
Rectified Factor Networks and Dropout

Djork-Arné Clevert

Thomas Unterthiner

Sepp Hochreiter

Institute of Bioinformatics

Johannes Kepler University, Linz, Austria

{okko, unterthiner, hochreit}@bioinf.jku.at

Abstract

The success of deep learning techniques is based on their robust, effective and abstract representations of the input. In particular, sparse representations that are obtained from rectified linear units and dropout increased classification performance at various tasks. Deep architectures are often constructed by unsupervised pretraining and stacking of either restricted Boltzmann machines (RBMs) or autoencoders. We propose rectified factor networks (RFNs) for pretraining of deep networks. In contrast to RBMs and autoencoders, RFNs (1) estimate the noise of each input component, (2) aim at decorrelating the hidden units (factors), (3) estimate the precision of hidden units by the posterior variance. In the E-step of an EM algorithm, RFN learning (i) enforces non-negative posterior means, (ii) allows dropout of hidden units, and (iii) normalizes the signal part of the hidden units. In the M-step, RFN learning applies gradient descent along the Newton direction to allow rectifying, dropout, and fast GPU implementations. RFN learning can be considered as a variational EM algorithm with unknown prior which is estimated during maximizing the likelihood. Using a fixed point analysis, we show RFNs explain the data variance like factor analysis.

RFNs produce sparse and non-linear input representations for new data by a linear mapping and subsequent rectification, therefore can be readily used for pretraining of deep networks. It is tailored to making full use of large hidden layers with respect to both using all of them to code the input and computational complexity.

We tested and compared RFNs for unsupervised pretraining of deep learning on nine different benchmark datasets: MNIST, basic MNIST, bg-rand MNIST, bg-img MNIST, rect (tall vs. wide rectangles), rect-img, convex (convex vs. concave shapes), NORB and CIFAR-10. Competitors were support vector machines, deep learning with stacking denoising autoencoders, deep learning with stacking regular autoencoders, and deep learning with stacking restricted Boltzmann machines. Only on bg-rand RFN pretraining was significantly worse than the best performing method, though it was second. On basic MNIST, bg-img MNIST, rect-img, rect, convex and NORB, RFNs pretraining outperformed all other methods — in four cases significantly.

1 Introduction

The advent of deep learning [17, 4] and its success in both academic challenges and industrial applications [6, 22, 12, 5] is based on better input representations compared to previous approaches. Input codes at higher levels have more abstract representations that capture complex and non-linear explanatory factors underlying the observed input [3]. Rectified linear units (ReLU) [26, 10, 32] and dropout [18, 31] improved the performance of deep learning methods in several supervised benchmark datasets. ReLU lead to sparse representations [10] which are further sparsified by dropout [1].

Therefore, beyond complex and abstract representation on different levels, sparse representations enhance the performance of subsequent supervised techniques. Sparse representations were originally motivated by findings in sensory neural systems and are well known objectives in machine learning [28, 16, 9]. In bioinformatics sparse codes excelled in biclustering of gene expression data [20] and in finding DNA sharing patterns between humans, Neanderthals and Denisovans [19]. However constructing a sparse code can be computationally expensive [13], while ReLU and dropout are fast and, therefore, are perfectly suited for deep learning with many representational units.

Pretraining deep networks with unsupervised sparse coding methods is very appealing because unlabeled data can be exploited to enhance data representation and to extract robust structures in the data. Conventional sparse coding methods adjust the code of new data and, therefore, are not appropriate for pretraining deep networks, which have weights that are independent of the input. These methods either construct a code for the observed data but not for new data or require an iterative update to compute a code for new data [28, 13]. Also generative models with sparse priors have for each input a different posterior variance which again makes the coding weights dependent on the input [9, 29]. Furthermore, generative models with sparse priors do not ensure sparse posteriors, though the latter represent the inputs. For example, rectified factor analysis, which rectifies Gaussian priors and selects models using a variational Bayesian learning procedure, does not yield sparse posteriors [14, 15].

ReLU and dropout produce sparse codes for the test data without iterative code updates. Currently, ReLU and dropout are used in connection with restricted Boltzmann machines (RBMs) [26, 18] and neural networks [10, 32, 31]. However during unsupervised pretraining, neither RBMs nor neural autoencoders (1) estimate the noise of the visible units, (2) explicitly decorrelate the hidden (code) units, (3) estimate the precision of hidden units. First, without input noise estimates, very noisy inputs are forced to be reconstructed via large errors in autoencoders or large input-hidden co-activations in RBMs. Furthermore, errors of highly predictive inputs with low noise levels should be up-weighted to detect more structures at which they participate. Consequently, input noise estimations can hint at interesting parts in the input. Secondly, many hidden units may capture the same, most dominant input structures while rare or small structures are missed. Decorrelation of hidden units will force more variety into hidden units and helps to detect rare or small input structures. Thirdly, RBMs and autoencoders estimate neither the precision nor the information content of hidden units. Thus, larger activations of the ReLU hidden units need not indicate more information, more precision, or more consistent input structures. We suggest to address issues (1)–(3) by a factor analysis model as an alternative to RBMs and autoencoders.

2 Rectified Factor Network

We propose to use factor analysis with many hidden units (factors) for unsupervised pretraining of deep networks to estimate the noise of each visible unit, to decorrelate the hidden units, and to estimate the information content, precision, and consistency of input structures captured by the hidden units. The factor analysis model is

$$\mathbf{v} = \mathbf{W}\mathbf{h} + \boldsymbol{\epsilon}. \quad (1)$$

The prior $\mathbf{h} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ on the hidden units $\mathbf{h} \in \mathbb{R}^l$, i.e. the factors, and the noise $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Psi})$ on visible units $\mathbf{v} \in \mathbb{R}^m$, i.e. the observations, are independent. The model parameters are the weight (factor loading) matrix $\mathbf{W} \in \mathbb{R}^{m \times l}$ and the noise covariance matrix $\boldsymbol{\Psi} \in \mathbb{R}^{m \times m}$. $\boldsymbol{\Psi}$ is assumed to be diagonal in order to explain correlations between input components by the hidden units (signal) and not by correlated noise. The linear factor analysis model with independent additive Gaussian noise is depicted in Fig. 1.

Given the mean-centered data $\{\mathbf{v}\} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ and the covariance matrix $\mathbf{C} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbf{v}_i^T$, the

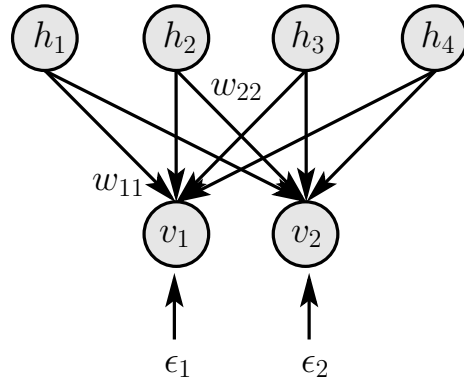


Figure 1: Factor analysis model with hidden units (factors) \mathbf{h} , visible units (observations) \mathbf{v} , weight (factor loading) matrix \mathbf{W} , and noise $\boldsymbol{\epsilon}$.

expectation maximization (EM) update for maximizing the likelihood of the factor analysis model is

$$\begin{aligned}
\text{E-step:} \quad & \mu_{h_i|v_i} = \mathbf{W}^T (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \mathbf{v}_i = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \mathbf{W}^T \Psi^{-1} \mathbf{v}_i, \\
& \Sigma_{h_i|v_i} = \mathbf{I} - \mathbf{W}^T (\mathbf{W} \mathbf{W}^T + \Psi)^{-1} \mathbf{W} = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1}, \\
& \mathbb{E}_{h_i|v_i}(\mathbf{h}_i) = \mu_{h_i|v_i}, \quad \mathbb{E}_{h_i|v_i}(\mathbf{h}_i \mathbf{h}_i^T) = \mu_{h_i|v_i} \mu_{h_i|v_i}^T + \Sigma_{h_i|v_i} \\
& \mathbf{U} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_{h_i|v_i}^T(\mathbf{h}_i), \quad \mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{h_i|v_i}(\mathbf{h}_i \mathbf{h}_i^T) \\
\text{M-step:} \quad & \mathbf{W} = \mathbf{U} \mathbf{S}^{-1}, \quad \Psi = \text{diag}(\mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U}^T + \mathbf{W} \mathbf{S} \mathbf{W}^T).
\end{aligned}$$

In this update rule the hidden unit posterior $p(\mathbf{h}_i | \mathbf{v}_i)$ is Gaussian with mean $\mu_{h_i|v_i}$ and covariance matrix $\Sigma_{h_i|v_i}$. In the EM algorithm, the posterior is computed in the E-step for the old parameters (\mathbf{W}, Ψ) which are then updated in the M-step. \mathbf{U} is the basic update term for weight matrix \mathbf{W} which is essentially Hebb's rule between hidden and visible units. As mentioned in the introduction, factor analysis (1) estimates the noise Ψ of visible units, (2) decorrelates the hidden units by multiplying \mathbf{U} with \mathbf{S}^{-1} , (3) estimates the information content and the precision of hidden units via $\Sigma_{h_i|v_i}$.

We aim to introduce rectifying and dropout of the hidden unit posterior into this EM framework to construct sparse codes. Besides that, we also need the rectifier to create a non-linear input representations, which allows for stacking RFNs, as without any non-linear activation function on the hidden unit every stacked linear model would collapse to one layer.

Rectifying, that is non-negative posterior means, can be implemented by the posterior constraint method [7, 11]. However dropout and normalization, which are essential for our approach, cannot be treated by posterior constraints. Therefore, we modify above EM algorithm with respect to four items:

- (i) **rectifying**: $[\mu_{h|v}]_j = \max\{0, [\mu_{h|v}]_j\}$ to enforce non-negative posterior means.
- (ii) **dropout**: $[\mu_{h|v}]_j = \delta [\mu_{h|v}]_j$ with $\delta \in \{0, 1\}$ as dropout variable.
- (iii) **normalizing**: $\frac{1}{n} \sum_i [\mu_{h_i|v_i}]_j^2 = 1$.
- (iv) **gradient descent**: the M-step is a gradient descent step in the Newton direction.

Rectifying and dropout in steps (i) and (ii), respectively, enforce sparse codes $\mu_{h|v}$. Normalizing the signal part $\mu_{h|v}$ of the hidden units in step (iii) is an important step for models with many hidden units. If the signal-to-noise ratio is varying across samples or if the noise is correlated with the signal, then factor analysis tends to explain signals away by noise. Normalizing the hidden unit's signal part, forces the model to explain correlation among visible units by signal which otherwise would be explained by noise and thereby finds a wider representation of the visible units. Step (iv), the gradient descent update, is important to keep the current solution and to avoid that extreme cases of rectifying or dropout corrupt it. Furthermore, this step is advantageous for using stochastic gradient and fast GPU implementations. The RFN learning algorithm with learning rates η_Ψ and η_W , dropout rate d , and the lower bound Ψ_{\min} on Ψ_{kk} is

$$\text{E-step (1):} \quad \mu_{h|v} = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} \mathbf{W}^T \Psi^{-1} \mathbf{v}, \quad (2)$$

$$\text{Dropout:} \quad \Pr(\delta = 0) = d, \quad \Pr(\delta = 1) = 1 - d$$

$$\text{Rectifier:} \quad [\mu_{h|v}]_j = \delta \max\{0, [\mu_{h|v}]_j\} \quad (3)$$

$$\text{Normalizer:} \quad \frac{1}{n} \sum_i [\mu_{h_i|v_i}]_j^2 = 1$$

$$\text{E-step (2): } \Sigma_{\mathbf{h}|\mathbf{v}} = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1}, \quad (4)$$

$$\mathbb{E}_{\mathbf{h}_i|\mathbf{v}_i}(\mathbf{h}_i) = \boldsymbol{\mu}_{\mathbf{h}_i|\mathbf{v}_i}, \quad \mathbb{E}_{\mathbf{h}_i|\mathbf{v}_i}(\mathbf{h}_i \mathbf{h}_i^T) = \boldsymbol{\mu}_{\mathbf{h}_i|\mathbf{v}_i} \boldsymbol{\mu}_{\mathbf{h}_i|\mathbf{v}_i}^T + \Sigma_{\mathbf{h}_i|\mathbf{v}_i} \quad (5)$$

$$\mathbf{U} = \frac{1}{n} \sum_{i=1}^n \mathbf{v}_i \mathbb{E}_{\mathbf{h}_i|\mathbf{v}_i}^T(\mathbf{h}_i), \quad \mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{h}_i|\mathbf{v}_i}(\mathbf{h}_i \mathbf{h}_i^T) \quad (6)$$

$$\text{M-step: } \mathbf{C} = \frac{1}{n} \sum_{k=1}^n \mathbf{v}_k \mathbf{v}_k^T, \quad c_k = [\mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U}^T + \mathbf{W} \mathbf{S} \mathbf{W}^T]_{kk} \quad (7)$$

$$\mathbf{W} = \mathbf{W} + \eta_W \Delta \mathbf{W}, \quad \Delta \mathbf{W} = \mathbf{U} \mathbf{S}^{-1} - \mathbf{W} \quad (7)$$

$$\Psi_{kk} = \max\{\Psi_{\min}, \Psi_{kk} + \eta_\Psi \Delta \Psi_{kk}\}, \quad \Delta \Psi_{kk} = c_k - \Psi_{kk}. \quad (8)$$

3 Analysis of RFN Learning

In the next subsection RFN learning is formulated as a variational EM algorithm that maximizes the data likelihood for an unknown prior. In the following subsection we show via a fixed point analysis that RFNs explain the data variance like factor analysis. Finally, we derive the Newton update rule of the RFN learning.

3.1 Variational EM

The variational EM framework [8, 27, 21, 2, 9, 29] is utilized to describe RFN learning. We parametrize a Gaussian prior by the variational parameter $\boldsymbol{\xi}$:

$$p(\mathbf{h}; \boldsymbol{\xi}) = \mathcal{N}(\mathbf{h}; \boldsymbol{\xi}, \mathbf{I}) = (2\pi)^{-\frac{1}{2}} \exp\left(-\frac{1}{2} (\mathbf{h} - \boldsymbol{\xi})^T (\mathbf{h} - \boldsymbol{\xi})\right). \quad (9)$$

Using a variational distribution $Q(\mathbf{h} | \mathbf{v})$, the data likelihood can be bounded via the following equation:

$$\begin{aligned} \log p(\mathbf{v}) - D_{\text{KL}}(Q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v})) &= - \int Q(\mathbf{h} | \mathbf{v}) \log \frac{Q(\mathbf{h} | \mathbf{v})}{p(\mathbf{h}, \mathbf{v})} d\mathbf{h} \\ &= \int Q(\mathbf{h} | \mathbf{v}) \log p(\mathbf{v} | \mathbf{h}) d\mathbf{h} - D_{\text{KL}}(Q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h})) = \mathcal{F}(\mathbf{W}, \Psi, \boldsymbol{\xi} | \hat{\mathbf{W}}, \hat{\Psi}, \hat{\boldsymbol{\xi}}), \end{aligned} \quad (10)$$

where D_{KL} denotes the Kullback-Leibler divergence. $(\hat{\mathbf{W}}, \hat{\Psi}, \hat{\boldsymbol{\xi}})$ are the actual parameter estimates that describe $Q(\mathbf{h} | \mathbf{v})$, (\mathbf{W}, Ψ) are the parameters of the posterior $p(\mathbf{v} | \mathbf{h})$, and $\boldsymbol{\xi}$ are the parameters of the prior $p(\mathbf{h})$. The variational EM maximizes \mathcal{F} in its M-step with respect to the parameters and in its E-step with respect to Q . The E-step is equivalent to minimizing $D_{\text{KL}}(Q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}))$, therefore it tightens the lower bound \mathcal{F} on the log-likelihood $\log p(\mathbf{v})$. For the M-step, the objective can be divided into $\int Q(\mathbf{h} | \mathbf{v}) \log p(\mathbf{v} | \mathbf{h}) d\mathbf{h}$ which only depends on (\mathbf{W}, Ψ) and must be maximized and $D_{\text{KL}}(Q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h}))$ which only depends on $\boldsymbol{\xi}$ and must be minimized.

Instead of optimizing Q as in the variational EM, we consider Q as being given via rectifying, dropout, and normalization. However, neither the variational parameter nor the posterior are known. Therefore we minimize $D_{\text{KL}}(Q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}))$ with respect to $\boldsymbol{\xi}$ for the given $Q(\mathbf{h} | \mathbf{v})$. The posterior $p(\mathbf{h} | \mathbf{v})$ is a Gaussian with mean and covariance matrix:

$$\boldsymbol{\mu}_{\mathbf{h}|\mathbf{v}} = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1} (\mathbf{W}^T \Psi^{-1} \mathbf{v} + \boldsymbol{\xi}), \quad \Sigma_{\mathbf{h}|\mathbf{v}} = (\mathbf{I} + \mathbf{W}^T \Psi^{-1} \mathbf{W})^{-1}. \quad (11)$$

Minimizing $D_{\text{KL}}(Q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}))$ gives

$$\boldsymbol{\xi} = (\mathbf{W}^T \Psi^{-1} \mathbf{W} + \mathbf{I}) \boldsymbol{\mu}_Q - \mathbf{W}^T \Psi^{-1} \mathbf{v}, \quad (12)$$

where $\boldsymbol{\mu}_Q$ is the modified $\boldsymbol{\mu}_{\mathbf{h}|\mathbf{v}}$ according to above steps (i)–(iii) (rectifying, dropout, normalizing). $Q(\mathbf{h} | \mathbf{v})$ and $p(\mathbf{h} | \mathbf{v})$ have the same covariance matrix and we enforced $\boldsymbol{\mu}_{\mathbf{h}|\mathbf{v}} = \boldsymbol{\mu}_Q$, thus the Kullback-Leibler divergence between them is zero. Therefore $\boldsymbol{\xi}$ makes the bound \mathcal{F} on $\log p(\mathbf{v})$

tight and can even be used to estimate the model prior. The performance of the model in explaining the data depends on the amount of information in \mathbf{v} that is conveyed to $\boldsymbol{\mu}_Q$. Dropout is covered by this framework if $\boldsymbol{\xi}$ is viewed as a random variable. Consequently, RFN learning maximizes the likelihood for a prior that is constructed during learning. Besides the architecture, the complexity of the model class depends on how much input information is conveyed to the hidden units after the posterior modification.

3.2 Fixed Point

For rectified factor networks (RFNs) the objective is unknown during learning because the prior and, therefore, also the likelihood are not given. We derive a fixed point of the algorithm which is independent of the posterior modification. The fixed point solution explains the variation in the data like factor analysis.

The fixed point equation for the \mathbf{W} update is

$$\Delta \mathbf{W} = \mathbf{U} \mathbf{S}^{-1} - \mathbf{W} = \mathbf{0} \Rightarrow \mathbf{U} - \mathbf{W} \mathbf{S} = \mathbf{0} \Rightarrow \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{S} \mathbf{W}^T = \mathbf{0}, \quad (13)$$

where for the last equation we multiplied the previous from the right hand side by \mathbf{W}^T . The fixed point equation for the full (not only diagonal) Ψ update is

$$\Psi = \mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U}^T + \mathbf{W} \mathbf{S} \mathbf{W}^T = \mathbf{C} - \mathbf{W} \mathbf{U}^T, \quad (14)$$

where we inserted Eq. (13). Since $\mathbf{W} \mathbf{S} \mathbf{W}^T$ in Eq. (13) is symmetric, we can insert $\mathbf{W} \mathbf{U}^T = \mathbf{U} \mathbf{W}^T = \mathbf{C} - \Psi$ in last Eq. (13):

$$\mathbf{C} = \Psi + \mathbf{W} \mathbf{S} \mathbf{W}^T. \quad (15)$$

Therefore the model that corresponds to the fixed point explains the data covariance matrix \mathbf{C} by a noise part Ψ and a signal part $\mathbf{W} \mathbf{S} \mathbf{W}^T$. Like factor analysis the data variance is explained by the model via the parameters Ψ (noise) and \mathbf{W} (signal). $\Psi = \frac{1}{n} \sum_{i=1}^n \epsilon_i \epsilon_i^T + \mathbf{W} \Sigma_{\mathbf{h}_i | \mathbf{v}_i} \mathbf{W}^T$ explains both reconstruction errors $\epsilon_i = \mathbf{v}_i - \boldsymbol{\mu}_{\mathbf{h}_i | \mathbf{v}_i} \mathbf{W}$ and uncertainties of hidden units captured by $\Sigma_{\mathbf{h}_i | \mathbf{v}_i}$. If $\boldsymbol{\mu}_{\mathbf{h}_i | \mathbf{v}_i}$ contains sufficient information on the visible units \mathbf{v}_i , then the fixed point equation holds also for restricting Ψ to a diagonal matrix. In this case the reconstruction errors approach zero (forced by the M-step) and $\mathbf{W} \Sigma_{\mathbf{h}_i | \mathbf{v}_i} \mathbf{W}^T$ is diagonal dominant and explained by Ψ . Dropout is covered by this fixed point analysis. For dropout $\boldsymbol{\mu}_{\mathbf{h}_i | \mathbf{v}_i}$ become random variables, thus \mathbf{U} and \mathbf{S} are replaced by their expectations. How much data variation is explained by signal $\mathbf{W} \mathbf{S} \mathbf{W}^T$ and how much is explained by noise Ψ depends on (a) the model complexity like the number of factors, (b) the measurement noise on the visible units, and (c) the amount of information about the visible units \mathbf{v} that is coded in the hidden units $\boldsymbol{\mu}_{\mathbf{h} | \mathbf{v}}$.

3.3 Newton Direction

We now derive the Newton update rule and address step (iv) from above. The objective to maximize is

$$\begin{aligned} \mathcal{L} = \int_{\mathbb{R}^l} Q_i(\mathbf{h}_i | \mathbf{v}_i) \log(p(\mathbf{v}_i | \mathbf{h}_i; \mathbf{W}, \Psi)) d\mathbf{h}_i, \quad \log \mathcal{L} = -\frac{m n}{2} \log(2\pi) - \frac{n}{2} \log |\Psi| \\ - \frac{1}{2} \sum_{i=1}^n \mathbf{v}_i^T \Psi^{-1} \mathbf{v}_i + \sum_{i=1}^n \mathbb{E}_{\mathbf{h}_i | \mathbf{v}_i} (\mathbf{v}_i^T \Psi^{-1} \mathbf{W} \mathbf{h}_i) \frac{1}{2} \mathbb{E}_{\mathbf{h}_i | \mathbf{v}_i} \left(\sum_{i=1}^n \mathbf{h}_i^T \mathbf{W}^T \Psi^{-1} \mathbf{W} \mathbf{h}_i \right) \end{aligned} \quad (16)$$

The Hessian $\mathbf{H}_{\mathbf{W}}$ of $(\frac{2}{n} \log \mathcal{L})$ with respect to \mathbf{W} as a vector is:

$$\mathbf{H}_{\mathbf{W}} = \frac{\partial \text{vec}(\frac{2}{n} \nabla_{\mathbf{W}} \log \mathcal{L})}{\partial \text{vec}(\mathbf{W})^T} = \frac{\partial \text{vec}(\Psi^{-1} \mathbf{U} - \Psi^{-1} \mathbf{W} \mathbf{S})}{\partial \text{vec}(\mathbf{W})^T} = -\mathbf{S} \otimes \Psi^{-1}, \quad (17)$$

where \otimes is the Kronecker product of matrices. For the product of the negative inverse Hessian with the gradient we have:

$$\begin{aligned} -\mathbf{H}_{\mathbf{W}}^{-1} \text{vec}(\Psi^{-1} \mathbf{U} - \Psi^{-1} \mathbf{W} \mathbf{S}) &= (\mathbf{S}^{-1} \otimes \Psi) \text{vec}(\Psi^{-1} \mathbf{U} - \Psi^{-1} \mathbf{W} \mathbf{S}) \\ &= \text{vec}(\Psi (\Psi^{-1} \mathbf{U} - \Psi^{-1} \mathbf{W} \mathbf{S}) \mathbf{S}^{-1}) = \text{vec}(\mathbf{U} \mathbf{S}^{-1} - \mathbf{W}) . \end{aligned} \quad (18)$$

Thus, if we apply a Newton update then the update direction for \mathbf{W} in the M-step is

$$\Delta \mathbf{W} = \mathbf{U} \mathbf{S}^{-1} - \mathbf{W} . \quad (19)$$

This is the exact EM update if the stepsize is 1. Since the objective is a quadratic function in \mathbf{W} , one Newton update would lead to the exact solution. However for rectified linear posteriors and for dropout we want to have small changes of the matrix \mathbf{W} . The Newton update converges fast, but comes at the cost of inverting \mathbf{S} , which is computational expensive, thus for large \mathbf{S} matrices $\Delta \mathbf{W} = \mathbf{U} - \mathbf{W} \mathbf{S}$ is a reasonable alternative update rule.

The Hessian \mathbf{H}_{Ψ} of $(\frac{2}{n} \log \mathcal{L})$ with respect to Ψ as a vector is a diagonal matrix with $[H_{\Psi}]_{ij} = 0$ for $i \neq j$ and for $i = j$:

$$[H_{\Psi}]_{ii} = \frac{1}{\Psi_{ii}^2} - \frac{2 c_i}{\Psi_{ii}^3} = \frac{1}{\Psi_{ii}^2} \left(1 - \frac{2 c_i}{\Psi_{ii}} \right) , \quad c_i = [\mathbf{C} - \mathbf{U} \mathbf{W}^T - \mathbf{W} \mathbf{U}^T + \mathbf{W} \mathbf{S} \mathbf{W}^T]_{ii} \quad (20)$$

Since we maximize the bound, we have to ensure that the Hessian is negative definite, that is, $[H_{\Psi}]_{ii} < 0$ or $\Psi_{ii} < 2 c_i$. Therefore we replace c_i by Ψ_{ii} and get $[H_{\Psi}]_{ii} = -1/\Psi_{ii}^2$. The Newton update direction is:

$$\Delta \Psi_{ii} = - \frac{[\nabla_{\Psi} \log \mathcal{L}]_{ii}}{[H_{\Psi}]_{ii}} = c_i - \Psi_{ii} . \quad (21)$$

This is the EM update for learning rate equal to 1.

4 Experiments

In this section, we assess the performance of rectified factor networks (RFNs) as a pretraining procedure to construct deep networks. We stacked RFNs in the same way as described by Vincent et al. [30], namely first training a single layer RFN and then passing the resulting representation as input for training the next RFN.

We conducted experiments with two deep network architectures constructed from pretrained RFNs: a 1-hidden layer network (RFN-1) and a 3-hidden layer network. The classification performance of the RFN pretrained deep networks was compared to (i) support vector machines, (ii) deep networks pretrained by stacking denoising autoencoders (SDAE), (iii) stacking regular autoencoders (SAE), (iv) restricted Boltzmann machines (RBM), and (v) stacking restricted Boltzmann machines (DBN).

The benchmark datasets are from previous publications [25, 30, 24, 23] and contain: (i) *MNIST* (original MNIST), (ii) *basic* (a smaller subset of MNIST for training), (iii) *bg-rand* (MNIST digits with random noise background), (iv) *bg-img* (MNIST digits with random image background), (v) *rect* (discrimination between tall and wide rectangles), (vi) *rect-img* (discrimination between tall and wide rectangular images overlayed on different background images), (vii) *convex* (discrimination between convex and concave shapes), (viii) *CIFAR-10* (60k 32x32 colour images in 10 classes, with 6k images per class), and (ix) *NORB* (29,160 stereo image pairs of 50 toys belonging to 5 generic categories). The dataset characteristics in terms of size of training, validation and test set is given in the second column of Tab. 1. In all experiments we apply no further data preprocessing except median centering and learn the representation in an unsupervised fashion. We followed Vincent et al. [30] and selected the models based on the validation set performance. The RFNs have as possible hyperparameters: (i) the number of units in the layers from $\{1024, 2048, 4096\}$ and (ii) the dropout rate from $\{0.0, 0.25, 0.5, 0.75\}$. The learning rates were fixed to $\eta_W = 0.1$ and $\eta_{\Psi} = 0.01$ which we found to be well suited for RFN learning on separate artificial toy datasets, where the task was to find biclusters [20]. For supervised fine-tuning with stochastic gradient descent, we selected the learning rate from $\{0.1, 0.01, 0.001\}$, the masking noise from $\{0.0, 0.25\}$, and the number of layers from $\{1, 3\}$. Again following Vincent et al. [30], fine-tuning was stopped early, where the stopping time was selected based on the validation set performance.

The results of the comparison of deep networks pretrained with RFNs and other models are given in Tab. 1, while Fig. 2 shows some examples of the learnt filters. The test error rate is reported for seven classification problems. A 95% confidence interval is computed according to Vincent et al. [30]. The result of the best performing method is given in bold, as well as the result of those

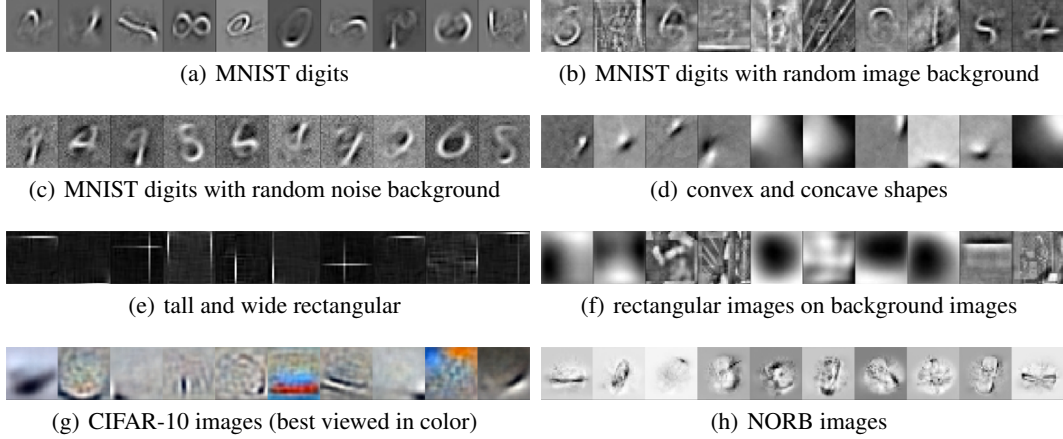


Figure 2: Various randomly selected filters trained on benchmark data sets using a RFN with 1024 hidden units. The panels (a)-(c) show filters for MNIST variations, while (d) and (e)-(f) show filters for convex and rectangle, respectively. Filters learnt from CIFAR-10 and NORB are shown in panel (g)-(h). RFN learnt various kinds of filters, such as stroke, local and global blob detectors. RFN extracts complex features and as in panel (c) can be seen, are the RFN filters unaffected by the background noise.

Table 1: Comparison of deep networks pretrained with RFNs and other models. Test error rate on all considered classification problems is reported together with a 95% confidence interval. The result of the best performing method is given in bold, as well as those for which confidence intervals overlap. The first column gives the data set, the second the size of training, validation and test set, the last column indicates the number of hidden layers of the RFN pretrained deep network which is chosen according to the validation set performance. In only one case RFN pretraining was significantly worse than the best method but still the second best method. In six out of the nine experiments RFN pretraining performed best, where in four cases it was significantly better than all other methods.

Dataset		SVM	RBM	DBN	SAE	SDAE	RFN
MNIST	50k-10k-10k	1.40 \pm 0.23	1.21 \pm 0.21	1.24 \pm 0.22	1.40 \pm 0.23	1.28 \pm 0.22	1.27 \pm 0.22 (1)
basic	10k-2k-50k	3.03 \pm 0.15	3.94 \pm 0.17	3.11 \pm 0.15	3.46 \pm 0.16	2.84 \pm 0.15	2.66 \pm 0.14 (1)
bg-rand	10k-2k-50k	14.58 \pm 0.31	9.80 \pm 0.26	6.73 \pm 0.22	11.28 \pm 0.28	10.30 \pm 0.27	7.94 \pm 0.24 (3)
bg-img	10k-2k-50k	22.61 \pm 0.37	16.15 \pm 0.32	16.31 \pm 0.32	23.00 \pm 0.37	16.68 \pm 0.33	15.66 \pm 0.32 (1)
rect	1k-0.2k-50k	2.15 \pm 0.13	4.71 \pm 0.19	2.60 \pm 0.14	2.41 \pm 0.13	1.99 \pm 0.12	0.63 \pm 0.06 (1)
rect-img	10k-2k-50k	24.04 \pm 0.37	23.69 \pm 0.37	22.50 \pm 0.37	24.05 \pm 0.37	21.59 \pm 0.36	20.77 \pm 0.36 (1)
convex	10k-2k-50k	19.13 \pm 0.34	19.92 \pm 0.35	18.63 \pm 0.34	18.41 \pm 0.34	19.06 \pm 0.34	16.41 \pm 0.32 (1)
NORB	19k-5k-24k	11.6 \pm 0.40	8.31 \pm 0.35	-	10.10 \pm 0.38	9.50 \pm 0.37	7.00 \pm 0.32 (1)
CIFAR	40k-10k-10k	62.7 \pm 0.95	40.39 \pm 0.96	43.38 \pm 0.97	43.25 \pm 0.97	-	41.29 \pm 0.95 (1)

methods for which confidence intervals overlap. In only one case RFN pretraining was significantly worse than the best method but still the second best method. In six out of the nine experiments RFN pretraining performed best, where in four cases it was significantly better than all other methods.

5 Conclusion

We have introduced rectified factor networks (RFNs) for pretraining deep networks. RFNs produce sparse and non-linear input representations of the data on different levels and extract factors underlying the data. RFNs are tailored to pretraining of deep networks and to many hidden units. We have shown that RFN learning is a variational EM algorithm with unknown prior and that the fixed point of RFN training explains the data variance like factor analysis.

On nine different benchmark datasets we compare deep networks with RFN pretraining with support vector machines and deep networks pretrained with regular / denoising autoencoders and restricted Boltzmann machines. Only on one dataset RFN pretraining was significantly worse than the best performing method. On six datasets RFNs pretraining outperformed all other methods — in four cases significantly.

RFNs are an interesting alternative to RBMs and autoencoders in the context of deep learning. RFNs have high potential as unsupervised methods for semi-supervised learning and producing sparse codes as they are geared to large datasets and many representational units.

Acknowledgment

The Tesla K40 used for this research was donated by the NVIDIA Corporation.

References

- [1] P. Baldi and P. Sadowski. The dropout learning algorithm. *Artificial Intelligence*, 210C:78–122, 2014.
- [2] M. J. Beal. *Variational algorithms for approximate Bayesian inference*. PhD thesis, University College London, The Gatsby Computational Neuroscience Unit, 17 Queen Square London WC1N 3AR, 2003.
- [3] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013.
- [4] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle. Greedy layer-wise training of deep networks. In *Advances in Neural Information Processing Systems 19 (NIPS 2006)*, pages 153–160. MIT Press, 2007.
- [5] D. C. Ciresan, U. Meier, and J. Schmidhuber. Multi-column deep neural networks for image classification. In *IEEE Conference on Computer Vision and Pattern Recognition CVPR 2012*, 2012. Long preprint arXiv:1202.2745v1 [cs.CV].
- [6] G. E. Dahl, D. Yu, L. Deng, and A. Acero. Context-dependent pre-trained deep neural networks for large vocabulary speech recognition. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(1):33–42, 2012.
- [7] K. Ganchev, J. Graca, J. Gillenwater, and B. Taskar. Posterior regularization for structured latent variable models. *Journal of Machine Learning Research*, 11:2001–2049, 2010.
- [8] Z. Ghahramani and G. E. Hinton. Variational learning for switching state-space models. *Neural Computation*, 12:963–996, 1998.
- [9] M. Girolami. A variational method for learning sparse and overcomplete representations. *Neural Comput.*, 13(11):2517–2532, 2001.
- [10] X. Glorot, A. Bordes, and Y. Bengio. Deep sparse rectifier neural networks. In G. Gordon, D. Dunson, and M. Dudk, editors, *JMLR W&CP: Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics (AISTATS 2011)*, volume 15, pages 315–323, 2011.
- [11] J. V. Graca, K. Ganchev, and B. Taskar. Expectation maximization and posterior constraints. In J.C. Platt, D. Koller, Y. Singer, and S.T. Roweis, editors, *Advances in Neural Information Processing Systems*, volume 20, pages 569–576, 2007.
- [12] A. Graves, A.-R. Mohamed, and G. E. Hinton. Speech recognition with deep recurrent neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pages 6645–6649. IEEE, 2013.
- [13] K. Gregor and Y. LeCun. Learning fast approximations of sparse coding. In J. Fürnkranz and T. Joachims, editors, *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 399–406. Omnipress 2010, ISBN 978-1-60558-907-7, 2010.
- [14] M. Harva and A. Kaban. A variational bayesian method for rectified factor analysis. In *Proc. Int. Joint Conf. on Neural Networks (IJCNN’05)*, pages 185–190, 2005.
- [15] M. Harva and A. Kaban. Variational learning for rectified factor analysis. 2006.
- [16] G. E. Hinton and Z. Ghahramani. Generative models for discovering sparse distributed representations. *Philos. T. R. Soc. B*, 352:1177–1190, 1997.
- [17] G. E. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- [18] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *CoRR*, abs/1207.0580, 2012.

- [19] S. Hochreiter. HapFABIA: Identification of very short segments of identity by descent characterized by rare variants in large sequencing data. *Nucleic Acids Res.*, 41(22):e202, 2013.
- [20] S. Hochreiter, U. Bodenhofer, M. Heusel, A. Mayr, A. Mitterecker, A. Kasim, S. VanSanden, D. Lin, W. Talloen, L. Bijmens, et al. FABIA: factor analysis for bicluster acquisition. *Bioinformatics*, 26(12):1520–1527, 2010.
- [21] M. I. Jordan, Z. Ghahramani, T. S. Jaakkola, and L. K. Saul. An introduction to variational methods for graphical models. *Machine Learning*, 37:183–233, 1999.
- [22] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS 2012)*, page 4, 2012.
- [23] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, Master’s thesis, Dept. of Comp. Sci., University of Toronto, 2009.
- [24] H. Larochelle, D. Erhan, A. Courville, J. Bergstra, and Y. Bengio. An empirical evaluation of deep architectures on problems with many factors of variation. In *Proceedings of the 24th International Conference on Machine Learning*, ICML ’07, pages 473–480. ACM, 2007.
- [25] Yann LeCun, Fu-Jie Huang, and Leon Bottou. Learning methods for generic object recognition with invariance to pose and lighting. In *Proceedings of CVPR’04*. IEEE Press, 2004.
- [26] V. Nair and G. E. Hinton. Rectified linear units improve restricted Boltzmann machines. In J. Fürnkranz and T. Joachims, editors, *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 807–814. Omnipress 2010, ISBN 978-1-60558-907-7, 2010.
- [27] R. Neal and G. E. Hinton. A view of the EM algorithm that justifies incremental, sparse, and other variants. In M. I. Jordan, editor, *Learning in Graphical Models*, pages 355–368. MIT Press, Cambridge, MA, 1998.
- [28] B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381:607–609, 1996.
- [29] J. Palmer, D. Wipf, K. Kreutz-Delgado, and B. Rao. Variational EM algorithms for non-Gaussian latent variable models. In *Advances in Neural Information Processing Systems*, volume 18, pages 1059–1066, 2006.
- [30] P. Vincent, H. Larochelle, I. Lajoie, Y. Bengio, and P.-A. Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, 11:3371–3408, 2010.
- [31] D. Warde-Farley, I. J. Goodfellow, A. Courville, and Y. Bengio. An empirical analysis of dropout in piecewise linear networks. *ArXiv e-prints*, 2013.
- [32] M. D. Zeiler, M. Ranzato, R. Monga, M. Mao, K. Yang, Q. V. Le, P. Nguyen, A. Senior, V. Vanhoucke, J. Dean, and G. E. Hinton. On rectified linear units for speech processing. In *38th International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3517–3521. IEEE, 2013.