i\mu}^{(1)} = \mathbf{g}'((\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu}) \mathbf{V}_{i\mu} / \mathbf{H}_{i\mu}^{(0)}$
$\sum_{i\mu} \mathbf{V}_{i\mu} / \mathbf{H}_{i\mu}^{(0)} + \log \mathbf{H}_{i\mu}^{(0)}$	$\mathbf{N}_{i\mu}^{(1)} = \mathbf{g}'((\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu}) / \mathbf{H}_{i\mu}^{(0)}$	$\mathbf{R}_{i\mu}^{(1)} = \mathbf{g}'((\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu}) \mathbf{V}_{i\mu} / (\mathbf{H}_{i\mu}^{(0)})^2$

and

$$\begin{aligned}
\mathbf{W}_{ia} &\leftarrow \mathbf{W}_{ia} \frac{\sum_{\mu} \frac{\mathbf{V}_{i\mu}}{(\mathbf{W}\mathbf{H})_{i\mu}} \mathbf{H}_{a\mu}}{\sum_{\nu} \mathbf{H}_{a\nu}} \\
&= \mathbf{W}_{ia} \frac{\sum_{\mu} \mathbf{R}_{i\mu} \mathbf{H}_{a\mu}}{\sum_{\nu} \mathbf{N}_{i\nu} \mathbf{H}_{a\nu}} \\
&= \mathbf{W}_{ia} \frac{(\mathbf{R}\mathbf{H}^T)_{ia}}{(\mathbf{N}\mathbf{H}^T)_{ia}} \quad (12)
\end{aligned}$$

where  $\mathbf{R} = \mathbf{R}^{(1)}$  is a matrix of ratios of  $\mathbf{V}$  to  $\mathbf{W}\mathbf{H}$  and  $\mathbf{N} = \mathbf{N}^{(1)}$  is a matrix of normalizing factors with all elements 1. If we put  $\mathbf{R} = \mathbf{V}$  and  $\mathbf{N} = \mathbf{W}\mathbf{H}$ , the above rules will be for least squared error under non-negativity. In fact such formulations can be applied to all cost functions F that are defined as follows:

**Definition 1** F is a cost function that quantify the quality of the approximation between  $\mathbf{V}$  and  $\mathbf{H}^{(0)} = \mathbf{g}(\mathbf{W}^{(1)} \mathbf{g}(\mathbf{W}^{(2)} \dots))$ , if the conditions

$$\begin{aligned}
\frac{\partial \mathbf{F}}{\partial \mathbf{W}^{(l)}} &= (\mathbf{N}^{(l)} - \mathbf{R}^{(l)}) (\mathbf{H}^{(l)})^T \quad (13) \\
\frac{\partial \mathbf{F}}{\partial \mathbf{H}^{(L)}} &= (\mathbf{W}^{(L)})^T (\mathbf{N}^{(L)} - \mathbf{R}^{(L)})
\end{aligned}$$

are satisfied, where  $\mathbf{R}^{(l)}$  and  $\mathbf{N}^{(l)}$  should be nonnegative matrices for arbitrary nonnegative matrices  $\mathbf{W}^{(l)}$  and  $\mathbf{H}^{(l)}$ .

If we obtain the matrix  $\mathbf{R}^{(l)}$  and  $\mathbf{N}^{(l)}$  at  $l$ th layer of the networks, we can optimize the multi-layer networks by using the similar update rules at all layers:

$$\mathbf{W}_{ia}^{(l)} \leftarrow \mathbf{W}_{ia}^{(l)} \left[ \frac{(\mathbf{R}^{(l)} \mathbf{H}^{(l)T})_{ia}}{(\mathbf{N}^{(l)} \mathbf{H}^{(l)T})_{ia}} \right]^{\eta} \quad (14)$$

and

$$\mathbf{H}_{a\mu}^{(L)} \leftarrow \mathbf{H}_{a\mu}^{(L)} \left[ \frac{(\mathbf{W}^{(L)T} \mathbf{R}^{(L)})_{a\mu}}{(\mathbf{W}^{(L)T} \mathbf{N}^{(L)})_{a\mu}} \right]^{\eta}, \quad (15)$$

where the learning rate  $0 < \eta \ll 1$  should be considered in nonlinear multilayer networks. Fortunately,

the  $\mathbf{R}^{(l)}$  and  $\mathbf{N}^{(l)}$  matrices are calculated by the following up-propagation rules:

$$\begin{aligned}
\mathbf{R}_{i\mu}^{(l+1)} &= (\mathbf{W}^{(l)T} \mathbf{R}^{(l)})_{i\mu} \mathbf{g}'(\mathbf{W}^{(l)} \mathbf{H}^{(l)})_{i\mu}, \\
\mathbf{N}_{i\mu}^{(l+1)} &= (\mathbf{W}^{(l)T} \mathbf{N}^{(l)})_{i\mu} \mathbf{g}'(\mathbf{W}^{(l)} \mathbf{H}^{(l)})_{i\mu} \quad (16)
\end{aligned}$$

where  $l = 1, 2, \dots, L-1$  and

$$\begin{aligned}
\mathbf{R}_{i\mu}^{(1)} &= \frac{\mathbf{V}_{i\mu}}{(\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu}} \mathbf{g}'((\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu}), \\
\mathbf{N}_{i\mu}^{(1)} &= \mathbf{g}'((\mathbf{W}^{(1)} \mathbf{H}^{(1)})_{i\mu}). \quad (17)
\end{aligned}$$

Another form of initializations is given in Table 1.

**Theorem 1** F is nonincreasing under the update rules of eqs. (14) and (15). The cost is invariant under these updates if and only if  $\mathbf{W}^{(l)}$  and  $\mathbf{H}^{(L)}$  are at a stationary point of the cost function.

To prove Theorem 1, we will make use of the auxiliary function that was used in (Lee & Seung, 2000).

**Definition 2**  $\mathbf{G}(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)})$  is an auxiliary function for  $\mathbf{F}(\mathbf{W}^{(l)})$  if the conditions

$$\mathbf{G}(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)}) \geq \mathbf{F}(\mathbf{W}^{(l)}), \quad \mathbf{G}(\mathbf{W}_t^{(l)}, \mathbf{W}_t^{(l)}) = \mathbf{F}(\mathbf{W}_t^{(l)})$$

are satisfied.  $t$  is a discrete time index.

The auxiliary function is a useful concept because of the following lemma.

**Lemma 1** If  $\mathbf{G}$  is an auxiliary function, then  $\mathbf{F}$  is non-increasing under the update

$$\mathbf{W}_{t+1}^{(l)} = \arg \min_{\mathbf{W}^{(l)}} \mathbf{G}(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)})$$

**Proof:**  $\mathbf{F}(\mathbf{W}_{t+1}^{(l)}) \leq \mathbf{G}(\mathbf{W}_{t+1}^{(l)}, \mathbf{W}_t^{(l)}) \leq \mathbf{G}(\mathbf{W}_t^{(l)}, \mathbf{W}_t^{(l)}) = \mathbf{F}(\mathbf{W}_t^{(l)})$ . ■

**Lemma 2** If the derivatives of  $G(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)})$  are

$$\begin{aligned} & \frac{\partial G(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)})}{\partial (\mathbf{W}^{(l)})_{ia}} \\ &= \frac{(\mathbf{W}^{(l)})_{ia}^{1/\eta}}{(\mathbf{W}_t^{(l)})_{ia}^{1/\eta}} \left( \mathbf{N}^{(l)}(\mathbf{H}^{(l)T})_{ia} - \mathbf{R}^{(l)}(\mathbf{H}^{(l)T})_{ia} \right), \end{aligned}$$

there exists an auxiliary function  $G(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)})$  for  $F(\mathbf{W}_t^{(l)})$ .

**Proof:** Let  $E$  be the difference of  $G$  from  $F$ ,  $G - F$ . Then the derivatives of  $E(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)})$  are

$$\begin{aligned} & \frac{\partial E(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)})}{\partial (\mathbf{W}^{(l)})_{ia}} \\ &= \frac{(\mathbf{W}^{(l)})_{ia}^{1/\eta} - (\mathbf{W}_t^{(l)})_{ia}^{1/\eta}}{(\mathbf{W}_t^{(l)})_{ia}^{1/\eta}} \{ \mathbf{N}^{(l)}(\mathbf{H}^{(l)T})_{ia} \}. \end{aligned}$$

Since the derivatives are nonnegative for  $(\mathbf{W}^{(l)})_{ia} \geq (\mathbf{W}_t^{(l)})_{ia}$  and negative for  $(\mathbf{W}^{(l)})_{ia} < (\mathbf{W}_t^{(l)})_{ia}$ , there exists the auxiliary function  $E(\mathbf{W}^{(l)}, \mathbf{W}_t^{(l)}) \geq 0$  by adding a constant so that  $E$  can be zero at  $\mathbf{W}^{(l)} = \mathbf{W}_t^{(l)}$ . ■

**Proof of Theorem 1** By lemma 1, putting the derivatives zero in lemma 2 results in the update rule:

$$\begin{aligned} (\mathbf{W}_{t+1}^{(l)})_{ia} &= (\mathbf{W}_t^{(l)})_{ia} \frac{\left\{ \mathbf{R}_{t+1}^{(l)}(\mathbf{H}^{(l)T})_{ia} \right\}^\eta}{\left\{ \mathbf{N}_{t+1}^{(l)}(\mathbf{H}^{(l)T})_{ia} \right\}^\eta} \\ &\approx (\mathbf{W}_t^{(l)})_{ia} \frac{\left\{ \mathbf{R}_t^{(l)}(\mathbf{H}^{(l)T})_{ia} \right\}^\eta}{\left\{ \mathbf{N}_t^{(l)}(\mathbf{H}^{(l)T})_{ia} \right\}^\eta}, \quad (18) \end{aligned}$$

from  $\mathbf{R}_{t+1}^{(l)} \approx \mathbf{R}_t^{(l)}$  and  $\mathbf{N}_{t+1}^{(l)} \approx \mathbf{N}_t^{(l)}$ . Eq. (15) can be shown similarly.

On the other hand, we can derive the up-propagation rules from the definition 1. By considering

$$\frac{\partial F}{\partial (\mathbf{W}^{(l)} \mathbf{H}^{(l)})} = \mathbf{N}^{(l)} - \mathbf{R}^{(l)}$$

from the chain rule  $\frac{\partial F}{\partial \mathbf{W}^{(l)}} = \frac{\partial F}{\partial (\mathbf{W}^{(l)} \mathbf{H}^{(l)})} (\mathbf{H}^{(l)T})$  and

the eqn. (13), we get

$$\begin{aligned} & \frac{\mathbf{N}^{(l+1)} - \mathbf{R}^{(l+1)}}{\partial F} \\ &= \frac{\partial F}{\partial (\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)})} \\ &= \frac{\partial F}{\partial \mathbf{g}(\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)})} \odot \mathbf{g}'(\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)}) \\ &= \frac{\partial F}{\partial \mathbf{H}^{(l)}} \odot \mathbf{g}'(\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)}) \\ &= \left[ (\mathbf{W}^{(l)})^T \frac{\partial F}{\partial (\mathbf{W}^{(l)} \mathbf{H}^{(l)})} \right] \odot \mathbf{g}'(\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)}) \\ &= \left[ (\mathbf{W}^{(l)})^T (\mathbf{N}^{(l)} - \mathbf{R}^{(l)}) \right] \odot \mathbf{g}'(\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)}), \end{aligned}$$

where the 3rd and 5th equations are expanded by using the chain rules.

Thus,  $\mathbf{N}^{(l+1)} - \mathbf{R}^{(l+1)} = \left[ (\mathbf{W}^{(l)})^T (\mathbf{N}^{(l)} - \mathbf{R}^{(l)}) \right] \odot \mathbf{g}'(\mathbf{W}^{(l+1)} \mathbf{H}^{(l+1)})$  or equation (16) can be obtained. ■

The proposed algorithm is very similar with the error up-propagation algorithm for multi-layer neural network model. We can find the rule by which error up-prop algorithm is modified into multiplicative up-prop algorithm: Elementwise addition(or subtraction) is changed into elementwise multiplication(or division) and the multiplication of the learning rate  $\eta$  is modified into the power of  $\eta$ . Both of them also use the same type of up-prop rules (2) and (16). But they are also different in a couple of aspects. The error up-propagation algorithm propagates only error, whereas the proposed algorithm should separately propagate both normalizing factors and ratios of inputs to reconstructed values.

## 4. Experiments

We applied the three-layer network model to a set of facial images and trained it by using the multiplicative up-propagation rules. As with all gradient-based algorithms, the multiplicative up-propagation rules are susceptible to local optima. Crucial to the success of our experiments is a pertinent network growing algorithm.

### 4.1. Network Growing

We begin the training of the networks from a single-layer networks. In the first stage, the given data matrix  $\mathbf{V}$  are factorized into  $\mathbf{g}(\mathbf{W}^{(1)} \mathbf{H}^{(1)})$ . Then we add one layer on that and train the second layer so that the  $\mathbf{H}^{(1)}$  can be factorized into  $\mathbf{g}(\mathbf{W}^{(2)} \mathbf{H}^{(2)})$ . The two-layer networks are then trained by using the multiplicative up-prop rules for the input matrix  $\mathbf{V}$  to

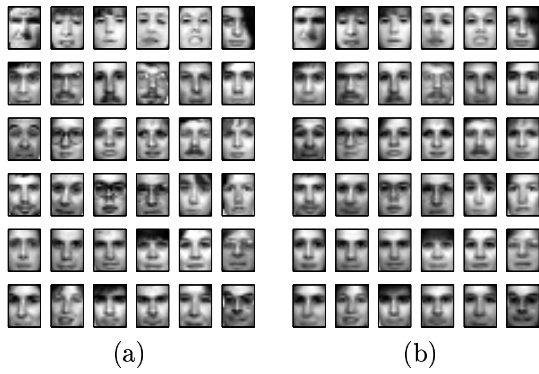


Figure 3. (a) original face images; (b) reconstructed face images

be approximated by the matrix  $\mathbf{H}^{(0)}$ . The process to grow the networks from the two-layer to three-layer is the same as the previous steps. However, the method does not assure the globally optimal solution to us. It only helps avoid some worst solutions. Also note that applying NMF recursively to the hidden values  $\mathbf{H}^{(l)}$  is not the same as our multiplicative up-propagation rules.

#### 4.2. Hierarchical Feature Analysis

We obtain a hierarchical set of features of facial images through the three-layer generative networks. By applying the model and learning algorithm to a set of size-normalized and eye-centered 1608 facial images with size  $19 \times 14$  (see Fig. 3), we obtained the visually hierarchical features. The transformed weights of the first layer represent the highly localized features of Fig. 4c, whereas the features of Fig. 4a correspond to the nodes of top layer. The features of Fig. 4b are similar to those of NMF. The complex features in a layer are defined by the less complex features in the subordinate layers, which is the reason why we name the application a 'hierarchical feature analysis'.

This property is specially helpful in the case that we should find the hierarchical relations or dependencies of the features. We've already made an application to the case and obtained some successful results (Ahn et al., 2004). We can develop the algorithm and perform the hierarchical clustering of a set of nonnegative data by constraining the algorithm to a condition such that the summation of column vectors is one.

#### 4.3. Representational Learning in IT Cortex

We are already familiar with the fact that mammalian visual system is hierarchically constructed, through which it acquires the flexible and invariant recognition

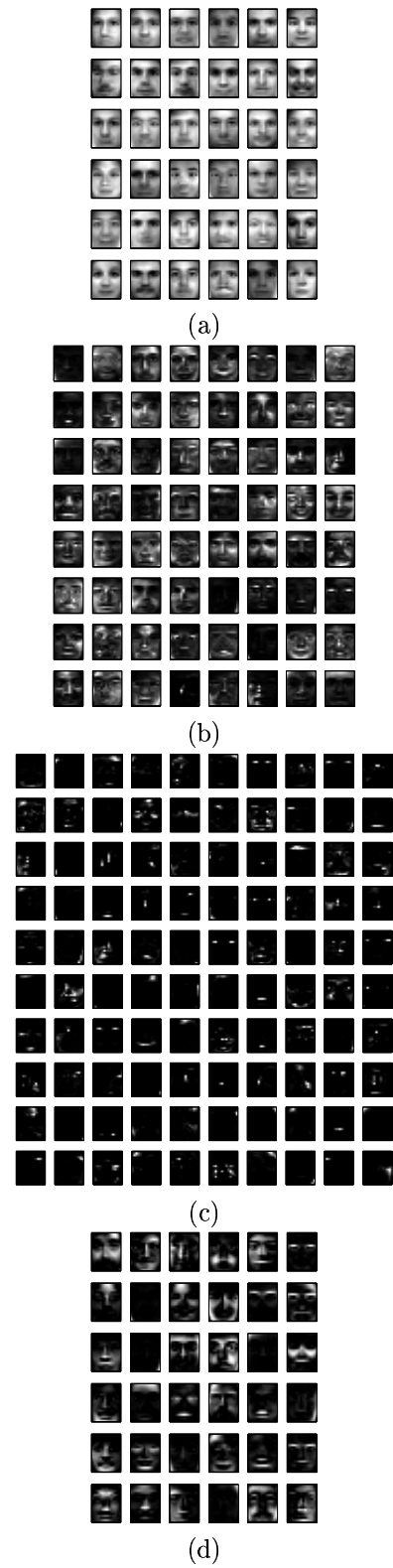


Figure 4. Three sets of facial features in a hierarchical relationship : (a) column vectors of  $\mathbf{g}(\mathbf{W}^{(1)}\mathbf{g}(\mathbf{W}^{(2)}\mathbf{g}(\mathbf{W}^{(3)})))$ ; (b) column vectors of  $\mathbf{g}(\mathbf{W}^{(1)}\mathbf{g}(\mathbf{W}^{(2)}))$ ; (c) column vectors of  $\mathbf{g}(\mathbf{W}^{(1)})$ ; (d) features obtained by NMF

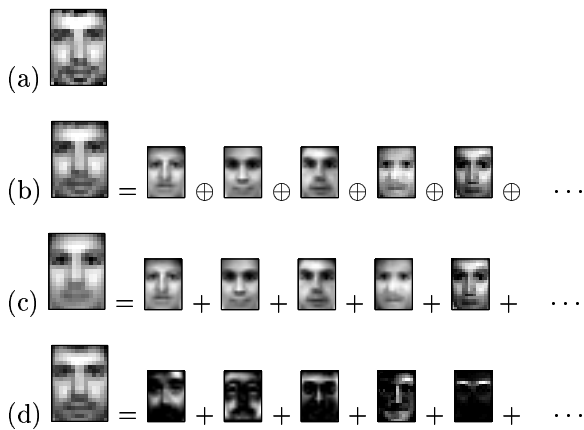


Figure 5. Whole-based representation vs. Parts-based representation : (a) original face; (b) our method (nonlinear summation,  $\sum_k g(\mathbf{W}^{(1)}g(\mathbf{W}^{(2)}g(\mathbf{W}^{(3)}h_k^{(3)}))))$ ); (c) our method (linear summation,  $\sum_k g(\mathbf{W}^{(1)}g(\mathbf{W}^{(2)}g(\mathbf{W}^{(3)})))h_k^{(3)}$ ); (d) NMF

ability. What we are interested in here is IT cortex, and its representational learning. The areas receive their inputs via a number of cortico-cortical stages from the primary visual cortex, striate cortex, through pre-striate visual areas. Neurons that are found in certain area of IT cortex preferentially or selectively to faces. When they were found in the pioneer days, some researchers conceived that they could respond to the whole face of a particular individual, that is 'Grandmother Cell' hypothesis. But it resulted in a few of problems such that generalization, noise robustness and excess of number of neurons. Instead, ensemble coding theory was popularized for representational learning of IT cortex. So to speaking, when a face is considered, a special set of parts is required from a pool of various types of parts such that eyes, nose and mouth. NMF showed a principle for the parts-based representation and it was noticed. But what is a role of the neurons that respond only to faces? Here we argue that there could exist a whole-based representation of faces in Fig. 5. We guess together that representations of facial information would be different according to the level in hierarchy, and the top layer of IT cortex could use the whole-based representations.

## 5. Conclusion

We proposed a new algorithm, multiplicative up-propagation, for training the multilayer generative network model with the nonnegativity. It could be considered as a nonlinear and multilayer extension of NMF. We applied the three-layer generative network model to the eye-aligned facial images and showed details of

parts and composites in hierarchical relations. We argue that these results better explain the representational learning of IT cortex.

## 6. Acknowledgment

This work was supported by Korea Ministry of Science and Technology under Brain Science and Engineering Research Program, KOSEF 2000-2-20500-009-5, and Brain Korea 21 in POSTECH.

## References

- Ahn, J.-H., Kim, S., Oh, J.-H., & Choi, S. (2004). Multiple nonnegative-matrix factorization of dynamic pet images. *Asian Conference on Computer Vision*.
- Cho, Y.-C., Choi, S., & Bang, S.-Y. (2003). Non-negative component parts of sound for classification. *Proc. IEEE Int. Symp. Signal Processing and Information Technology*. Darmstadt, Germany.
- Gazzaniga, M. S., Ivry, R. B., & Mangum, G. R. (2001). *Cognitive neuroscience: The biology of the mind*. New York: W. W. Norton & Company.
- Lee, D. D., & Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401, 788–791.
- Lee, D. D., & Seung, H. S. (2000). Algorithms for non-negative matrix factorization. *Advances in Neural Information Processing Systems*.
- Lee, J. S., Lee, D. D., Choi, S., & Lee, D. S. (2001). Application of non-negative matrix factorization to dynamic positron emission tomography. *Proc. ICA '01* (pp. 629–632).
- Li, S. Z., Hou, X. W., Zhang, H. J., & Cheng, Q. S. (2001). Learning spatially localized, parts-based representation. *Proc. IEEE Conf. Computer Vision and Pattern Recognition* (pp. 207–212). Kauai, Hawaii.
- Oh, J.-H., & Seung, H. S. (1997). Learning generative models with the up-propagation algorithm. *Advances in Neural Information Processing Systems* (pp. 605–611). MIT.
- Saul, L. K., Sha, F., & Lee, D. D. (2003). Statistical signal processing with nonnegativity constraints. *Proc. EUROSPEECH* (pp. 1001–1004). Geneva, Switzerland.