# Pushing Stochastic Gradient towards Second-Order Methods – Backpropagation Learning with Transformations in Nonlinearities

**Tommi Vatanen, Tapani Raiko, Harri Valpola**
Department of Information and Computer Science
Aalto University School of Science
P.O.Box 15400, FI-00076, Aalto, Espoo, Finland
`first.last@aalto.fi`

**Yann LeCun**
New York University
715 Broadway, New York, NY 10003, USA
`firstname@cs.nyu.edu`

## Abstract

Recently, we proposed to transform the outputs of each hidden neuron in a multi-layer perceptron network to have zero output and zero slope on average, and use separate shortcut connections to model the linear dependencies instead. We continue the work by firstly introducing a third transformation to normalize the scale of the outputs of each hidden neuron, and secondly by analyzing the connections to second order optimization methods. We show that the transformations make a simple stochastic gradient behave closer to second-order optimization methods and thus speed up learning. This is shown both in theory and with experiments. The experiments on the third transformation show that while it further increases the speed of learning, it can also hurt performance by converging to a worse local optimum, where both the inputs and outputs of many hidden neurons are close to zero.

## 1 Introduction

Learning deep neural networks has become a popular topic since the invention of unsupervised pretraining [4]. Some later works have returned to traditional back-propagation learning in deep models and noticed that it can also provide impressive results [6] given either a sophisticated learning algorithm [9] or simply enough computational power [2]. In this work we study back-propagation learning in deep networks with up to five hidden layers, continuing on our earlier results in [10].

In learning multi-layer perceptron (MLP) networks by back-propagation, there are known transformations that speed up learning [8, 11, 12]. 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 [8]. Schraudolph [12, 11] proposed centering all factors in the gradient to have zero mean, and further adding linear shortcut connections that bypass the nonlinear layer. Gradient factor centering changes the gradient as if the nonlinear activation functions had zero mean and zero slope on average. As such, it does not change the model itself. It is assumed that the discrepancy between the model and the gradient is not an issue, since the errors will be easily compensated by the linear shortcut connections in the proceeding updates. Gradient factor centering leads to a significant speed-up in learning.

In this paper, we transform the nonlinear activation functions in the hidden neurons such that they have on average 1) zero mean, 2) zero slope, and 3) unit variance. Our earlier results in [10] included the first two transformations and here we introduce the third one. cdsaasd We explain the usefulness of these transformations by studying the Fisher information matrix and the Hessian, e.g. by measuring the angle between the traditional gradient and a second order update direction with and without the transformations.

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 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 or the Hessian. Gradient descent can be seen as an approximation of the second-order methods, where the matrix is approximated by a scalar constant times a unit matrix. Our transformations aim at making the Fisher information matrix as close to such matrix as possible, thus diminishing the difference between first and second order methods. Matlab code for replicating the experiments in this paper is available at

https://github.com/tvatanen/ltmlp-neuralnet

## 2 Proposed Transformations

Let us study a MLP-network with a single hidden layer and shortcut mapping, that is, the output column vectors $\mathbf{y}_t$ for each sample $t$ are modeled 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 argument 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 $\alpha_i$, $\beta_i$, and $\gamma_i$ for each nonlinearity $f_i$. They are updated before each gradient evaluation in order to help learning of the other parameters $\mathbf{A}$, $\mathbf{B}$, and $\mathbf{C}$. We define

$$f_i(\mathbf{b}_i\mathbf{x}_t) = \gamma_i\left[\tanh(\mathbf{b}_i\mathbf{x}_t) + \alpha_i\mathbf{b}_i\mathbf{x}_t + \beta_i\right], \tag{2}$$

where $\mathbf{b}_i$ is the $i$th row vector of matrix $\mathbf{B}$. We will ensure that

$$\sum_{t=1}^{T} f_i(\mathbf{b}_i\mathbf{x}_t) = 0 \tag{3}$$

$$\sum_{t=1}^{T} f_i'(\mathbf{b}_i\mathbf{x}_t) = 0 \tag{4}$$

$$\left[\sum_{t=1}^{T} \frac{f_i(\mathbf{b}_i\mathbf{x}_t)^2}{T}\right]\left[\sum_{t=1}^{T} \frac{f_i'(\mathbf{b}_i\mathbf{x}_t)^2}{T}\right] = 1 \tag{5}$$

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

$$\alpha_i = -\frac{1}{T}\sum_{t=1}^{T} \tanh'(\mathbf{b}_i\mathbf{x}_t) \tag{6}$$

$$\beta_i = -\frac{1}{T}\sum_{t=1}^{T} \left[\tanh(\mathbf{b}_i\mathbf{x}_t) + \alpha_i\mathbf{b}_i\mathbf{x}_t\right] \tag{7}$$

$$\gamma_i = \left\{\frac{1}{T}\sum_{t=1}^{T} \left[\tanh(\mathbf{b}_i\mathbf{x}_t) + \alpha_i\mathbf{b}_i\mathbf{x}_t + \beta_i\right]^2\right\}^{1/4}\left\{\frac{1}{T}\sum_{t=1}^{T} \left[\tanh'(\mathbf{b}_i\mathbf{x}_t) + \alpha_i\right]^2\right\}^{1/4}. \tag{8}$$

One way to motivate the first two transformations in Equations (3) and (4), is to study the expected output $\mathbf{y}_t$ and its dependency of the input $\mathbf{x}_t$:

$$\frac{1}{T}\sum_t \mathbf{y}_t = \mathbf{A}\frac{1}{T}\sum_t \mathbf{f}(\mathbf{B}\mathbf{x}_t) + \mathbf{C}\frac{1}{T}\sum_t \mathbf{x}_t \tag{9}$$

$$\frac{1}{T}\sum_t \frac{\partial \mathbf{y}_t}{\partial \mathbf{x}_t} = \mathbf{A}\left[\frac{1}{T}\sum_t \mathbf{f}'(\mathbf{B}\mathbf{x}_t)\right]\mathbf{B}^T + \mathbf{C}. \tag{10}$$

We note that by making nonlinear activations $\mathbf{f}(\cdot)$ zero mean in Eq. (3), we disallow the nonlinear mapping $\mathbf{A}\mathbf{f}(\mathbf{B}\cdot)$ to affect the expected output $\mathbf{y}_t$, that is, to compete with the bias term. Similarly, by making the nonlinear activations $\mathbf{f}(\cdot)$ zero slope in Eq. (4), we disallow the nonlinear mapping $\mathbf{A}\mathbf{f}(\mathbf{B}\cdot)$ to affect the expected dependency of the input, that is, to compete with the linear mapping $\mathbf{C}$. In traditional neural networks, the linear dependencies (expected $\partial \mathbf{y}_t/\partial \mathbf{x}_t$) are modeled by many competing paths from an input to an output (e.g. via each hidden unit), whereas our architecture gathers the linear dependencies to be modeled only by $\mathbf{C}$. We argue that less competition between parts of the model will speed up learning. Another explanation for choosing these transformations is that they make the nondiagonal parts of the Fisher information matrix closer to zero (see Section 3).

The goal of Equation (5) is to normalize both the output signals (similarly as data is often normalized as a preprocessing step – see,e.g., [8]) and the slopes of the output signals of each hidden unit at the same time. This is motivated by observing that the diagonal of the Fisher information matrix contains elements with both the signals and their slopes. By these normalizations, we aim pushing these diagonal elements more similar to each other. As we cannot normalize both the signals and the slopes to unity at the same time, we normalize their geometric mean to unity.

The effect of the first two transformations can be compensated exactly by updating the shortcut mapping $\mathbf{C}$ by

$$\begin{aligned}\mathbf{C}_{\mathrm{new}} = \mathbf{C}_{\mathrm{old}} &- \mathbf{A}(\boldsymbol{\alpha}_{\mathrm{new}} - \boldsymbol{\alpha}_{\mathrm{old}})\mathbf{B} \\ &- \mathbf{A}(\boldsymbol{\beta}_{\mathrm{new}} - \boldsymbol{\beta}_{\mathrm{old}})[0 \ \ 0 \dots 1],\end{aligned} \tag{11}$$

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. The third transformation can be compensated by

$$\mathbf{A}_{\mathrm{new}} = \mathbf{A}_{\mathrm{old}}\boldsymbol{\gamma}_{\mathrm{old}}\boldsymbol{\gamma}_{\mathrm{new}}^{-1}, \tag{12}$$

where $\boldsymbol{\gamma}$ is a diagonal matrix with $\gamma_i$ as the diagonal elements.

Schraudolph [12, 11] proposed centering the factors of the gradient to zero mean. It was argued that deviations from the gradient fall into the linear subspace that the shortcut connections operate in, so they do not harm the overall performance. Transforming the nonlinearities as proposed in this paper has a similar effect on the gradient. Equation (3) corresponds to Schraudolph's *activity centering* and Equation (4) corresponds to *slope centering*.

## 3 Theoretical Comparison to a Second-Order Method

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 due to heavy computations with large matrices. The natural gradient is the basic gradient multiplied from the left by the inverse of the Fisher information matrix. Using basic gradient descent can thus be seen as using the natural gradient while approximating the Fisher information with a unit matrix multiplied by the inverse learning rate. We will show how the first two proposed transformations move the non-diagonal elements of the Fisher information matrix closer to zero, and the third transformation makes the diagonal elements more similar in scale, thus making the basic gradient behave closer to the natural gradient.

The Fisher information matrix contains elements

$$G_{ij} = \sum_t \left\langle \frac{\partial^2 \log p(\mathbf{y}_t \mid \mathbf{x}_t, \mathbf{A}, \mathbf{B}, \mathbf{C})}{\partial \theta_i \partial \theta_j}\right\rangle, \tag{13}$$

3

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 $\mathbf{C}$. Note that here $\mathbf{y}_t$ is a random variable and thus the Fisher information does not depend on the output data. The Hessian matrix is closely related to the Fisher information, but it does depend on the output data and contains more terms, and therefore we show the analysis on the simpler Fisher information matrix.

The elements in the Fisher information matrix are:

$$\frac{\partial}{\partial a_{ij}} \frac{\partial}{\partial a_{i'j'}} \log p = \begin{cases} 0 & i' \neq i \\ -\frac{1}{\sigma_i^2} \sum_t f_j(\mathbf{b}_j \mathbf{x}_t) f_{j'}(\mathbf{b}_{j'} \mathbf{x}_t) & i' = i, \end{cases} \tag{14}$$

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

$$\frac{\partial}{\partial b_{jk}} \frac{\partial}{\partial b_{j'k'}} \log p = -\sum_i \frac{1}{\sigma_i^2} a_{ij} a_{ij'} \sum_t f_j'(\mathbf{b}_j \mathbf{x}_t) f_{j'}'(\mathbf{b}_{j'} \mathbf{x}_t) x_{kt} x_{k't} \tag{15}$$

and

$$\frac{\partial}{\partial c_{ik}} \frac{\partial}{\partial c_{i'k'}} \log p = \begin{cases} 0 & i' \neq i \\ -\frac{1}{\sigma_i^2} \sum_t x_{kt} x_{k't} & i' = i. \end{cases} \tag{16}$$

The cross terms are

$$\frac{\partial}{\partial a_{ij}} \frac{\partial}{\partial b_{j'k}} \log p = -\frac{1}{\sigma_i^2} a_{ij'} \sum_t f_j(\mathbf{b}_j \mathbf{x}_t) f_{j'}'(\mathbf{b}_{j'} \mathbf{x}_t) x_{kt} \tag{17}$$

$$\frac{\partial}{\partial c_{ik}} \frac{\partial}{\partial a_{i'j}} \log p = \begin{cases} 0 & i' \neq i \\ -\frac{1}{\sigma_i^2} \sum_t f_j(\mathbf{b}_j \mathbf{x}_t) x_{kt} & i' = i \end{cases} \tag{18}$$

$$\frac{\partial}{\partial c_{ik}} \frac{\partial}{\partial b_{jk'}} \log p = -\frac{1}{\sigma_i^2} a_{ij} \sum_t f_j'(\mathbf{b}_j \mathbf{x}_t) x_{kt} x_{k't}. \tag{19}$$

Now we can notice that Equations (14–19) contain factors such as $f_j(\cdot)$, $f_j'(\cdot)$, and $x_{it}$. We argue that by making the factors as close to zero as possible, we help in making nondiagonal elements of the Fisher information closer to zero. For instance, $E[f_j(\cdot) f_{j'}(\cdot)] = E[f_j(\cdot)] E[f_{j'}(\cdot)] + \text{Cov}[f_j(\cdot), f_{j'}(\cdot)]$, so assuming that the hidden units $j$ and $j'$ are representing different things, that is, $f_j(\cdot)$ and $f_{j'}(\cdot)$ are uncorrelated, the nondiagonal element of the Fisher information in Equation (14) becomes exactly zero by using the transformations. When the units are not completely uncorrelated, the element in question will be only approximately zero. The same argument applies to all other elements in Equations (15–19), some of them also highlighting the benefit of making the input data $\mathbf{x}_t$ zero-mean. Naturally, it is unrealistic to assume that inputs $\mathbf{x}_t$, nonlinear activations $\mathbf{f}(\cdot)$, and their slopes $\mathbf{f}'(\cdot)$ are all uncorrelated, so the goodness of this approximation is empirically evaluated in the next section.

The diagonal elements of the Fisher can be found in Equations (14–16) when $i = i'$, $j = j'$, and $k = k'$. There we find $f(\cdot)^2$ and $f'(\cdot)^2$ that we aim to keep similar in scale by using the third transformation in Equation (5).

## 4   Empirical Comparison to a Second-Order Method

Here we investigate how linear transformations affect the gradient by comparing it to a second-order method, namely Newton's algorithm with a simple regularization to make the Hessian invertible.

We compute an approximation of the Hessian matrix using finite difference method, in which case $k$-th row vector $\mathbf{h}_k$ of the Hessian matrix $\mathbf{H}$ is given by

$$\mathbf{h}_k = \frac{\partial(\nabla E(\boldsymbol{\theta}))}{\partial \theta_k} \approx \frac{\nabla E(\boldsymbol{\theta} + \delta \boldsymbol{\phi}_k) - \nabla E(\boldsymbol{\theta} - \delta \boldsymbol{\phi}_k)}{2\delta}, \tag{20}$$

where $\boldsymbol{\phi}_k = (0, 0, \dots, 1, \dots, 0)$ is a vector of zeros and 1 at the $k$-th position, and the error function $E(\boldsymbol{\theta}) = -\sum_t \log p(\mathbf{y}_t \mid \mathbf{x}_t, \boldsymbol{\theta})$. The resulting Hessian might still contain some very small or
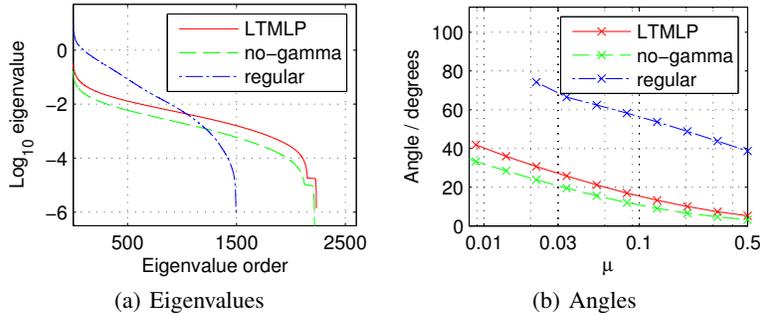
(a) Eigenvalues         (b) Angles

Figure 1: Comparison of (a) distributions of the eigenvalues of Hessians ($2600 \times 2600$ matrix) and (b) angles compared to the second-order update directions using LTMLP and regular MLP. In (a), the eigenvalues are distributed most evenly when using LTMLP. (b) shows that gradients of the transformed networks point to the directions closer to the second-order update.

even negative eigenvalues which cause its inversion to blow up. Therefore we do not use the Hessian directly, but include a regularization term similarly as in the Levenberg-Marquardt algorithm, resulting in a second-order update direction

$$\Delta\boldsymbol{\theta} = (\mathbf{H} + \mu\mathbf{I})^{-1}\nabla E(\boldsymbol{\theta}), \tag{21}$$

where $\mathbf{I}$ denotes the unit matrix. Basically, Equation (21) combines the steepest descent and the second-order update rule in such a way, that when $\mu$ gets small, the update direction approaches the Newton's method and vice versa.

Computing the Hessian is computationally demanding and therefore we have to limit the size of the network used in the experiment. We study the MNIST handwritten digit classification problem where the dimensionality of the input data has been reduced to 30 using PCA with a random rotation [10]. We use a network with two hidden layers with architecture 30–25–20–10. The network was trained using the standard gradient descent with weight decay regularization. Details of the training are given in the appendix.

In what follows, networks with all three transformations (*LTMLP*, linearly transformed multi-layer perceptron network), with two transformations (*no-gamma* where all $\gamma_i$ are fixed to unity) and a network with no transformations (*regular*, where we fix $\alpha_i = 0$, $\beta_i = 0$, and $\gamma_i = 1$) were compared. The Hessian matrix was approximated according to Equation (20) 10 times in regular intervals during the training of networks. All figures are shown using the approximation after 4000 epochs of training, which roughly corresponds to the midpoint of learning. However, the results were parallel to the reported ones all along the training.

We studied the eigenvalues of the Hessian matrix ($2600 \times 2600$) and the angles between the methods compared and second-order update direction. The distribution of eigenvalues in Figure 1a for the networks with transformations are more even compared to the regular MLP. Furthermore, there are fewer negative eigenvalues, which are not shown in the plot, in the transformed networks. In Figure 1b, the angles between the gradient and the second-order update direction are compared as a function of $\mu$ in Equation (21). The plots are cut when $\mathbf{H} + \mu\mathbf{I}$ ceases to be positive definite as $\mu$ decreases. Curiously, the update directions are closer to the second-order method, when $\gamma$ is left out, suggesting that $\gamma$s are not necessarily useful in this respect.

Figure 2 shows histograms of the diagonal elements of the Hessian after 4000 epochs of training. All the distributions are bimodal, but the distributions are closer to unimodal when transformations are used (subfigures (a) and (b))[1]. Furthermore, the variance of the diagonal elements in log-scale is smaller when using LTMLP, $\sigma_\mathrm{a}^2 = 0.90$, compared to the other two, $\sigma_\mathrm{b}^2 = 1.71$ and $\sigma_\mathrm{c}^2 = 1.43$. This suggests that when transformations are used, the second-order update rule in Equation (21) corrects different elements of the gradient vector more evenly compared to a regular back-propagation learning, implying that the gradient vector is closer to the second-order update direction when using all the transformations.

---

[1]It can be also argued whether (a) is more unimodal compared to (b).
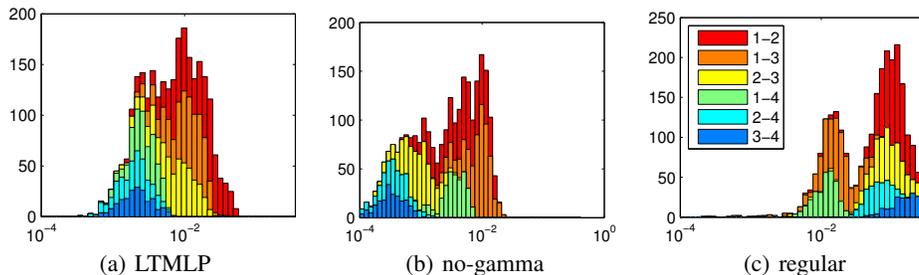
Figure 2: Comparison of distributions of the diagonal elements of Hessians. Coloring according to legend in (c) shows which layers to corresponding weights connect (1 = input, 4 = output). Diagonal elements are most concentrated in LTMLP and most spread in the regular MLP network. Notice the logarithmic x-axis.

To conclude this section, there is no clear evidence in way or another whether the addition of $\gamma$ benefits the back-propagation learning with only $\alpha$ and $\beta$. However, there are some differences between these two approaches. In any case, it seems clear that transforming the nonlinearities benefits the learning compared to the standard back-propagation learning.

## 5 Experiments: MNIST Classification

We use the proposed transformations for training MLP networks for MNIST classification task. Experiments are conducted without pretraining, weight-sharing, enhancements of the training set or any other known tricks to boost the performance. No weight decay is used and as only regularization we add Gaussian noise with $\sigma = 0.3$ to the training data. Networks with two and three hidden layers with architechtures 784–800–800–10 (solid lines) and 784–400–400–400–10 (dashed lines) are used. Details are given in the appendix.

Figure 3 shows the results as number of errors in classifying the test set of 10 000 samples. The results of the regular back-propagation without transformations, shown in blue, are well in line with previously published result for this task. When networks with same architecture are trained using the proposed transformations, the results are improved significantly. However, adding $\gamma$ in addition to previously proposed $\alpha$ and $\beta$ does not seem to affect results on this data set. The best results, 112 errors, is obtained by the smaller architecture without $\gamma$ and for the three-layer architecture with $\gamma$ the result is 114 errors. The learning seems to converge faster, especially in the three-layer case, with $\gamma$. The results are in line what was obtained in [10] where the networks were regularized more thoroughly. These results show that it is possible to obtain results comparable to dropout networks (see [5]) using only minimal regularization.

## 6 Experiments: MNIST Autoencoder

Previously, we have studied an auto-encoder network using two transformations, $\alpha$ and $\beta$, in [10]. Now we use the same auto-encoder architecture, 784–500–250–30–250–500–784. Adding the third transformation $\gamma$ for training the auto-encoder poses problems. Many hidden neurons in decoding layers (i.e., 4th and 5th hidden layers) tend to be relatively inactive in the beginning of training, which induces corresponding $\gamma$s to obtain very large values. In our experiments, auto-encoder with $\gamma$s eventually diverge despite simple constraint we experimented with, such as $\gamma_i \leq 100$. This behavior is illustrated in Figure 4. The subfigure (a) shows the distribution of variances of outputs of all hidden neurons in MNIST classification network used in Section 5 given the MNIST training data. The corresponding distribution for hidden neurons in the decoder part of the auto-encoder is shown in the subfigure (b). The "dead neurons" can be seen as a peak in the origin. The corresponding $\gamma$s, constrained $\gamma_i \leq 100$, can be seen in the subfigure (c). We hypothesize that this behavior is due to the fact, that in the beginning of the learning there is not much information reaching the bottleneck layer through the encoder part and thus there is nothing to learn for the decoding neurons. According to our tentative experiments, the problem described above may be overcome by disabling $\gamma$s in the
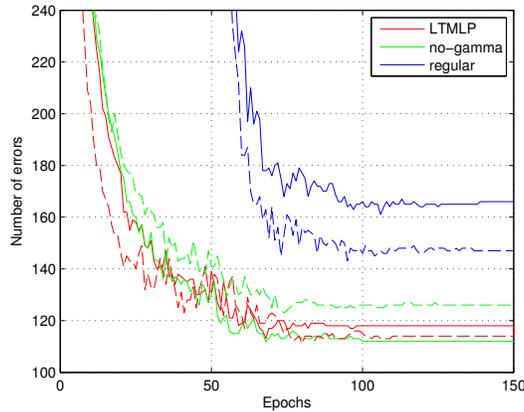
Figure 3: The error rate on the MNIST test set for LTMLP training, LTMLP without $\gamma$ and regular back-propagation. The solid lines show results for networks with two hidden layers of 800 neurons and the dashed lines for networks with three hidden layers of 400 neurons.
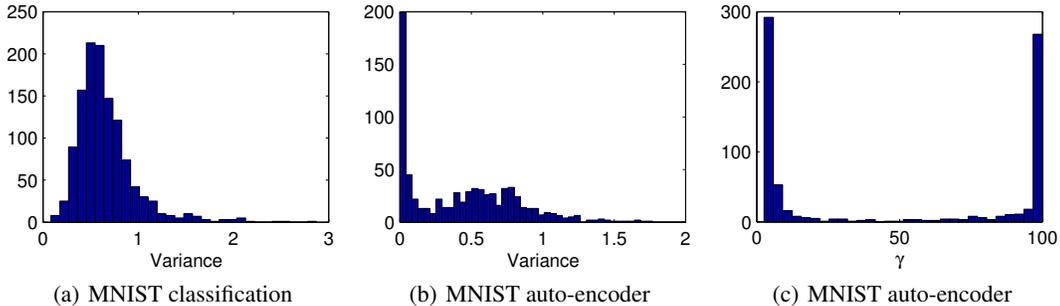


(a) MNIST classification   (b) MNIST auto-encoder   (c) MNIST auto-encoder

Figure 4: Histograms of (a-b) variation of output of hidden neurons given the MNIST training data and (c) $\gamma$s of the decoder part (4th and 5th hidden layer) in the MNIST auto-encoder. (a) shows a healthy distributions of variances, whereas in (b), which includes only variances of the decoder part, there are many "dead neurons". These neurons induce corresponding $\gamma$s, histogram of which is shown in (c), to blow up which eventually lead to divergence.

decoder network (i.e., fix $\gamma = 1$). However, this does not seem to speed up the learning compared to our earlier results with only two transformations in [10]. It is also be possible to experiment with weight-sharing or other constraints to overcome the difficulties with $\gamma$s.

## 7   Discussion and Conclusions

We have shown that introducing linear transformation in nonlinearities significantly improves the back-propagation learning in (deep) MLP networks. In addition to two transformation proposed earlier in [10], we propose adding a third transformation in order to push the Fisher information matrix closer to unit matrix (apart from its scale). The hypothesis proposed in [10], that the transformations actually mimic a second-order update rule, was confirmed by experiments comparing the networks with transformations and regular MLP network to a second-order update method. However, in order to find out whether the third transformation, $\gamma$, we proposed in this paper, is really useful, more experiments ought to be conducted. It might be useful to design experiments where convergence is usually very slow, thus revealing possible differences between the methods. As hyperparameter selection and regularization are usually nuisance in practical use of neural networks, it would be interesting to see whether combining dropouts [5] and our transformations can provide a robust framework enabling training of robust neural networks in reasonable time.

The effect of the first two transformations is very similar to gradient factor centering [12, 11], but transforming the model instead of the gradient makes it easier to generalize to other contexts: When learning by by MCMC, variational Bayes, or by genetic algorithms, one would not compute the basic gradient at all. For instance, consider using the Metropolis algorithm on the weight matrices, and expecially matrices $\mathbf{A}$ and $\mathbf{B}$. Without transformations, the proposed jumps would affect the expected output $\mathbf{y}_t$ and the expected linear dependency $\partial \mathbf{y}_t / \partial \mathbf{x}_t$ in Eqs. (9)–(10), thus often leading to low acceptance probability and poor mixing. With the proposed transformations included, longer proposed jumps in $\mathbf{A}$ and $\mathbf{B}$ could be accepted, thus mixing the nonlinear part of the mapping faster. For further discussion, see [10], Section 6. The implications of the proposed transformations in these other contexts are left as future work.

# References

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

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

[3] 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.

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

[5] Geoffrey E. Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *CoRR*, abs/1207.0580, 2012.

[6] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. 2012.

[7] N. Le Roux, P. A. 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, G. B. Orr, and K.-R. Müller. Efficient backprop. In *Neural Networks: tricks of the trade*. Springer-Verlag, 1998.

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

[10] Tapani Raiko, Harri Valpola, and Yann LeCun. Deep learning made easier by linear transformations in perceptrons. *Journal of Machine Learning Research - Proceedings Track*, 22:924–932, 2012.

[11] N. N. Schraudolph. Accelerated gradient descent by factor-centering decomposition. Technical Report IDSIA-33-98, Istituto Dalle Molle di Studi sull'Intelligenza Artificiale, 1998.

[12] N. N. Schraudolph. Centering neural network gradient factors. In Genevieve Orr and Klaus-Robert Mller, editors, *Neural Networks: Tricks of the Trade*, volume 1524 of *Lecture Notes in Computer Science*, pages 548–548. Springer Berlin / Heidelberg, 1998.

# Appendix

**Details of Section 4**

In experiments of Section 4, networks with all three transformations (LTMLP), only $\alpha$ and $\beta$ (nogamma) and network with no transformations (regular) were compared. Full batch training without momentum was used to make things as simple as possible. The networks were regularized using weight decay and adding Gaussian noise to the training data. Three hyperparameters, weight decay term, input noise variance and learning rate, were validated for all networks separately. The input data was normalized to zero mean and the network was initialized as proposed in [3], that is, the weights were 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.

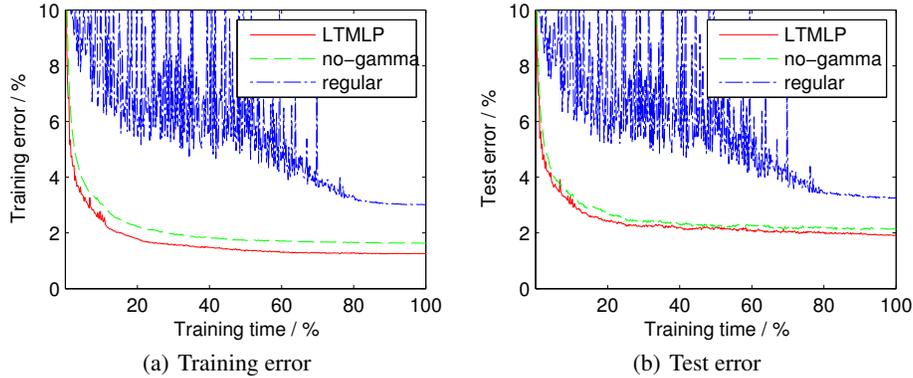(a) Training error                        (b) Test error

Figure 5: Comparison of (a) training and (b) test errors of the algorithms using the MNIST data in the experiment comparing them to the second-order method. Note how the best learning for the regular MLP is relatively high, leading to oscillations until it is annealed towards the end.

We sampled the three hyperparameters randomly (given our best guess intervals) for 500 runs and selected the median of the runs that resulted in the best 50 validation errors as the hyperparameters. Resulting hyperparameters are listed in Table 1. Notable differences occur in step sizes, as it seems that networks with transformations allow using significantly larger step size which in turn results in more complete search in the weight space.

Our weight updates are given by

$$\boldsymbol{\theta}^\tau = \boldsymbol{\theta}^{\tau-1} - \varepsilon^\tau \nabla \boldsymbol{\theta}^\tau. \tag{22}$$

where the learning rate on iteration $\tau$, $\varepsilon^\tau$, is given by

$$\varepsilon^\tau = \left\{ \begin{array}{ll} \varepsilon_0 & \tau \leq T/2 \\ 2(1 - \frac{\tau}{T})\varepsilon_0 & \tau > T/2 \end{array} \right. \tag{23}$$

that is, the learning rate starts decreasing linearly after the midpoint of the given training time $T$. Furthermore, the learning rate $\varepsilon^\tau$ is dampened for shortcut connection weights by multiplying with $\left(\frac{1}{2}\right)^s$, where $s$ is number of skipped layers as proposed in [10].[2] Figure 5 shows training and test errors for the networks. The LTMLP obtains the best results although there is no big difference compared to training without $\gamma$.

**Details of Section 5**

The MNIST dataset consists of $28 \times 28$ images of hand-drawn digits. There are 60 000 training samples and 10 000 test samples. We experimented with two networks with two and three hidden layers and number of hidden neurons by arbitrary choice. Training was done in minibatch mode with 1000 samples in each batch and transformations are updated on every iteration using the current minibatch with using (6)–(8). This seems to speed up learning compared to the approach in [10] where transformations were updated only occasionally with the full training data. Random Gaussian noise with $\sigma = 0.3$ was injected to the training data in the beginning of each epoch.

---

[2]This heuristic is not well supported by analysis of Figure 2 and could be re-examined.

Table 1: Hyperparameters for the neural networks

|  | LTMLP | no-gamma | regular |
|---|---|---|---|
| weight decay | $4.6 \times 10^{-5}$ | $1.3 \times 10^{-5}$ | $3.9 \times 10^{-5}$ |
| noise | 0.31 | 0.36 | 0.29 |
| step size | 1.2 | 2.5 | 0.45 |

Our weight update equations are given by:

$$\Delta\boldsymbol{\theta}^\tau = \nabla\boldsymbol{\theta} + p^\tau \Delta\boldsymbol{\theta}^{t-1}, \tag{24}$$

$$\boldsymbol{\theta}^\tau = \boldsymbol{\theta}^{\tau-1} - \varepsilon^\tau \Delta\boldsymbol{\theta}^\tau, \tag{25}$$

where

$$\varepsilon^\tau = \begin{cases} \varepsilon_0 & \tau \leq T \\ \varepsilon_0 f^{\tau-T} & \tau > T \end{cases} \tag{26}$$

$$p^\tau = \begin{cases} \frac{\tau}{T}p_f + (1 - \frac{\tau}{T})p_0 & \tau \leq T \\ p_f & \tau > T \end{cases} \tag{27}$$

In the equations above, $T$ is a "burn-in time" where momentum $p^\tau$ is increased from starting value $p_0 = 0.5$ to $p_f = 0.9$ and learning rate $\varepsilon = \varepsilon_0$ is kept constant. When $\tau > T$ momentum is kept constant and learning rate starts decreasing exponentially with $f = 0.9$. Hyperparameters were not validated but chosen by arbitrary guess such that learning did not diverge. For the regular training, $\varepsilon_0 = 0.05$ was selected since it diverged with higher learning rates. Then according to lessons learned, e.g. in Section 4, $\varepsilon_0 = 0.3$ was set for LTMLP with $\gamma$ and $\varepsilon_0 = 0.7$ for the variant with no $\gamma$. Basically, it seems that transformations allow using higher learning rates and thus enable faster convergence.