

# Backstitch: Counteracting Finite-sample Bias via Negative Steps

Yiming Wang<sup>1</sup>, Vijayaditya Peddinti<sup>1,2</sup>, Hainan Xu<sup>1</sup>,  
Xiaohui Zhang<sup>1</sup>, Daniel Povey<sup>1,2</sup>, Sanjeev Khudanpur<sup>1,2</sup>

<sup>1</sup>Center for Language and Speech Processing

<sup>2</sup>Human Language Technology Center of Excellence

The Johns Hopkins University, Baltimore, MD 21218, USA

{yiming.wang, vijay.p, hxu31, xiaohui}@jhu.edu, dpovey@gmail.com, khudanpur@jhu.edu

## Abstract

In this paper we describe a modification to Stochastic Gradient Descent (SGD) that improves generalization to unseen data. It consists of doing two steps for each minibatch: a backward step with a small negative learning rate, followed by a forward step with a larger learning rate. The idea was initially inspired by ideas from adversarial training, but we show that it can be viewed as a crude way of canceling out certain systematic biases that come from training on finite data sets. The method gives  $\sim 10\%$  relative improvement over our best acoustic models based on lattice-free MMI, across multiple datasets with 100-300 hours of data.

**Index Terms:** deep learning, stochastic gradient descent, speech recognition.

## 1. Introduction

Recently, the concept of training on adversarial examples has been proposed [1, 2, 3, 4]. We had previously attempted to use adversarial training for speech-recognition tasks, but failed to obtain any improvement. It occurred to us that a “model-space” version of adversarial training might work better, and this became the *backstitch* method, which is as follows. When processing a minibatch, instead of taking a single SGD step, we first take a step with  $-\alpha$  times the current learning rate, for a small  $\alpha$  (e.g. 0.3), and then a step with  $1 + \alpha$  times the learning rate, with the same minibatch (and a recomputed gradient). So we are taking a small negative step, then a larger positive step. This resulted in unexpectedly large improvements – around 10% relative improvements for our best speech recognition models based on lattice-free MMI (LF-MMI) [5], and the improvement was consistent across datasets.

In this paper we will develop the outlines of a theory why backstitch training might be working, based on the notion that it counteracts finite-sample bias. We will show results on a number of LVCSR tasks and find consistent improvements compared with conventional SGD.

In Section 2 we will discuss finite-sample bias and how it affects optimization tasks where we are training on samples from an underlying distribution. Section 3 analyzes backstitch

training from this perspective and discusses the conditions under which it can be considered to be counteracting finite-sample bias. Section 4 discusses the interaction of backstitch training with various other aspects of our training framework. Section 5 shows our experimental results on multiple datasets, and we conclude in Section 6.

## 2. Finite-sample bias

By finite-sample bias what we specifically mean is that any time we train using gradients from samples that we have seen before, for most model types this will cause the gradient estimate to be systematically biased. We first illustrate this via an example and will then show what we mean more generally.

### 2.1. Example

Probably the simplest example of finite-sample bias is the case where we are estimating the mean and variance of a Gaussian distribution to maximize the likelihood of some data. Maximizing the likelihood of the data from finite samples leads to variances which are systematically smaller than the real variance. This could be formulated as maximizing a function

$$f(x; \theta) = N(x; \mu, \sigma^2) \quad (1)$$

where  $\theta = (\mu, \sigma^2)$ . Although there are easier ways to estimate a mean and variance, it’s quite possible to do so via SGD; and using a likelihood-based objective function, this will converge to the biased maximum likelihood estimate.

### 2.2. Finite-sample bias (more general case)

We assume we are minimizing  $E_{\mathbf{x} \sim P}[f(\mathbf{x}; \theta)]$ , and that there is a global maximum  $\theta^*$ . We also assume that the Hessian  $\mathbf{H}$  equals  $\mathbf{I}$ , which will simplify some of the following expressions. This can be achieved via a change of variables, as long as the original Hessian is full rank. Suppose that we have access only to a finite number of samples from  $P$ :  $\mathbf{x}_i$  for  $i = 1, \dots, I$ .

To establish some compact notation for the gradients and Hessians for the training samples  $\mathbf{x}_i$ , let:

$$\mathbf{g}(\mathbf{x}) \triangleq \left. \frac{\partial f(\mathbf{x}; \theta)}{\partial \theta} \right|_{\theta=\theta^*} \quad (2)$$

$$\mathbf{H}(\mathbf{x}) \triangleq \left. \frac{\partial^2 f(\mathbf{x}; \theta)}{\partial^2 \theta} \right|_{\theta=\theta^*} \quad (3)$$

and as a shorthand:

$$\mathbf{g}_i \triangleq \mathbf{g}(\mathbf{x}_i) \quad (4)$$

$$\mathbf{H}_i \triangleq \mathbf{H}(\mathbf{x}_i) \quad (5)$$

---

This work was partially supported by DARPA BOLT Contract No HR0011-12-C-0015, NSF Grant No IIS 0963898, CRI-1513128 and IARPA BABEL Contract No W911NF-12-C-0015, 2012-12050800010. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of DARPA, IARPA, DoD/ARL or the U.S. Government.

and let the averages of these be:

$$\mathbf{g}_{\text{avg}} \triangleq \frac{1}{I} \sum_i \mathbf{g}_i \quad (6)$$

$$\mathbf{H}_{\text{avg}} \triangleq \frac{1}{I} \sum_i \mathbf{H}_i \quad (7)$$

The argmin of the objective function given the sampled data can be approximated as:

$$\boldsymbol{\theta}^{\text{est}} \simeq \boldsymbol{\theta}^* - \mathbf{H}_{\text{avg}}^{-1} \mathbf{g}_{\text{avg}} \quad (8)$$

and this is an exact expression if  $f(\mathbf{x}; \boldsymbol{\theta})$  is quadratic in  $\boldsymbol{\theta}$ . Since we previously specified that the ‘‘true’’ expected Hessian  $\mathbf{H}$  is identity, we can write

$$\mathbf{H}_{\text{avg}}^{-1} \simeq 2\mathbf{I} - \mathbf{H}_{\text{avg}} \quad (9)$$

which comes from a first-order Taylor expansion around  $\mathbf{H}_{\text{avg}} = \mathbf{I}$ . This gives us:

$$\begin{aligned} \boldsymbol{\theta}^{\text{est}} &\simeq \boldsymbol{\theta}^* - 2\mathbf{g}_{\text{avg}} + \mathbf{H}_{\text{avg}} \mathbf{g}_{\text{avg}} \\ &= \boldsymbol{\theta}^* - 2\mathbf{g}_{\text{avg}} + \frac{1}{I^2} \sum_{i,j} \mathbf{H}_i \mathbf{g}_j \end{aligned} \quad (10)$$

We are interested in the expected bias in the estimate of  $\boldsymbol{\theta}$ , i.e. in  $E[\boldsymbol{\theta}^{\text{est}} - \boldsymbol{\theta}^*]$ , where the expectations are taken over the set of  $I$  training samples generated from their underlying distribution. In (10), the term  $-2\mathbf{g}_{\text{avg}}$  does not lead to any expected bias, because the expected value of  $\mathbf{g}$  is zero. Nor do the terms in summation on the right for  $i \neq j$ , because (due to independence) the expectation can be decomposed into the product of two terms, one of which (involving  $\mathbf{g}$ ) is zero. The only potential bias comes from the ‘‘self-terms’’ (for  $i = j$ ), i.e. from the quantity  $\frac{1}{I^2} \sum_i \mathbf{H}_i \mathbf{g}_i$ , so we can write:

$$E_{\mathbf{x}}[\boldsymbol{\theta}^{\text{est}} - \boldsymbol{\theta}^*] = \frac{1}{I} E_{\mathbf{x}}[\mathbf{H}(\mathbf{x})\mathbf{g}(\mathbf{x})] \quad (11)$$

Since we previously ensured that the expected Hessian is identity ( $\mathbf{H} = \mathbf{I}$ ), adding an extra term  $\mathbf{c}$  to the derivative of  $f(\mathbf{x}; \boldsymbol{\theta})$  during training would approximately become an offset  $-\mathbf{c}$  in the estimated parameter, so in order to cancel out the bias of Equation (11) it’s reasonable to add a correction term:

$$\mathbf{c} = \frac{1}{I} E_{\mathbf{x}}[\mathbf{H}(\mathbf{x})\mathbf{g}(\mathbf{x})] \quad (12)$$

to the gradients we use to train the network. The quantities  $\mathbf{H}(\mathbf{x})$  and  $\mathbf{g}(\mathbf{x})$  are of course defined w.r.t. the optimal parameter  $\boldsymbol{\theta}^*$  which is unknown, and in practice any expressions will use the current  $\boldsymbol{\theta}$ .

### 2.3. Finite-sample bias (non-identity Hessian)

Above, we assumed that via a change of variables, the Hessian was identity. In this section we find the equivalent expressions for when the Hessian is positive-definite but non-identity. Suppose we have a non-identity  $\mathbf{H}_{\text{orig}}$  w.r.t. an original parameter  $\boldsymbol{\psi}$ , then we can define

$$\boldsymbol{\theta} = \mathbf{H}_{\text{orig}}^{0.5} \boldsymbol{\psi}. \quad (13)$$

We’ll use the notation  $\mathbf{g}_{\boldsymbol{\theta}}$  and  $\mathbf{H}_{\boldsymbol{\theta}}$  for a derivative or Hessian w.r.t.  $\boldsymbol{\theta}$ , and  $\mathbf{g}_{\boldsymbol{\psi}}$  and  $\mathbf{H}_{\boldsymbol{\psi}}$  for a derivative or Hessian w.r.t.  $\boldsymbol{\psi}$ . The conversion is:

$$\mathbf{g}_{\boldsymbol{\theta}} = \mathbf{H}_{\text{orig}}^{-0.5} \mathbf{g}_{\boldsymbol{\psi}} \quad (14)$$

$$\mathbf{H}_{\boldsymbol{\theta}} = \mathbf{H}_{\text{orig}}^{-0.5} \mathbf{H}_{\boldsymbol{\psi}} \mathbf{H}_{\text{orig}}^{-0.5} \quad (15)$$

By adding  $\boldsymbol{\theta}$  subscripts to Equation (11), and expanding out the expressions for  $\mathbf{H}_{\boldsymbol{\theta}}$  and  $\mathbf{g}_{\boldsymbol{\theta}}$  above, we would get:

$$E[\boldsymbol{\theta}^{\text{est}} - \boldsymbol{\theta}^*] = \frac{1}{I} E_{\mathbf{x}}[\mathbf{H}_{\text{orig}}^{-0.5} \mathbf{H}_{\boldsymbol{\psi}}(\mathbf{x}) \mathbf{H}_{\text{orig}}^{-0.5} \mathbf{H}_{\text{orig}}^{-0.5} \mathbf{g}_{\boldsymbol{\psi}}(\mathbf{x})], \quad (16)$$

and because we want the expectation in terms of the difference in  $\boldsymbol{\psi}$ , we multiply Equation (11) by  $\mathbf{H}_{\text{orig}}^{-0.5}$  on the left, and then replace  $\mathbf{H}_{\text{orig}}^{-0.5}(\boldsymbol{\theta}^{\text{est}} - \boldsymbol{\theta}^*)$  with  $\boldsymbol{\psi}^{\text{est}} - \boldsymbol{\psi}^*$ , giving us a bias:

$$E[\boldsymbol{\psi}^{\text{est}} - \boldsymbol{\psi}^*] = \frac{1}{I} E_{\mathbf{x}}[\mathbf{H}_{\text{orig}}^{-1} \mathbf{H}_{\boldsymbol{\psi}}(\mathbf{x}) \mathbf{H}_{\text{orig}}^{-1} \mathbf{g}_{\boldsymbol{\psi}}(\mathbf{x})], \quad (17)$$

Suppose we add an offset/correction term  $\mathbf{c}_{\boldsymbol{\psi}}$  to the derivative of  $f(\mathbf{x}; \boldsymbol{\psi})$  during training, this would lead to an approximate offset  $-\mathbf{H}_{\text{orig}}^{-1} \mathbf{c}_{\boldsymbol{\psi}}$  to the parameter  $\boldsymbol{\psi}$ , so the correction term that we need to add to the derivatives of  $f$  w.r.t.  $\boldsymbol{\psi}$  in order to cancel out the bias of Equation (17), is:

$$\mathbf{c}_{\boldsymbol{\psi}} = \frac{1}{I} E_{\mathbf{x}}[\mathbf{H}_{\boldsymbol{\psi}}(\mathbf{x}) \mathbf{H}_{\text{orig}}^{-1} \mathbf{g}_{\boldsymbol{\psi}}(\mathbf{x})]. \quad (18)$$

Comparing with the expression (12) for when the average Hessian is unit, the difference is the factor  $\mathbf{H}_{\text{orig}}^{-1}$ , which is the inverse of the average Hessian.

## 3. Backstitch training and relation to finite-sample bias

We will now introduce some notation for backstitch training and show how, under certain conditions, it acts to counter the finite-sample bias of Equation (17).

Suppose a single iteration of conventional SGD is written

$$\boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \nu \mathbf{g}(\mathbf{x}_n, \boldsymbol{\theta}_n) \quad (19)$$

where  $\mathbf{x}_n$  is whichever sample we choose for the  $n$ ’th SGD iteration and  $\mathbf{g}(\mathbf{x}, \boldsymbol{\theta}_n)$  is the derivative of  $f(\mathbf{x}, \boldsymbol{\theta})$  w.r.t.  $\boldsymbol{\theta}$ , evaluated at  $\boldsymbol{\theta} = \boldsymbol{\theta}_n$ . Although in practice we use minibatches, we will not develop special notation for that, since it doesn’t affect the math– in principle we could define the task as being over minibatches. In backstitch training, we do two steps:

$$\boldsymbol{\theta}'_{n+1} \leftarrow \boldsymbol{\theta}_n + \alpha \nu \mathbf{g}(\mathbf{x}_n, \boldsymbol{\theta}_n) \quad (20)$$

$$\boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}'_{n+1} - (1 + \alpha) \nu \mathbf{g}(\mathbf{x}_n, \boldsymbol{\theta}'_{n+1}) \quad (21)$$

where the constant  $\alpha > 0$  determines how strongly we are applying the backstitch training. View this as a small backwards step followed by a larger forwards step. For purposes of analysis we can telescope these two iterations into one by making a quadratic approximation around  $\boldsymbol{\theta} = \boldsymbol{\theta}_n$ :

$$\boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \nu \mathbf{g}(\mathbf{x}_n, \boldsymbol{\theta}_n) - \nu^2 (\alpha + \alpha^2) \mathbf{H}(\mathbf{x}_n, \boldsymbol{\theta}_n) \mathbf{g}(\mathbf{x}_n, \boldsymbol{\theta}_n) \quad (22)$$

If we were to view this as a correction term to  $\mathbf{g}$ , it would look like:

$$\boldsymbol{\theta}_{n+1} \leftarrow \boldsymbol{\theta}_n - \nu (\mathbf{g}(\mathbf{x}_n, \boldsymbol{\theta}_n) + \mathbf{c}_n) \quad (23)$$

where

$$\mathbf{c}_n = \nu (\alpha + \alpha^2) \mathbf{H}(\mathbf{x}_n, \boldsymbol{\theta}_n) \mathbf{g}(\mathbf{x}_n, \boldsymbol{\theta}_n). \quad (24)$$

Comparing this with Equation (12), this would make sense as an exact bias-correction method if  $\nu(\alpha + \alpha^2) = 1/I$ , where  $I$  is the number of samples, and if we were in a space where the expected Hessian were identity. Of course these are quite strong assumptions, and we haven’t developed a theory of how this method is expected to perform when the assumptions aren’t quite met.

## 4. Other aspects of backstitch training

### 4.1. Efficiency and backstitch interval

Naively implemented, backstitch training would take twice as long per epoch as regular training because we need to train twice on each minibatch. This is inconvenient. To speed it up, we have experimented with versions where we do the backstitch training every  $n$  minibatches, and do normal SGD updates in between. We call this  $n$  value the “backstitch interval”;  $n = 1$  means doing it every update,  $n = 4$  means doing it every 4th update. We have found that the performance improvement we obtained with  $\alpha = 0.3$  and  $n = 1$  could be more efficiently obtained with  $\alpha = 1.0$  and  $n = 4$ .

### 4.2. Interaction with natural gradient

The assumption that the expected Hessian is identity is actually not too unreasonable because we are using a Natural Gradient (NG) method for training [6, 7, 8]. Natural Gradient can be equivalently viewed as a change of variables into a space where a certain factored estimate of the Fisher matrix is identity; and if the objective function can be interpreted as a log probability or likelihood of some kind, under certain regularity conditions the Fisher and Hessian should have the same value [9]. So from a theoretical perspective there is reason to expect that backstitch training should work particularly well with NG. Our preliminary experiments show backstitch training is also effective without NG, but its improvement is less than with NG, as expected.

We should mention regarding the interaction with NG more generally: we have used terminology such as “a parameter space where the Hessian is identity”, as if we were actually performing a change of variables, and we have mentioned how NG takes us closer to such a situation. NG training can indeed be *formulated* as a change of variables, but actually we formulate it as a modification to the derivatives that is more like a matrix multiplication (by the inverse of a factored Fisher matrix). This equivalence with a change of variables holds for backstitch training too, so the interaction is straightforward; we just mention it as it could be a source of confusion.

### 4.3. Interaction with natural gradient updates

We are using the “online NG” method described in [8]. This involves updating an estimate of the Fisher matrix on each minibatch. It is important for the convergence of the overall SGD, that the estimate of the Fisher matrix should be obtained from *previous* minibatches. Otherwise it could cause a bias due to effects similar to the finite-sample bias we are discussing in this paper. For backstitch training we modified our NG implementation to ensure that the Fisher-matrix estimates used in the NG code are not contaminated with the current minibatch. With reference to the two-step training procedure of Equations (20) and (21), we freeze the Fisher-matrix estimates during the first step and allow them to be updated only during the second. We experimented with doing it the “wrong way” (updating the Fisher matrix on the first step, allowing the second step to use contaminated Fisher-matrix estimates), and as expected this degraded the results.

### 4.4. Backstitch training and parameter maximum changes

One aspect of our training procedure is that to prevent instability, we enforce a maximum parameter-change. This is done at two levels: per component (which very roughly means: for each

layer of the neural network), and globally. There is no maximum change per individual parameter. Our normal defaults are a max-change of 0.75 (in Euclidean distance in parameter space) per component per minibatch, and a max-change of 2.0 globally per minibatch, enforced first per component and then globally.

When implementing backstitch training we tried to ensure that the same “effective” learning-rate matrix (after imposing the max-change) was used in both the first and second steps. The way we did this was by scaling the max-change constraints before applying them, by  $\alpha$  in the first pass and  $1 + \alpha$  in the second pass.

### 4.5. Backstitch training and momentum

The interaction of backstitch training with momentum is non-trivial and we have not implemented the combination. Most of our systems do not use momentum anyway, as NG is effective in preventing instability and we usually find that momentum hurts performance. None of our LF-MMI systems [5] use momentum; the only systems where we use momentum as a matter of course are our cross-entropy (CE) trained LSTM-based systems where we use the moderate momentum value of 0.5. We are currently working on a momentum-friendly version, so that it can be used in toolkits that rely on momentum for good performance.

### 4.6. Slow start backstitch training

We noticed that when backstitch is introduced partway through training, there is a sharp degradation in both training-set and validation-set objective function, which is then quite rapidly reversed. If we start training with backstitch enabled, this degradation is visible at the start of training. (These objective function values are measured separately from the normal training process, and are independent from the expected objective-function degradation when we process each minibatch the second time). We don’t yet fully understand why this happens, but in order to prevent any possible bad effects, we always introduce backstitch gradually over about 10 iterations<sup>1</sup>.

## 5. Experiments

We conducted experiments with the Kaldi speech recognition toolkit [10], with techniques including speed perturbation [11], i-vector adaptation [12] and pronunciation and silence probability modeling [13]. We did extensive experiments on three different LVCSR tasks using backstitch SGD training with different setups, including different criteria (CE and LF-MMI [5]), different network architectures (Time-Delay Neural Network [14, 15] (TDNN) with ReLU nonlinearities [16, 17], Bidirectional LSTM [18, 19] (BLSTM) consisting of stacked LSTM layers [20], and a mixture of TDNN and unidirectional LSTM (TDNN-LSTM<sup>2</sup>) which we recently found not only outperforms BLSTM, but is also computationally more efficient

<sup>1</sup>An iteration is how long it takes for each job to process a fixed amount of data in our framework, tuned to take about 2 to 5 minutes’ worth of GPU time. For a typical acoustic model this is  $\sim 1600$  parameter updates.

<sup>2</sup>The network architecture is basically several unidirectional LSTM layers interleaved with TDNN layers, and having one or more densely spliced TDNN layers preceding the first LSTM layer is crucial to achieve good performance. An example recipe can be found at [https://github.com/kaldi-asr/kaldi/blob/master/egs/swbd/s5c/local/chain/tuning/run\\_tdnn\\_lstm\\_1e.sh](https://github.com/kaldi-asr/kaldi/blob/master/egs/swbd/s5c/local/chain/tuning/run_tdnn_lstm_1e.sh).

than BLSTM). The backstitch scale  $\alpha$ , backstitch interval  $n$  (see Sec. 4.1) and number of epochs were tuned for each individual backstitch training experiment, therefore these hyper-parameters vary across different setups. The baseline systems' number of epochs had been tuned previously. A common configuration is,  $\alpha = 1.0, n = 4$ , meaning we use  $\alpha = 1.0$  once every 4 minibatches. Some older experiments, before we introduced the faster version, use  $\alpha = 0.3, n = 1$ . The number of epochs with backstitch training is 1.5  $\sim$  2 times the number of epochs with normal training, because we found that WER continues to improve for more epochs with this method. All the other hyper-parameters are kept unchanged from those in the normal training experiment<sup>3</sup>. The results of LF-MMI and CE systems are shown in Table 1 and Table 2 respectively ( $\alpha = 0.0$  means the baseline SGD training). While improvements are consistently observed across all the setups, they are most prominent with LF-MMI+TDNN-LSTM (which happen to be our best systems) at around 10% relative WER improvement.

Table 1: Comparison of Backstitch SGD training with normal SGD training in LF-MMI systems.

		LF-MMI			
Dataset	$\alpha[/math>/n]$	WER (%)			
		TDNN-LSTM		BLSTM	
AMI-SDM (use IHM alignment)	0.0	dev	eval	dev	eval
	1.0/4	37.9	41.1	39.7	42.9
	Rel. Gain (%)	<b>10.0</b>	<b>8.5</b>	<b>7.6</b>	<b>5.8</b>
Switchboard	0.0	fsh_fg	tg	fsh_fg	tg
	0.3/1 0.2/1	14.1	15.5	14.3	15.8
	Rel. Gain (%)	<b>12.8</b>	<b>10.3</b>	<b>7.0</b>	<b>5.8</b>
TED-LIUM [21]	0.0	dev	test	dev	test
	1.0/4	9.4	8.8	9.4	9.0
	Rel. Gain (%)	<b>11.7</b>	<b>11.4</b>	<b>7.4</b>	<b>10.0</b>

Table 2: Comparison of Backstitch SGD training with normal SGD training in CE systems.

		CE			
Dataset	$\alpha[/math>/n]$	WER (%)			
		TDNN-LSTM		BLSTM	
AMI-SDM (use IHM alignment)	0.0	dev	eval	dev	eval
	0.3/1	37.0	41.0	38.0	41.3
	Rel. Gain (%)	<b>2.7</b>	<b>2.9</b>	<b>2.9</b>	<b>2.4</b>
TED-LIUM [21]	0.0	dev	test	dev	test
	1.0/4	N/A	N/A	9.9	9.1
	Rel. Gain (%)	N/A	N/A	<b>7.5</b>	<b>4.2</b>

Figure 1 compares the training log-probabilities between regular SGD and the backstitch method, on our LF-MMI+TDNN-LSTM system on the AMI-SDM dataset. The main observation, which we consistently see, is that the difference between train and validation objective functions is smaller when using backstitch.

<sup>3</sup>Except that momentum is disabled when doing backstitch training in CE systems. See Section 4.5 for more explanations.

<sup>4</sup>To the best of our knowledge, this is the best number ever reported on AMI-SDM.

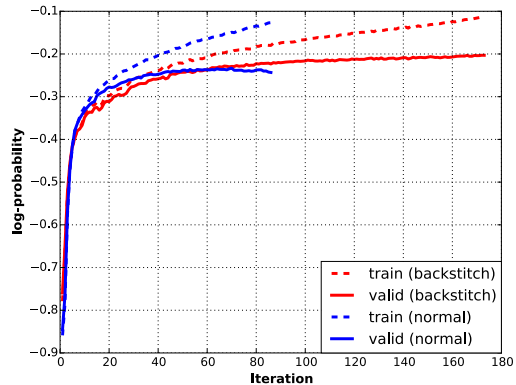


Figure 1: Plot of log-probability vs iterations on AMI-SDM using LF-MMI+TDNN-LSTM (We continued backstitch training for twice the epochs of normal training).

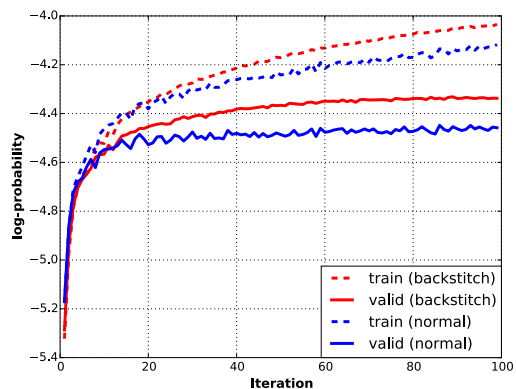


Figure 2: Plot of log-probability vs iterations on AMI-SDM RNNLMs.

We also tested the backstitch method in language modeling tasks with training recurrent neural network language models [22] (RNNLMs) on the AMI-SDM text data. We used nnet3 implementation of neural networks of Kaldi to train RNNLMs with CE objective functions and the backstitch scale is tuned to be 0.8. Figure 2 shows the comparisons between regular SGD and the backstitch method, where we see similar trends. We are encouraged that the backstitch method's usefulness might not be limited to acoustic modeling tasks. However, we have less confidence in our RNNLM results than with our ASR results, since the setup is much newer and may not be as well tuned.

## 6. Conclusion and Future work

In this paper we proposed a modified Stochastic Gradient Descent method consisting of two steps of update for each minibatch. We showed that the proposed method can be considered as a way to approximately eliminate the systemic bias that comes from training on finite data. We observed around 10% relative improvement in WER on multiple LVCSR datasets, versus best systems. The improvement on language modeling also suggests its potential effectiveness in other tasks.

As future work, we would like to evaluate the proposed method in other tasks such as object recognition/classification and machine translation, etc. We are also working on the combination of backstitch with momentum.

## 7. References

- [1] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *ICLR*, 2015.
- [2] R. Huang, B. Xu, D. Schuurmans, and C. Szepesvári, "Learning with a strong adversary," *ICLR*, 2016.
- [3] A. Nokland, "Improving back-propagation by adding an adversarial gradient," *arXiv preprint arXiv:1510.04189*, 2015.
- [4] P. Tabacof and E. Valle, "Exploring the space of adversarial images," in *Neural Networks (IJCNN), 2016 International Joint Conference on*. IEEE, 2016, pp. 426–433.
- [5] D. Povey, V. Peddinti, D. Galvez, P. Ghahmani, V. Manohar, X. Na, Y. Wang, and S. Khudanpur, "Purely sequence-trained neural networks for asr based on lattice-free mmi," *Proc. Interspeech*, 2016.
- [6] H. H. Yang and S.-i. Amari, "Complexity issues in natural gradient descent method for training multilayer perceptrons," *Neural Computation*, vol. 10, no. 8, pp. 2137–2157, 1998.
- [7] N. L. Roux, P.-A. Manzagol, and Y. Bengio, "Topmoumoute online natural gradient algorithm," in *Advances in neural information processing systems*, 2008, pp. 849–856.
- [8] D. Povey, X. Zhang, and S. Khudanpur, "Parallel training of deep neural networks with natural gradient and parameter averaging," *ICLR: Workshop track*, 2015.
- [9] Y. Pawitan, *In all likelihood: statistical modelling and inference using likelihood*. Oxford University Press, 2001.
- [10] D. Povey, A. Ghoshal, G. Boulianne, L. Burget, O. Glembek, N. Goel, M. Hannemann, P. Motlíček, Y. Qian, P. Schwarz *et al.*, "The kaldı speech recognition toolkit," *Proc. ASRU*, 2011.
- [11] T. Ko, V. Peddinti, D. Povey, and S. Khudanpur, "Audio augmentation for speech recognition," *Proc. Interspeech*, pp. 3586–3589, 2015.
- [12] G. Saon, H. Soltau, D. Nahamoo, and M. Picheny, "Speaker adaptation of neural network acoustic models using i-vectors," in *ASRU*, 2013, pp. 55–59.
- [13] G. Chen, H. Xu, M. Wu, D. Povey, and S. Khudanpur, "Pronunciation and silence probability modeling for asr," *Proc. Interspeech*, pp. 533–537, 2015.
- [14] A. Waibel, T. Hanazawa, G. Hinton, K. Shikano, and K. J. Lang, "Phoneme recognition using time-delay neural networks," *IEEE transactions on acoustics, speech, and signal processing*, vol. 37, no. 3, pp. 328–339, 1989.
- [15] V. Peddinti, D. Povey, and S. Khudanpur, "A time delay neural network architecture for efficient modeling of long temporal contexts," *Proc. Interspeech*, 2015.
- [16] M. D. Zeiler, M. Ranzato, R. Monga, M. Mao, K. Yang, Q. V. Le, P. Nguyen, A. Senior, V. Vanhoucke, J. Dean *et al.*, "On rectified linear units for speech processing," in *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*. IEEE, 2013, pp. 3517–3521.
- [17] A. L. Maas, A. Y. Hannun, and A. Y. Ng, "Rectifier nonlinearities improve neural network acoustic models," in *Proc. ICML*, vol. 30, no. 1, 2013.
- [18] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [19] A. Graves and J. Schmidhuber, "Framewise phoneme classification with bidirectional lstm and other neural network architectures," *Neural Networks*, vol. 18, no. 5, pp. 602–610, 2005.
- [20] H. Sak, A. W. Senior, and F. Beaufays, "Long short-term memory recurrent neural network architectures for large scale acoustic modeling," *Proc. Interspeech*, pp. 338–342, 2014.
- [21] A. Rousseau, P. Deléglise, and Y. Estève, "Ted-lium: an automatic speech recognition dedicated corpus," in *LREC*, 2012, pp. 125–129.
- [22] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur, "Recurrent neural network based language model," in *Interspeech*, vol. 2, 2010, p. 3.