

H-DETACH: MODIFYING THE LSTM GRADIENT TOWARDS BETTER OPTIMIZATION

ICLR REPRODUCIBILITY CHALLENGE 2019

Anonymous authors

Paper under double-blind review

ABSTRACT

Recurrent neural networks have been widely used for processing sequences. Recurrent neural networks are known for their exploding and vanishing gradient problem (EVGP). This problem becomes more evident in tasks where the information needed to correctly solve them exist over long time scales, because EVGP prevents important gradient components from being back-propagated adequately over a large number of steps. This paper aims to reproduce the results of the paper Kanuparthi et al. (2019). Kanuparthi et al. (2019) introduces a stochastic algorithm (h-detach) to mitigate the EVGP problem in Long Short Term Memory networks. Long Short Term Memory networks usually just called LSTMs are a special kind of RNN, capable of learning long-term dependencies.

1 INTRODUCTION

Recurrent neural networks Rumelhart et al. (1988) are widely used for processing sequences of information. Recurrent neural networks are very hard to optimize due to the exploding and vanishing gradient problem Bengio et al. (1994). This problem becomes more evident in tasks where training data has dependencies that exist over long time scales.

A number of architectures have been proposed to mitigate this problem. These architectures include Long Short Term Memory Hochreiter & Schmidhuber (1997), Gated Recurrent Unit Chung et al. (2014). These networks introduce a linear temporal path that allow gradients to flow freely across time-steps.

The LSTM is the most widely used architecture. It has been shown that LSTMs have greater representational power as compared to the other architectures Weiss et al. (2018). It is still hard to optimize LSTMs to learn long term dependencies.

The authors say that when the weights of the LSTM are large, the gradients through the linear temporal path get suppressed. This path is important as it was introduced to carry information about long term dependencies. The authors prove this empirically. To prevent this, the authors provide a simple stochastic algorithm (h-detach). The authors experiment with this algorithm on tasks that require capturing long-term dependencies. Examples of such tasks include the copying task Bengio et al. (1994), Pascanu et al. (2012) and the sequential-mnist task Le et al. (2015).

The authors claim that their approach results in faster convergence and stable training. In this work, we attempt to reproduce the results for the copying task and the sequential mnist task (permuted and non-permuted). We also conduct experiments to prove their claims mentioned in the ablation study.

2 PROPOSED METHOD:H-DETACH

We will first introduce the LSTM equations and then analyze the problems with back-propagation according to the authors.

2.1 LONG SHORT TERM MEMORY NETWORKS

LSTM was introduced to provide free-flowing paths for gradients to propagate. The equations of the LSTM can be defined as follows-

$$c_t = f_t \odot c_{t-1} + i_t \odot g_t \quad (1)$$

$$h_t = o_t \odot \tanh(c_t) \quad (2)$$

Here, \odot denotes the pointwise product and the gates f_t , i_t , g_t and o_t are defined as

$$f_t = \sigma(W_{fh}h_{t-1} + W_{fx}x_t + b_f) \quad (3)$$

$$i_t = \sigma(W_{ih}h_{t-1} + W_{ix}x_t + b_i) \quad (4)$$

$$g_t = \tanh(W_{gh}h_{t-1} + W_{gx}x_t + b_g) \quad (5)$$

$$o_t = \sigma(W_{oh}h_{t-1} + W_{ox}x_t + b_o) \quad (6)$$

Here, c_t and h_t are the cell state and hidden state respectively. The output of the LSTM at each time step is $\phi(h_t)$. The loss is calculated as $l_t = l(\phi(h_t))$. In an auto-regressive LSTM networks, the output of one time step is fed as input to the next time step.

There is a linear recursive relation between the cell states c_t . This path is the linear temporal path which allows the free flowing of gradients backwards. More information about LSTMs can be obtained from colah's blog.

2.2 BACKPROPAGATION

Here we will present the authors analysis of the backpropagation equations in brief.

Let w be a weight in one of the weight matrices ($W_{fh}, W_{fx}, W_{ih}, W_{ix}, W_{gh}, W_{gx}, W_{oh}$ or W_{ox}). Let,

$$z_t = \begin{vmatrix} dh_t/dw \\ dc_t/dw \\ 1_n \end{vmatrix} \quad (7)$$

Then, $z_t = (A_t + B_t)z_{t-1}$. The exact definitions of A_t and B_t can be obtained from Kanuparthi et al. (2019). The equation for z_t can be expanded as

$$z_t = (A_t + B_t)(A_{t-1} + B_{t-1}) \dots (A_2 + B_2)z_1 \quad (8)$$

Here A_t carries the gradients that arise due to the linear temporal path. The matrix A_t is bounded given that the number of time steps is finite and the initialization of c_0 is finite. On the other hand, the matrix B_t is obtained through the repeated multiplication of weights. If the weights are large, the matrix B_t can also become very large in magnitude. In this case, the magnitude of A_t will be very small as compared to B_t , hence the gradients through the linear temporal path (which carry information of long term dependency) will have no effect on the weights.

2.3 ALGORITHM:H-DETACH

The above problem can be mitigated by multiplying B_t by a positive number between 0 and 1. The authors provide a simple stochastic algorithm to achieve this which is presented as Algorithm 1.

Algorithm 1 h-detach algorithm

```

1: Input :  $p, h_0, c_0, [x_0, x_1, \dots, x_T]$ 
2:  $l \leftarrow 0$ 
3: for  $1 \leq t \leq T$  do
4:   if  $\text{bernoulli}(p) == 1$  then
5:      $\tilde{h}_{t-1} \leftarrow \text{stop\_gradient}(h_{t-1})$ 
6:   else
7:      $\tilde{h}_{t-1} \leftarrow h_{t-1}$ 
8:    $h_t, c_t \leftarrow \text{LSTM}(x_t, \tilde{h}_{t-1}, c_{t-1})$ 
9:    $l_t \leftarrow L(\phi(h_t))$ 
10:   $l \leftarrow l + l_t$ 
11: return  $l$ 

```

3 IMPLEMENTATION DETAILS

The authors have made their code public¹. The codebase used for the reproducibility experiments is also publicly available². The implementation is in Pytorch³. The author provides instructions to run the code to perform the experiments. We could get the code running immediately without any changes. The stop-gradient in algorithm 1 can be implemented using the `detach()`⁴ method in pytorch.

The authors claim that LSTMs trained using h-detach converge in lesser number of epochs as compared to normal LSTMs. We have been able to confirm this claim in our experiments. Even though the convergence is faster in terms of the number of epochs, the implementation (provided by the authors) is slower as compared to the vanilla LSTM implementation provided by pytorch⁵. This is because, in the implementation, each input is fed sequentially which makes it slower compared to feeding inputs at all time steps at once (as is usually done for efficiency). On communication with the authors, it was confirmed that this was done to ensure the correctness of the h-detach algorithm. The authors have themselves implemented the LSTM cell in pytorch. In our experiments, we observed that implementing the LSTM cell in CUDA code. The code can be found in the reproducibility repo⁶. Implementing the LSTM cell in CUDA code can result in upto 2x speed-up. This CUDA code can be added as an extension in pytorch and used.

4 EXPERIMENTS

4.1 COPYING TASK

This task requires the recurrent network to memorize the network inputs provided at the first few time steps and output them in the same order after a large time delay. Thus the only way to solve this task is for the network to capture the long term dependency between inputs and targets which requires gradient components carrying this information to flow through many time steps.

The authors define 10 token $\{a_i\}_{i=0}^{i=9}$. The input to the LSTM is a sequence of length $T + 20$ formed using one of the ten tokens at each time step. Input for the first 10 time steps are sampled i.i.d.

¹<https://github.com/bhargav104/h-detach>

²<https://github.com/dido1998/h-detach>

³<https://pytorch.org/>

⁴<https://pytorch.org/docs/stable/autograd.html#torch.Tensor.detach>

⁵<https://pytorch.org/docs/stable/nn.html#torch.nn.LSTM>

⁶<https://github.com/dido1998/h-detach/tree/master/h-detach-cuda-extension>

(uniformly) from $\{a_i\}_{i=0}^{i=7}$. The next T1 entries are set to a_8 , which constitutes a delay. The next single entry is a_9 , which represents a delimiter, which should indicate to the algorithm that it is now required to reproduce the initial 10 input tokens as output. The remaining 10 input entries are set to a_8 . The target sequence consists of $T + 10$ repeated entries of a_8 , followed by the first 10 entries of the input sequence in exactly the same order. This copying task is similar to the one described in Arjovsky et al. (2015).

We run each experiment for 200-300 epochs. The main point the authors wanted to prove is that using h-detach results in more stable training and faster convergence. We observe that these properties can be observed in the first 200-300 epochs. The authors have used the probability values of 0.25 and 0.50 for h-detach. The original paper does not mention how they arrived at these values. However, on the openreview site⁷, the authors mention that they tried the probability values 0.1, 0.25, 0.4, 0.5, 0.75 and 0.9. They found that the values between 0.25 and 0.50 work the best.

We present our experiments in Figure 1. The top row presents results for the copying task with $T=100$ and with the probabilities for h-detach set as 0, 0.25 and 0.5. The bottom row presents results for the copying task with $T=300$ and the same probabilities. As the authors mention, we see the main advantage of using h-detach in the task with $T=300$. We see that, the models trained with h-detach result in higher accuracies and faster convergence. Here, we also see that $p=0.50$ works better than $p=0.25$. For the task with $T=100$, we observe that both methods achieve high accuracy of close to 100% but h-detach results in faster convergence. The hyperparameters are exactly the same as provided by the authors in the source code.

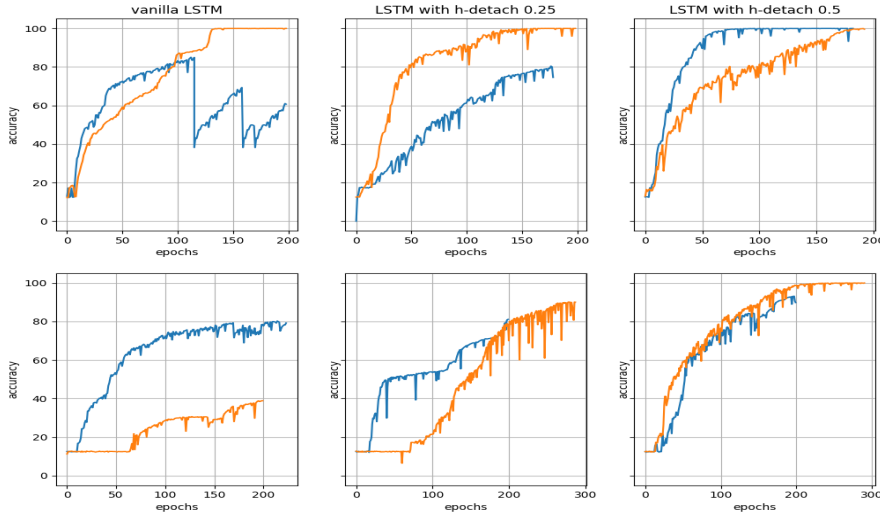


Figure 1: Validation accuracy curves for vanilla LSTM (left), h-detach with probability 0.25 (centre) and h-detach with probability 0.5(right) for the copying task. The top row shows the curves for $T=100$ and the bottom row shows the curves for $T=300$. For $T=300$, we can see that only the LSTMs with h-detach are able to cross the accuracy of 80%. It can be seen that training with h-detach is more stable and leads to faster convergence. Each experiment was tried with two different seeds.

4.2 SEQUENTIAL MNIST

This is the sequential version of the MNIST classification task. Each pixel is fed one at a time. The authors consider two versions of this task: one is which the pixels are fed in order (left to right and top to bottom) and the other is in which pixels are fed in random but fixed order. The second one is

⁷<https://openreview.net/forum?id=ryf7ioRqFX>

called permuted mnist. We tested with two different learning rates (0.0005 and 0.001). The results can be seen in Figure 2

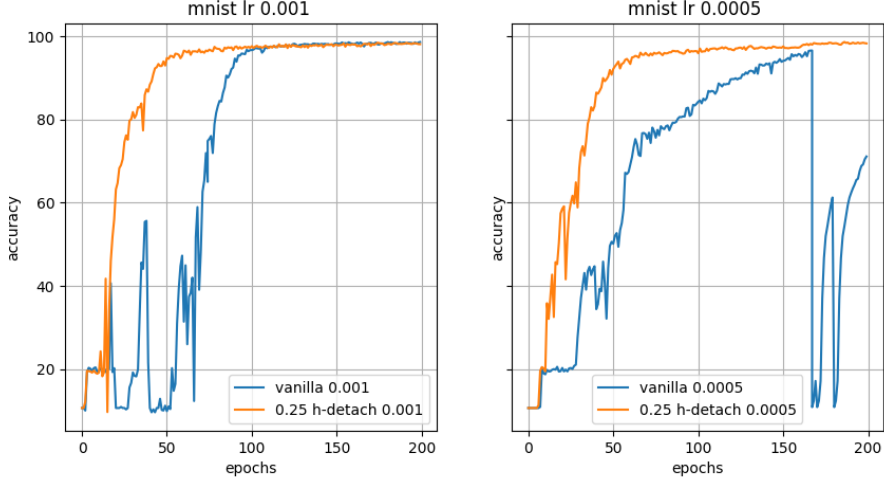


Figure 2: Validation accuracy curves for the sequential mnist task. Each plot shows the accuracy curves with and without h-detach for different learning rates. H-detach is more robust to different learning rates and converges faster.

The optimizer used was ADAM. The batch-size was set to 100. The number of examples in the training set was set to 50000, 10000 for validation and 10000 for testing. The number of epochs was set to 200. On the sequential MNIST task, vanilla LSTM and training with h-detach give an accuracy of 98.60%(authors report 98.60%) and 98.62%(authors report 98.80%) respectively. The numbers are not very different but we observed that training with h-detach gives faster convergence.

For the permuted mnist (pmnist) task we trained with h-detach and obtained the maximum accuracy of 91.96%(authors report 92.92%). The validation accuracy curve for the permuted mnist task can be found in Figure 3

4.3 ABLATION STUDIES

The authors first study the effect of removing gradient clipping from the sequential mnist task. This would be insightful as it would confirm that stochastically blocking gradients through the hidden state h_t prevents increase in gradient magnitude. As the authors say, we observe that removing gradient clipping makes the training extremely unstable in the case of a vanilla LSTM as compared to LSTM with h-detach. This is demonstrated in Figure 4.

The authors also conduct experiments to prove that the cell state carries information about long term dependencies. This is done by stochastically blocking gradients through the cell state instead of the hidden state. This new algorithm is tested on the copying task with $T=100$ and the sequential mnist task. The results are demonstrated in Figure 5. For the copying task, the accuracy does not cross 40% even after 200 epochs. For the sequential mnist task, the accuracy crosses 80% but the training is extremely unstable and the h-detach version completely out performs the c-detach version.

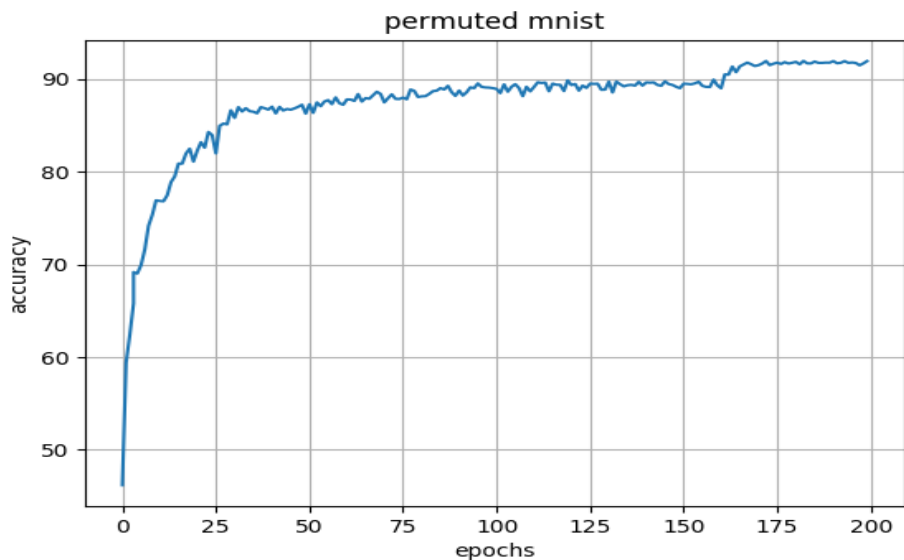


Figure 3: validation accuracy curve for permuted mnist with h-detach probability 0.25 and learning rate 0.001.

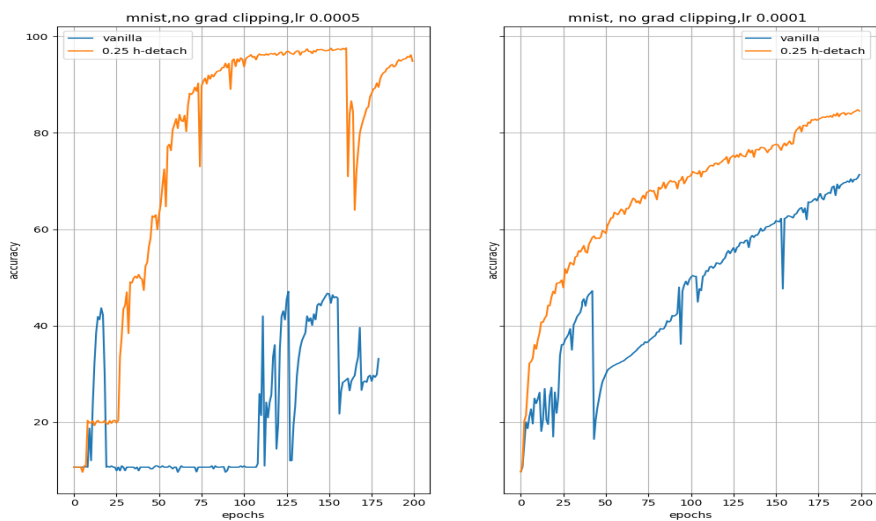


Figure 4: The effect of removing gradient clipping while training for vanilla LSTM vs LSTM with h-detach.

5 ENVIRONMENT DETAILS

The experiments were conducted on a NVIDIA GTX 1060 and on google colab⁸ (tesla K80).

⁸<https://colab.research.google.com/>

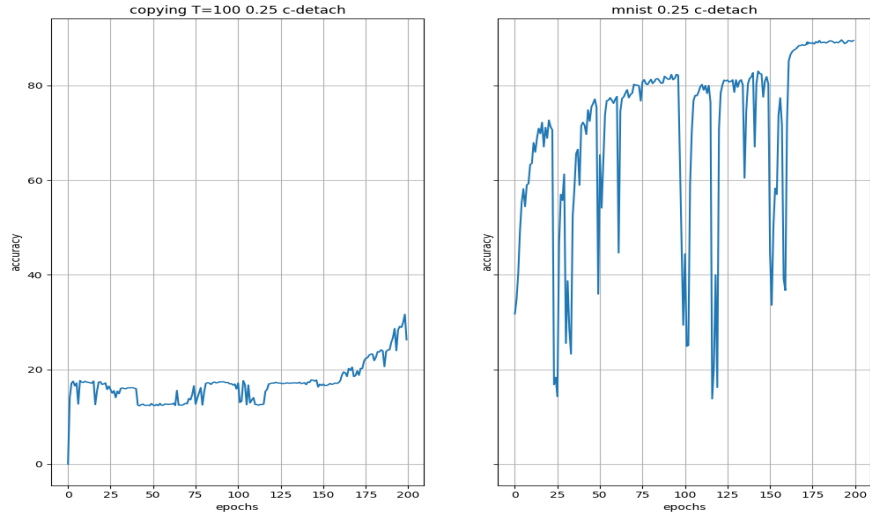


Figure 5: Validation accuracy curves for the copying task $T=100$ (left) and pixel mnist task(right) such that the gradient is stochastically blocked through the cell state(c-detach). The highest accuracy obtained for the pixel mnist task is 89.65%. As the authors say, c-detach performs worse than vanilla LSTM. This proves that the cell state path carries information about long term dependencies.

6 FUTURE WORK

Though the authors have thoroughly explored the advantages of the h-detach algorithm through their experiments, they have not evaluated their proposed algorithm on stacked LSTM networks and bidirectional RNNs with LSTM cells. It would be interesting to see if we get similar performance improvements on these architectures as in the single layer uni-directional case.

7 CONCLUSION

In this paper, we reproduce selected experiments of the original paper, Kanuparthi et al. (2019). Through communication with the authors, we understood that the main aim of this work was to prove that h-detach achieves better training stability, convergence speed, and it is robust to different learning rates and seeds. Though some of the numbers reported do not match exactly with the original paper, we have successfully proved that h-detach results in more stable and faster convergence.

REFERENCES

- Martín Arjovsky, Amar Shah, and Yoshua Bengio. Unitary evolution recurrent neural networks. *CoRR*, abs/1511.06464, 2015. URL <http://arxiv.org/abs/1511.06464>.
- Y Bengio, P Simard, and P Frasconi. Learning long-term dependencies with gradient descent is difficult. *IEEE transactions on neural networks / a publication of the IEEE Neural Networks Council*, 5(2):157–66, 1994.
- Junyoung Chung, Çağlar Gülçehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014. URL <http://arxiv.org/abs/1412.3555>.

- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Comput.*, 9(8):1735–1780, November 1997. ISSN 0899-7667. doi: 10.1162/neco.1997.9.8.1735. URL <http://dx.doi.org/10.1162/neco.1997.9.8.1735>.
- Bhargav Kanuparthi, Devansh Arpit, Giancarlo Kerg, Nan Rosemary Ke, Ioannis Mitliagkas, and Yoshua Bengio. h-detach: Modifying the LSTM gradient towards better optimization. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=ryf7ioRqFX>.
- Quoc V. Le, Navdeep Jaitly, and Geoffrey E. Hinton. A simple way to initialize recurrent networks of rectified linear units. *CoRR*, abs/1504.00941, 2015. URL <http://arxiv.org/abs/1504.00941>.
- Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio. Understanding the exploding gradient problem. *CoRR*, abs/1211.5063, 2012. URL <http://arxiv.org/abs/1211.5063>.
- David E. Rumelhart, Geoffrey E. Hinton, and Ronald J. Williams. Neurocomputing: Foundations of research. chapter Learning Representations by Back-propagating Errors, pp. 696–699. MIT Press, Cambridge, MA, USA, 1988. ISBN 0-262-01097-6. URL <http://dl.acm.org/citation.cfm?id=65669.104451>.
- Gail Weiss, Yoav Goldberg, and Eran Yahav. On the practical computational power of finite precision rnns for language recognition. *CoRR*, abs/1805.04908, 2018. URL <http://arxiv.org/abs/1805.04908>.