# Deep Learning Made Easier
# by Linear Transformations in Perceptrons

**Tapani Raiko**
Aalto University School of Science
Dept. of Information and Computer Science
Espoo, Finland
`firstname.lastname@aalto.fi`

**Harri Valpola**
Aalto University School of Science
Dept. of Information and Computer Science
Espoo, Finland
`firstname.lastname@aalto.fi`

**Yann LeCun**
New York University
Courant Institute of Mathematical Sciences
New York, NY
`firstname@cs.nyu.edu`

## Abstract

We transform the outputs of each hidden neuron in a multi-layer perceptron network to be zero mean and zero slope, and use separate shortcut connections to model the linear dependencies instead. This transformation aims at separating the problems of learning the linear and nonlinear parts of the whole input-output mapping, which has many benefits. We study the theoretical properties of the transformation by noting that they make the Fisher information matrix closer to a diagonal matrix, and thus standard gradient closer to the natural gradient. We experimentally confirm the usefulness of the transformations by noting that they make basic stochastic gradient learning competitive with state-of-the-art learning algorithms in speed, and that they seem also to help find solutions that generalize better. The experiments include both classification of handwritten digits with a 3-layer network and learning a low-dimensional representation for images by using a 6-layer auto-encoder network. The transformations were beneficial in all cases, with and without regularization.

## 1 Introduction

Learning deep neural networks has become a popular topic since the invention of unsupervised pretraining [5]. Some later works have returned to traditional back-propagation learning and noticed that it can also provide impressive results given either a sophisticated learning algorithm [10] or simply enough computational power [3]. In this work we study back-propagation learning in deep networks with up to five hidden layers.

In learning multi-layer perceptron (MLP) networks by back-propagation, there are known transformations that speed up learning [9]. For instance, inputs are recommended to be centered to zero mean (or even whitened), and nonlinear functions are proposed to have a range from -1 to 1 rather than 0 to 1 [9]. Schraudolph [12] proposed a transformation of gradients of the error signals to have zero slope on average. In this paper we achieve equivalent transformations to the gradient by transforming the nonlinearities in the hidden units.

It is well known that second-order optimization methods such as the natural gradient [1] or Newton's method decrease the number of required iterations compared to the basic gradient descent, but they cannot be easily used with high-dimensional models in practice due to heavy computations with

large matrices. In practice, it is possible to use a diagonal or block-diagonal approximation [7] of the Fisher information matrix. If we will be approximating most of the matrix with zeros anyway, we should use transformations that make those elements to be as close to zero as possible.

## 2 Proposed Transformations

Let us study an MLP-network with a single hidden layer[1] and a shortcut mapping, that is, the output column vectors $\mathbf{y}_t$ for each sample $t$ are modelled as a function of the input column vectors $\mathbf{x}_t$ with

$$\mathbf{y}_t = \mathbf{A}\mathbf{f}\left(\mathbf{B}\mathbf{x}_t\right) + \mathbf{C}\mathbf{x}_t + \boldsymbol{\epsilon}_t, \tag{1}$$

where $\mathbf{f}$ is a nonlinearity (such as $\tanh$) applied to each component of the vector separately, $\mathbf{A}$, $\mathbf{B}$, and $\mathbf{C}$ are the weight matrices, and $\boldsymbol{\epsilon}_t$ is the noise which is assumed to be zero mean and Gaussian, that is, $p(\epsilon_{it}) = \mathcal{N}\left(\epsilon_{it}; 0, \sigma_i^2\right)$. In order to avoid separate bias vectors that complicate formulas, the input vectors are assumed to have been supplemented with an additional component that is always one.

Let us supplement the $\tanh$ nonlinearity with auxiliary scalar variables $\alpha_i$ and $\beta_i$ for each nonlinearity $f_i$. They are not learnt, but instead they will be set in a manner to help learn the other parameters. We define

$$f_i(\mathbf{b}_i\mathbf{x}_t) = \tanh(\mathbf{b}_i\mathbf{x}_t) + \alpha_i\mathbf{b}_i\mathbf{x}_t + \beta_i. \tag{2}$$

An example can be seen in Figure 1. We will ensure that

$$\sum_{t=1}^{T} f_i(\mathbf{b}_i\mathbf{x}_t) = 0 \qquad\qquad \sum_{t=1}^{T} f_i'(\mathbf{b}_i\mathbf{x}_t) = 0 \tag{3}$$

by setting $\alpha_i$ and $\beta_i$ to

$$\alpha_i = -\frac{1}{T}\sum_{t=1}^{T} \tanh'(\mathbf{b}_i\mathbf{x}_t) \qquad\qquad \beta_i = -\frac{1}{T}\sum_{t=1}^{T} \left[\tanh(\mathbf{b}_i\mathbf{x}_t) + \alpha_i\mathbf{b}_i\mathbf{x}_t\right]. \tag{4}$$

Schraudolph [12] proposed transformations to the gradient that has the same effect as Equation (3), which was shown to speed up learning significantly.

The effect of a change of the auxiliary parameters $\alpha$ and $\beta$ can be compensated exactly by updating the shortcut mapping $\mathbf{C}$ accordingly by

$$\begin{aligned}\mathbf{C}_{\text{new}} = \mathbf{C}_{\text{old}} &- \mathbf{A}(\boldsymbol{\alpha}_{\text{new}} - \boldsymbol{\alpha}_{\text{old}})\mathbf{B} \\ &- \mathbf{A}(\boldsymbol{\beta}_{\text{new}} - \boldsymbol{\beta}_{\text{old}})\left[0 \ \ 0 \ldots 1\right],\end{aligned} \tag{5}$$

where $\boldsymbol{\alpha}$ is a matrix with elements $\alpha_i$ on the diagonal and one empty row below for the bias term, and $\boldsymbol{\beta}$ is a column vector with components $\beta_i$ and one zero below for the bias term.

We also emphasize making the inputs $x_k$ zero mean (and similar in scale) as a preprocessing step (see e.g. [9]).

We expect that the linear part will be easier to learn than the nonlinear part. When learning the nonlinear part, without the transformations, the linear part would have to adapt to the changes. By using the proposed automatic transformations we will make learning of the nonlinear part less dependent on the linear part and also different hidden neurons less dependent on each other.

## 3 Intuitive Justification

Second-order optimization methods such as the natural gradient [1] or Newton's method decrease the number of required iterations compared to the basic gradient descent, but they cannot be easily used with large models in practice due to heavy computations with large matrices. Using basic gradient descent can be seen as approximating the Fisher information matrix with a unit matrix. We will see how the proposed transformations changes at least the non-diagonal elements to be as close to zero as possible.

---

[1]The assumption is done for notational simplicity only, the method is applied in the general deep case.

The Fisher information matrix contains elements

$$G_{ij} = \sum_t \left\langle \frac{\partial^2 \log p(\mathbf{y}_t \mid \mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{x}_t)}{\partial \theta_i \partial \theta_j} \right\rangle, \tag{6}$$

where $\langle \cdot \rangle$ is the expectation over the Gaussian distribution of noise $\epsilon_t$ in Equation (1), and vector $\boldsymbol{\theta}$ contains all the elements of matrices $\mathbf{A}$, $\mathbf{B}$, and $\mathbf{C}$. These elements are:

$$\frac{\partial}{\partial a_{ij}} \frac{\partial}{\partial a_{i'j'}} \log p = \begin{cases} 0 & i' \neq i \\ -\frac{1}{\sigma_i^2} \sum_t f_j(\mathbf{b}_j \mathbf{x}_t) f_{j'}(\mathbf{b}_{j'} \mathbf{x}_t) & i' = i, \end{cases} \tag{7}$$

where $a_{ij}$ is the $ij$th element of matrix $\mathbf{A}$, $f_j$ is the $j$th nonlinearity, and $\mathbf{b}_j$ is the $j$th row vector of matrix $\mathbf{B}$. Similarly

$$\frac{\partial}{\partial b_{jk}} \frac{\partial}{\partial b_{j'k'}} \log p = -\sum_i \frac{1}{\sigma_i^2} a_{ij} a_{ij'} \sum_t f'_j(\mathbf{b}_j \mathbf{x}_t) f'_{j'}(\mathbf{b}_{j'} \mathbf{x}_t) x_{kt} x_{k't} \tag{8}$$

and

$$\frac{\partial}{\partial c_{ik}} \frac{\partial}{\partial c_{i'k'}} \log p = \begin{cases} 0 & i' \neq i \\ -\frac{1}{\sigma_i^2} \sum_t x_{kt} x_{k't} & i' = i. \end{cases} \tag{9}$$

The cross terms are

$$\frac{\partial}{\partial a_{ij}} \frac{\partial}{\partial b_{j'k}} \log p = -\frac{1}{\sigma_i^2} a_{ij'} \sum_t f_j(\mathbf{b}_j \mathbf{x}_t) f'_{j'}(\mathbf{b}_{j'} \mathbf{x}_t) x_{kt} \tag{10}$$

$$\frac{\partial}{\partial c_{ik}} \frac{\partial}{\partial a_{i'j}} \log p = \begin{cases} 0 & i' \neq i \\ -\frac{1}{\sigma_i^2} \sum_t f_j(\mathbf{b}_j \mathbf{x}_t) x_{kt} & i' = i \end{cases} \tag{11}$$

$$\frac{\partial}{\partial c_{ik}} \frac{\partial}{\partial b_{jk'}} \log p = -\frac{1}{\sigma_i^2} a_{ij} \sum_t f'_j(\mathbf{b}_j \mathbf{x}_t) x_{kt} x_{k't}. \tag{12}$$

Now we can notice that Equations (7–12) contain factors such as $f_i(\cdot)$, $f'_i(\cdot)$, and $x_{it}$. By making them as close to zero as possible, we help in making nondiagonal terms of the Fisher information closer to zero, and thus, the basic gradient descent direction closer to the natural gradient. For instance, $E[f_i(\cdot)f_j(\cdot)] = E[f_i(\cdot)]E[f_j(\cdot)] + \text{Cov}[f_i(\cdot), f_j(\cdot)]$, so assuming that the hidden units $i$ and $j$ are representing different things, that is, $f_i(\cdot)$ and $f_j(\cdot)$ are uncorrelated, the nondiagonal element of the Fisher information in Equation (7) becomes zero by using the transformation. The same applies to all other elements in Equations (8–12), some of them also highlighting the benefit of making the input data $\mathbf{x}$ zero-mean.

### 3.1 Positive Side Effect

Having a non-zero $\alpha_i$ has a positive side effect of reducing plateaus in learning. Typical nonlinearities like the $\tanh$ function saturate exponentially on positive and negative sides. When the derivative of an activation $f'_i(\cdot)$ is about zero for most of the data samples, the gradient propagated through it also becomes almost zero, and learning can proceed very slowly or even seem to stop completely. This may explain plateaus in typical learning curves, where the learning proceeds slowly at times. To alleviate the problem, Glorot and Bengio [4] suggested to use the soft-sign nonlinearity that saturates more slowly, but having a non-zero $\alpha_i$ provides a nonlinearity that does not saturate at all. The difference is illustrated in Figure 1.

## 4 Experimental Protocol

There are many practical issues when learning MLP networks and they are addressed below.

**Learning:** Backpropagation learning is basically a gradient ascent algorithm to maximize the log likelihood $\log p(\{\mathbf{y}_t\}_{t=1}^T \mid \{\mathbf{x}_t\}_{t=1}^T, \boldsymbol{\theta})$ of the parameter vector $\boldsymbol{\theta}$, where the back propagation corresponds to using the chain rule of derivatives and dynamic programming to compute the gradient efficiently.
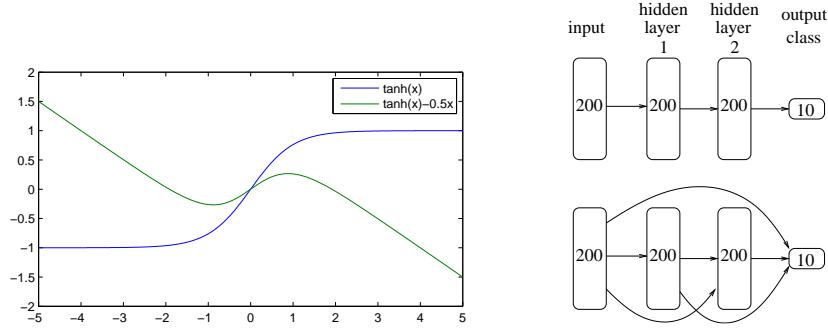
Figure 1: Left: As a positive side effect, the nonlinearity in Equation (2) does not saturate at all for example with a typical $\alpha = -0.5$ and $\beta = 0$. Right top: Traditional structure for a feed-forward multi-layer perceptron network with full connections. Right bottom: Network with shortcut connections included, also used in the proposed method with transformations.

The gradient is

$$g_i = \frac{\partial \langle \log p(\mathbf{y}_t \mid \mathbf{x}_t, \boldsymbol{\theta}) \rangle}{\partial \theta_i}, \tag{13}$$

where $\langle \cdot \rangle$ is the expectation over the data set. The update is

$$\theta_i \leftarrow \theta_i + \gamma g_i, \tag{14}$$

where $\gamma_i$ is a learning rate.

**Online Learning:** It is well known (see e.g. [2]) that looking at all the available data before each update is wasteful, because many samples possess redundant information. We will use a mini-batch learning algorithm, where each update is done based on the 1000 next samples. To reduce the effect of noise due to such a small sample, a momentum term is used. The update direction is thus set to

$$g_i \leftarrow 0.1 \frac{\partial \langle \log p(\mathbf{y}_t \mid \mathbf{x}_t, \boldsymbol{\theta}) \rangle}{\partial \theta_i} + 0.9 g_i. \tag{15}$$

The transformations parameters $\alpha_i$ and $\beta_i$, however are updated after the initialization and after every 1000 iterations thereafter, using the whole data set in Equation (4). At the same time, the changes are compensated by updating the shortcut mappings according to Equation (5) and the momentum $g_i$ for the gradient updates is reset to zero.

**Discrete Outputs:** In classification problems, the output $y_t$ is discrete and one can use the soft-max model. The Equation (1) is replaced by

$$P(y_t = i \mid \mathbf{x}_t, \boldsymbol{\theta}) = \frac{\exp \left[ \mathbf{A}_i \mathbf{f} \left( \mathbf{B} \mathbf{x}_t \right) + \mathbf{C}_i \mathbf{x}_t \right]}{\sum_j \exp \left[ \mathbf{A}_j \mathbf{f} \left( \mathbf{B} \mathbf{x}_t \right) + \mathbf{C}_j \mathbf{x}_t \right]}, \tag{16}$$

where $\mathbf{A}_i$ is the $i$th row of matrix $\mathbf{A}$. Backpropagation is done as before.

**Multiple Hidden Layers:** The extension of the proposed approach to multiple hidden layers is straightforward. Equations (2–4) apply to all nonlinear units with $\mathbf{b}_i \mathbf{x}_t$ replaced by the appropriate input signal. The shortcut mappings need to be included for skipping any number of layers (see Figure 1 for an example). The number of weights of course increases quadratically with the number of layers, but this can be avoided by including a layer without transformations as shown in Figure 5.

**Initialization:** We use the initialization proposed in [4] for weights between consecutive layers, that is, the weight is drawn from a uniform distribution between $\pm \sqrt{6}/\sqrt{n_j + n_{j+1}}$ where $n_j$ is the number of neurons on the $j$th layer. This normalized initialization is based on the objectives of maintaining activation variances and back-propagated gradient's variance throughout the layers. With unnormalized initialization, the gradient tends to vanish or explode exponentially with the number of layers [4]. Biases are drawn from a uniform distribution between $\pm 0.5$. Weights for all shortcut connections are initialized to zero.
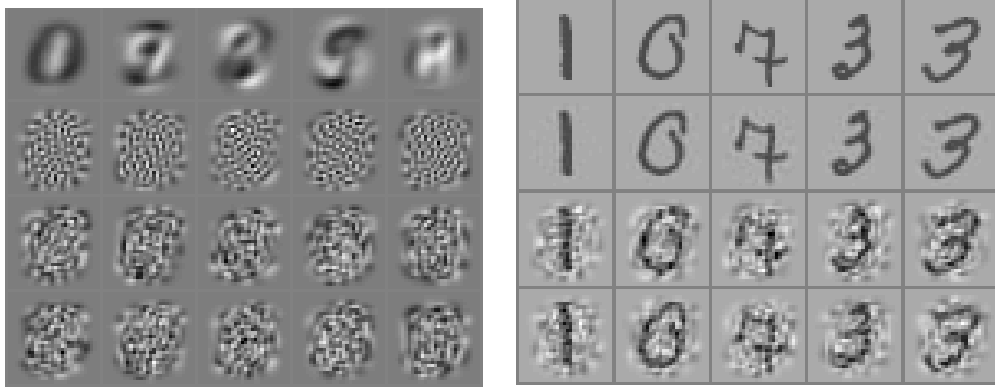
4

Figure 2: Left: Top row shows the PCA filters corresponding to the largest eigenvalues (1–5), second row to the smallest eigenvalues (196–200). The bottom rows show 10 filters after the random rotation in the principal subspace. Right: Top row shows 5 examples from the MNIST handwritten digit data set. The second row shows reconstructions from the 200 component principal subspace. The third and forth row show reconstructions when including the added noise to the training data. Two instantiations of the noise is shown to remind that the noise is resampled for each epoch.

**Learning Rate:** We use a hand-set learning rate $\gamma_i$ for each problem. The base learning rate $\gamma$ is halved for shortcut mappings once for each layer that the mapping skips. This is a heuristic to take into account that shortcut connections have a more direct influence to the final output of the network. Also we linearly decrease the learning rate to zero after half of the allocated computational time has been used.

**Regularization:** Overfitting is a crucial problem in large networks, and regularization is essential to make it perform well with new data. We use three different regularization methods. Firstly, the dimensionality of input data was decreased as a preprocessing step. We used principal component analysis (PCA) followed by a random rotation (see Figure 2). The motivation for the random rotation was to make each input approximately equally important, which works well with the uniform initialization of weights and with the following regularization procedures. Secondly, we use weight decay (see e.g. [6]), or equivalently a Gaussian prior on the parameters $\boldsymbol{\theta}$. The final update direction becomes

$$g_i \leftarrow 0.1 \left[ \frac{\partial \langle \log p(\mathbf{y}_t \mid \mathbf{x}_t, \boldsymbol{\theta}) \rangle}{\partial \theta_i} - \lambda \theta_i \right] + 0.9 g_i, \tag{17}$$

where $\lambda$ is the weight decay parameter. Thirdly, we add randomly generated Gaussian noise to the original input data each time they are presented to the learner, inspired by denoising autoencoders in [13] and also recently proposed in [11].

## 5 Experiments

We compare MLP learning with and without proposed transformations in three problems where we can also compare to other state-of-the-art learning algorithms. The two first experiments are image classification tasks, and the last one is an autoregressive model to find a low-dimensional representation of images. Even though all experiments use image data, we do not use or compare to any image-specific processing such as convolutional networks or elastic distortions. Our approach would work exactly the same even if the order of pixels was randomly permutated. All experiments were run on a desktop computer with Intel Core i7 2.93GHz and 8 gigabytes of memory.
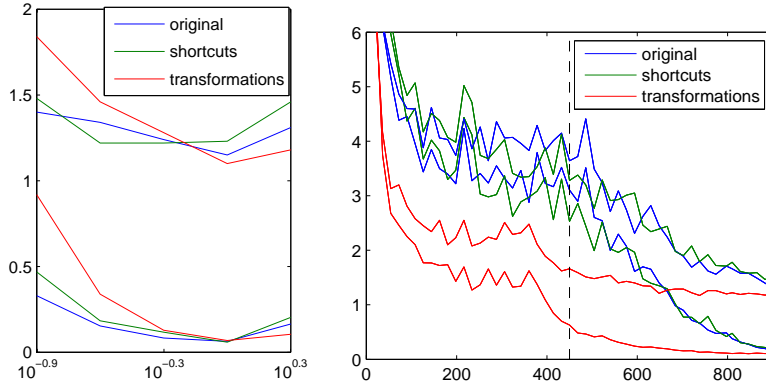
5

Figure 3: MNIST classification problem. Left: Classification error rate in percentage as a function of learning rate after 15 minutes of learning. Lower curves in each figure are the training error rates and higher curves are test error rates. Right: Error rates against learning time for the best learning rates for each method. The vertical dashed line shows the point at which the learning rate starts to be decreased.

## 5.1 MNIST Handwritten Digit Classification

The MNIST data set [8] consists of 28 by 28 pixel gray-scale images of handwritten digits with examples depicted in Figure 2. There are 60000 training samples and 10000 test samples. The mean activation of each pixel was subtracted from the data, and the dimensionality was dropped from $28 \times 28 = 784$ to 200 using PCA followed by a random rotation (see Figure 2). A classification network with layer sizes 200–200–200–10 was learned with a normal MLP model (original), one with shortcut mapping $\mathbf{C}$ included (shortcuts), and the proposed model with shortcut mappings and transformations in the nonlinearities (transformations) (see Figure 1). We also ran a simple 200–10 network (*linear*) for comparison. Training and test errors were tracked during learning and their computation was not included in learning time. Learning time was restricted to 15 minutes, weight decay parameter was $\lambda = 0.0001$ and the regularization noise was used with standard deviation 0.4 for each input (see Figure 2).

The results are shown in Table 1 and Figure 3. With a proper learning rate (1.0, see Figure 3), the proposed method gets the test error below 1.5% in six minutes and to 1.1% in fifteen minutes. Thus, the transformations make a simple gradient descent algorithm competitive in speed with complex state-of-the-art learning methods such as TONGA [7], which reaches test error 1.7% in around 30 minutes (although with an older computer and with a different network structure). Deep networks learned with backpropagation have reached 1.64% error in [4]. Deep belief network [5] gives 1.20% error.

The experiments were run with an 8-core desktop computer and Matlab implementation. The learning time was measured with Matlab function cputime, which counts time spent by each core, so wall-clock learning times were smaller, for example 233 seconds for a 900 cpu-second run including computations of the training and test errors. Table 1 reports the number of iterations reached in the allocated learning time.

Let us also study some of the properties of learning, for instance, how the transformations affect the Fisher information matrix. Equation (7) measures the covariance of the signals $f_i(\cdot)$. The ratio of mean square nondiagonal element of the covariance to mean square diagonal element drops from 0.051 to 0.007 in the first hidden layer and from 0.080 to 0.009 in the second hidden layer, when comparing models learned traditionally or with transformation. Figure 4 shows what the transformed nonlinearities look like in practice. The transformations also decrease the norm of the gradient with respect to non-shortcut weights. The decrease is about 2 to 3-fold in the initial phase of learning. This might also explain why the proposed model performed worse than the others with a too small learning rate (See left part of Figure 3). With a small norm of the gradient and a small learning rate, the optimization simply does not finish in the allocated time.

6

|  |  | linear | original | shortcuts | transformations | literature |
|---|---|---|---|---|---|---|
| MNIST classification | training error | 8.99 | 0.063 | 0.058 | 0.068 | - |
|  | test error | 8.58 | 1.15 | 1.22 | **1.10** | 1.64 |
|  | learning rate | - | 1.0 | 0.5 | 1.0 | - |
|  | # of iterations | 30k | 4717 | 3498 | 2674 | - |
| MNIST autoencoder | training error | 8.11 | 2.37 | 2.11 | 1.94 | 1.75 |
|  | test error | 7.85 | 2.76 | 2.61 | **2.44** | 2.55 |
|  | # of iterations | 92k | 49k | 38k | 37k | - |

Table 1: Results of the proposed method (transformations) compared against other methods run with the same settings and against results in the literature [4, 10]. The number of iterations in the allocated time is reported to compare the computational complexities.
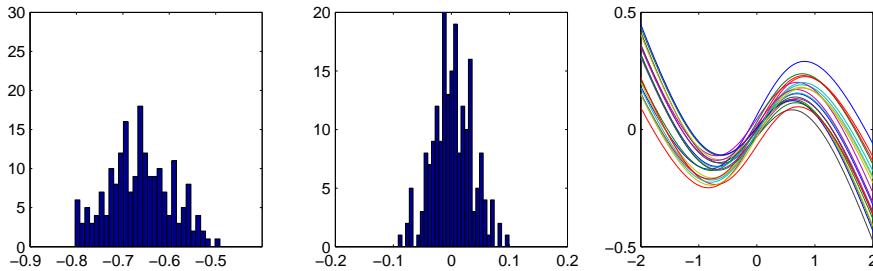


Figure 4: MNIST classification problem. Histograms of the transformation parameters $\alpha_i$ (left) and $\beta_i$ (middle) for the first hidden layer after learning. Right: Twenty examples of corresponding nonlinearities $f_i(\cdot)$.

To study which regularization method was important, the runs were repeated by including modifications one by one. Table 2 shows test errors evaluated by including regularization methods one by one using learning rate $\gamma = 1.0$. The final run was repeated with ten times the learning time and $\gamma = 0.5$. Adding noise to the inputs turned out to be the most important regularization method, followed by dimensionality reduction by PCA. Using the transformations improved all versions.

| regularization | none | weight decay | PCA | rotation | input noise | long run |
|---|---|---|---|---|---|---|
| original | 1.87 | 1.85 | 1.62 | 1.63 | 1.15 | 1.03 |
| shortcuts | 2.02 | 1.77 | 1.59 | 1.60 | 1.23 | 1.17 |
| transform. | 1.63 | 1.56 | 1.56 | 1.48 | 1.10 | **1.02** |

Table 2: Evaluation of regularization methods in MNIST classification.

## 5.2 MNIST Autoencoder

The third problem uses the same MNIST handwritten digit dataset [8], but in this case both the inputs and outputs of the network are the images. The network topology is 784–500–250–30–250–500–784 with tanh-nonlinearity at each layer except the bottleneck layer in the middle. The output is scaled from -1 to 1 to match the tanh nonlinearity. The goal in this problem is to find a good 30-dimensional representation from which the image can be reconstructed. Naturally the shortcut connections that skip the bottleneck are not used. The labels are not used in this experiment at all. Learning time was restricted to 60000 seconds, the base learning rate was $\gamma = 0.05$, weight decay parameter was $\lambda = 0.001$ and the regularization noise was used with standard deviation 0.1 for each input. To avoid early divergence of learning, the learning rate was increased from one hundredth to the full rate exponentially during the first one percent of learning time.

The error is measured by the average sum of squared reconstruction error on the test data when scaled back from [-1,1] to [0,1]. The final reconstruction error is 2.44 and some reconstructions are visualized in the left part of Figure 5. As a comparison, a linear 784–30–784 autoencoder gives an error 7.85. State-of-the-art comparison results are presented by Martens [10]: Hessian-free optimization gives 2.55 with a larger network. The results in [10] were further improved to 2.28 by initializing the network with layerwise pretraining, which could also be used here.
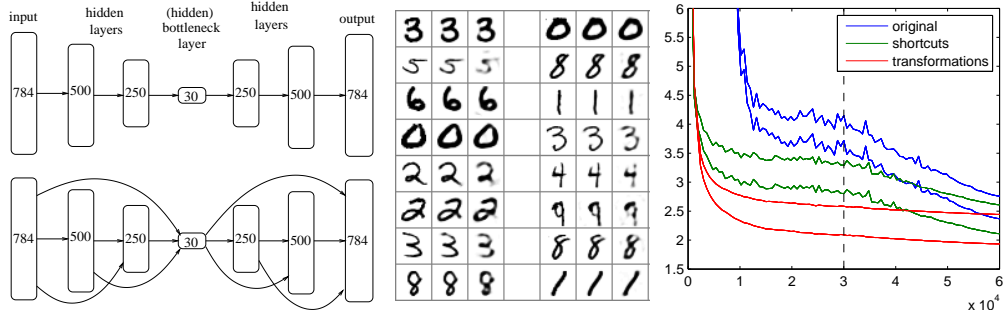
Figure 5: MNIST autoencoder. Left: Traditional autoencoder network above, shortcut connections included below. Note that the bottleneck layer does not have a nonlinearity or transformations. Middle: Each triplet shows an example digit from the validation data, its reconstruction with a deep autoencoder, and a reconstruction with a linear autoencoder as a comparison. Right: Error rate in average sum of squared reconstruction errors plotted against learning time in seconds. Higher curves are test errors, lower curves are training errors. The vertical dashed line shows the point at which the learning rate starts to be decreased.

## 6 Discussion

We proposed transformations that make learning MLP networks much easier. By transforming the nonlinearities of an MLP by Equations (2–4) such that they would become separated from the linear mapping which is modelled separately using shortcut weights. A basic stochastic gradient became faster than state-of-the-art learning algorithms. Those algorithms could also be tested with transformations for further improvements. The theory of the speed-up is based on making the standard gradient and the natural gradient more similar by having the nondiagonal terms of the Fisher information matrix closer to zero. As a side effect, the added linear part in the nonlinearity might also help in avoiding plateaus [9]. The experiments showed that these simple transformations helps learning deep structures at least up to 5 hidden layers using good-old backpropagation learning.

The transformations also seemed to help generalization when no regularization was used. We think that this is because the transformations help to separate the simpler and complex parts of each mapping. Let us think of a network with just one hidden layer. The linear part of the problem (the shortcut connections $\mathbf{C}$) do not suffer much from the overfitting and can be learned more reliably even without regularization. Overfitting the more complex parts $\mathbf{A}$ and $\mathbf{B}$ does not hurt the performance as much when they do not influence the linear part of the whole mapping.

The effect of the transformations could also be studied in other contexts. For variational Bayesian (VB) learning, the same transformations make both the signals in the network and the network weights less dependent a posteriori. Often they are assumed to be independent in the posterior approximation of VB anyway. Thus, the transformation makes the effect of the assumption smaller and the approximation more accurate. This should help in avoiding the problem of underfitting (or output weights of too many hidden neurons going to zero). One could also use MCMC methods for sampling network weights. The effectiveness of the sampling process depends heavily on how large jumps in the parameter space can be made. We can make longer jumps for the matrices $\mathbf{A}$ and $\mathbf{B}$ if we use the proposed transformations to ensure that even large changes in them do not affect the linear part of the input-output mapping.

The term deep learning refers either to networks with many layers (as in this work) but sometimes it is used for a network learned with unsupervised pretraining which allows for better performing deeper networks in practice. The proposed transformations could also be applied to initializations based on unsupervised pretraining. We could also study a network with both the class labels and the inputs are used as outputs.

# References

[1] S. Amari. Natural gradient works efficiently in learning. *Neural Computation*, 10(2):251–276, 1998.

[2] R. Battiti. First- and second-order methods for learning: Between steepest descent and Newton's method. *Neural Computation*, 4(2):141–166, 1992.

[3] D. Ciresan, U. Meier, L. Gambardella, and J. Schmidhuber. Deep big simple neural nets excel on handwritten digit recognition. *CoRR*, abs/1003.0358, 2010.

[4] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics (AISTATS)*, pages 249–256, 2010.

[5] G. E. Hinton and R. R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.

[6] A. Krogh and J. Hertz. A simple weight decay can improve generalization. In *Advances in Neural Information Processing Systems 4 (NIPS 1991)*, pages 950–957, 1992.

[7] N. Le Roux, P. Manzagol, and Y. Bengio. Topmoumoute online natural gradient algorithm. In *Advances in Neural Information Processing Systems 20 (NIPS*2007)*, 2008.

[8] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, volume 86, pages 2278–2324, 1998.

[9] Y. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller. Efficient backprop. In *Neural Networks: tricks of the trade*. Springer-Verlag, 1998.

[10] J. Martens. Deep learning via Hessian-free optimization. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, 2010.

[11] S. Rifai, X. Glorot, Y. Bengio, and P. Vincent. Adding noise to the input of a model trained with a regularized objective. Technical Report 1359, Université de Montréal, Montréal (QC), H3C 3J7, Canada, Apr. 2011.

[12] N. Schraudolph. Slope centering: Making shortcut weights effective. In *Proceedings of the Eight International Conference on Artificial Neural Networks*, 1998.

[13] P. Vincent, H. Larochelle, Y. Bengio, and P. Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the Twenty-fifth International Conference on Machine Learning (ICML08)*, pages 1096–1103, 2008.