

---

# Recurrent neural network regularization

---

**Wojciech Zaremba**

Google & New York University

**Ilya Sutskever**

Google

**Oriol Vinyals**

Google

WOJ.ZAREMBA@GMAIL.COM

ILYASU@GOOGLE.COM

VINYALS@GOOGLE.COM

## Abstract

We present a simple regularization technique of recurrent neural networks (RNNs) with long short term memory (LSTM) units. This technique is based on dropout and gives tremendous performance boost. We show that it is beneficial in variety sequence modelling problems like language modeling, speech recognition, and machine translation.

## 1. Introduction

RNNs yields the state-of-the-art performance on many sequence modelling tasks like language modelling, and speech recognition. Moreover, recent results in machine translation (Cho et al., 2014) shows their potential use in this field as well. However, up today there is no good techniques to regularize them. As a result, people tend to use too small models, due to over-fitting by a large one. So far, various regularization attempts applied to RNNs give just a small improvements (Graves, 2013). This work presents how to augment LSTMs with dropout. Resulting models give a tremendous improvement over the baseline in aforementioned tasks.

## 2. Related work

Tools, and tasks considered in this paper were broadly used before. Moreover, there is a great interest in understanding (Warde-Farley et al., 2013; Srivastava, 2013), and extending dropout (Wang & Manning, 2013; Wan et al., 2013) for feed forward networks. There was a little research so far in using dropout in recurrent network architectures. The only known for us paper in this area is (Bayer et al., 2013). It focuses on using “dropout marginalization” (from

(Wang & Manning, 2013)) instead of dropout, which has a smaller variance than dropout. They claim that dropout cannot be successfully used in RNNs due to large variance, which might cause divergence. We show that application of dropout on proper connections gives viable, high-accuracy network.

In terms of tools, there has been extensive work on using RNNs for sequence modelling (Mikolov, 2012; Sutskever, 2013). Moreover, there were considered various twists into common architecture, which potentially can capture long term dependencies (Hochreiter & Schmidhuber, 1997; Graves et al., 2009; Cho et al., 2014). This work choses LSTM architecture, however it is possible that conclusions presented here would extend to other models.

We focus here on tasks, which were previously considered in RNN literature like (1) language modelling, (2) speech recognition, and (3) machine translation. Language modelling is a classical task addressed by RNNs (Pascanu et al., 2013; Mikolov et al., 2010; 2011). Using RNNs for speech recognition was previously considered by (Robinson et al., 1996; Graves et al., 2013). Finally, machine translation is a novel task for which RNNs are used. Usually, they are used for language modelling, re-ranking, or phrase modelling (Cho et al., 2014).

## 3. Regularized RNN with LSTM cells

We describe in this section how LSTM works 3.1. Next, we show how to regularize them 3.2, and we give some intuitions why it works.

We use upper indexing for layer number, lower indexing for a time step. Our states are  $n$  dimensional.  $h_k^l \in \mathbb{R}^n$  is a hidden state in layer  $l$  in step  $k$ . Moreover, for simplicity we denote by  $T : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^n$  a linear transform with bias ( $Wx + b$ ).  $\odot$  is a element-wise multiplication.  $h_k^0$  is a input word-vector, and  $h_k^L$  is used to predict  $y_k$  ( $L$  is a number of layers).

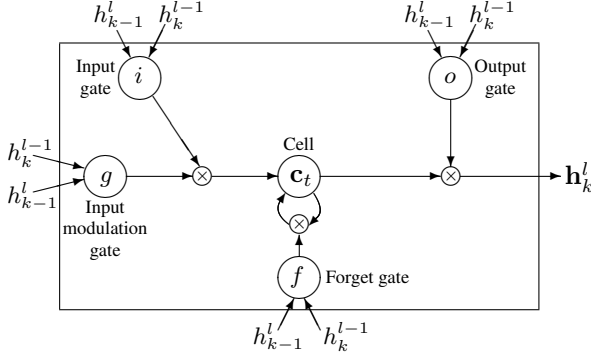


Figure 1. Graphical representation of LSTM memory cells used in this paper (there are minor differences in compare to (Graves, 2013)).

### 3.1. Long-short term memory units

Dynamics of RNNs can be described in terms of transitions from previous hidden states to the hidden states in the next step. That is the function

$$\text{RNN} : h_k^{l-1}, h_{k-1}^l \rightarrow h_k^l$$

In case of classical RNNs this function is described as:

$$h_k^l = f(Th_k^{l-1} + Th_{k-1}^l), \text{ where } f \in \{\sigma, \tanh\}$$

LSTM has more complicated dynamics in order to be able to “memorize” information for extended number of time steps. Memory is stored in a cell ( $c_k^l \in \mathbb{R}^n$ ) units. There has been proposed various LSTM architectures, which differ in terms of applied activations functions, connections etc. However, the common goal is to have a “memory cell”, which might store value for extended time. Network can decide to override this value, retrieve it, or keep it for a next time step. LSTM used in our experiments is described with following mapping.

$$\begin{aligned} \text{LSTM} : h_k^{l-1}, h_{k-1}^l, c_{k-1}^l &\rightarrow h_k^l, c_k^l \\ \begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} &= \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} T \begin{pmatrix} h_k^{l-1} \\ h_{k-1}^l \end{pmatrix} \\ c_k^l &= f \odot c_{k-1}^l + i \odot g \\ h_k^l &= o \odot \tanh(c_k^l) \end{aligned}$$

Where  $\sigma, \tanh$  are applied element-wise. Diagram 1 presents graphically LSTM equations.

### 3.2. Regularization with dropout

Contribution of this paper is finding out that carefully placed dropout improves generalization of LSTMs tremen-

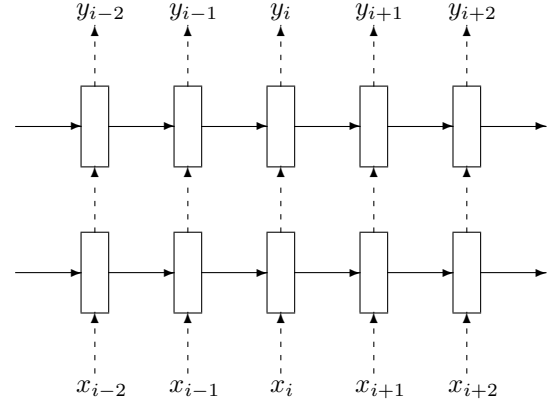


Figure 2. Regularized multilayer RNN. Dashed arrows indicate connections with applied dropout, while solid lines indicate connections where dropout is not applied.

dously. It is enough to place it on non-recurrent connections 2. Following equation describes it in more mathematical way, where  $D$  is a dropout operator:

$$\begin{aligned} \begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} &= \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} T \left( \mathbf{D} \begin{pmatrix} h_k^{l-1} \\ h_{k-1}^l \end{pmatrix} \right) \\ c_k^l &= f \odot c_{k-1}^l + i \odot g \\ h_k^l &= o \odot \tanh(c_k^l) \end{aligned}$$

Dropout removes part of information, and units have to become more robust while performing intermediate input-intermediate output mapping. Simultaneously we don’t want to erase entire information. Especially we would like to facilitate memory about events that happened long time ago. Figure 3 shows a flow of an exemplary information from the event  $x_{i-2}$  to the prediction in the step  $i+2$ . We can see that information is influenced with dropout only  $L+1$  times, and it is independent of how far in past event occurred. All the previous regularization techniques were constraining recurrent connections, which effectively exponentially fast “blurs” information about the past events.

## 4. Experiments

We present here results in various domains (1) language modeling 4.1, (2) speech recognition 4.2, and (3) machine translation 4.3.

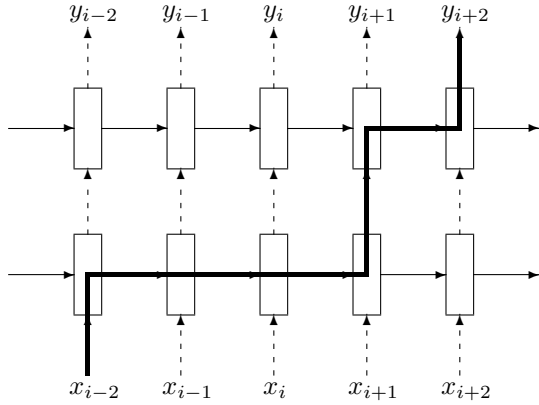


Figure 3. Thick line indicates an exemplary information flow in RNN. Information flow line is crossed  $L + 1$  times, where  $L$  is depth of network.

#### 4.1. Language modeling

We have conducted word-level prediction experiments on Penn tree bank (PTB) dataset. This dataset consists of 929k training examples, 73k validation examples, and 82k test examples. It has 10k words in vocabulary.

Our model is a two layer LSTM network, with 650 units initialized uniformly in  $[-0.05, 0.05]$ . We apply 50% of dropout on non-recurrent connections. We train for 39 epochs, starting with learning rate 1, and after 6 epochs we decrease it by 1.2 in every epoch. We unroll RNN for 35 steps, and clip gradients at 5. We set mini-batch to 20. Table 1 compares various models with our models.

#### 4.2. Speech recognition

#### 4.3. Machine translation

### 5. Discussion

We present amazing performance boosts with methods, which we only intuitively understand. We think, that it crucial to derive our models analytically, rather than based only on intuition. However, so far this problems seem to be very difficult to tackle.

### References

Bayer, Justin, Osendorfer, Christian, Chen, Nutan, Urban, Sebastian, and van der Smagt, Patrick. On fast dropout and its applicability to recurrent networks. *arXiv preprint arXiv:1311.0701*, 2013.

Model	Validation set	Test set
A single model		
(Pascanu et al., 2013)		107.5
(Cheng et al.)		100.0
Non-regularized LSTM	120.7	114.5
Regularized LSTM	86.2	<b>82.7</b>
Model averaging		
(Mikolov, 2012)	83.5	89.4
(Cheng et al.)		80.6
2 non-regularized LSTMs	100.4	96.1
5 non-regularized LSTMs	87.9	84.1
10 non-regularized LSTMs	83.5	80.0
2 regularized LSTMs	80.6	77.0
5 regularized LSTMs	76.7	73.3
10 regularized LSTMs	75.2	<b>72.0</b>

Table 1. Word-level perplexity on Penn-tree-bank dataset.

Cheng, Wei-Chen, Kok, Stanley, Pham, Hoai Vu, Chieu, Hai Leong, and Chai, Kian Ming A. Language modeling with sum-product networks.

Cho, Kyunghyun, van Merriënboer, Bart, Gulcehre, Caglar, Bougares, Fethi, Schwenk, Holger, and Bengio, Yoshua. Learning phrase representations using rnn encoder-decoder for statistical machine translation. *arXiv preprint arXiv:1406.1078*, 2014.

Graves, Alex. Generating sequences with recurrent neural networks. *arXiv preprint arXiv:1308.0850*, 2013.

Graves, Alex, Liwicki, Marcus, Fernández, Santiago, Bertolami, Roman, Bunke, Horst, and Schmidhuber, Jürgen. A novel connectionist system for unconstrained handwriting recognition. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 31(5):855–868, 2009.

Graves, Alex, Mohamed, Abdel-rahman, and Hinton, Geoffrey. Speech recognition with deep recurrent neural networks. In *Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on*, pp. 6645–6649. IEEE, 2013.

Hochreiter, Sepp and Schmidhuber, Jürgen. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.

Mikolov, Tomáš. *Statistical language models based on neural networks*. PhD thesis, Ph. D. thesis, Brno University of Technology, 2012.

Mikolov, Tomas, Karafiát, Martin, Burget, Lukas, Cernocký, Jan, and Khudanpur, Sanjeev. Recurrent neural network based language model. In *INTERSPEECH*, pp. 1045–1048, 2010.

- Mikolov, Tomas, Deoras, Anoop, Povey, Daniel, Burget, Lukas, and Cernocky, Jan. Strategies for training large scale neural network language models. In *Automatic Speech Recognition and Understanding (ASRU), 2011 IEEE Workshop on*, pp. 196–201. IEEE, 2011.
- Pascanu, Razvan, Gulcehre, Caglar, Cho, Kyunghyun, and Bengio, Yoshua. How to construct deep recurrent neural networks. *arXiv preprint arXiv:1312.6026*, 2013.
- Robinson, Tony, Hochberg, Mike, and Renals, Steve. The use of recurrent neural networks in continuous speech recognition. In *Automatic speech and speaker recognition*, pp. 233–258. Springer, 1996.
- Srivastava, Nitish. *Improving neural networks with dropout*. PhD thesis, University of Toronto, 2013.
- Sutskever, Ilya. *Training Recurrent Neural Networks*. PhD thesis, University of Toronto, 2013.
- Wan, Li, Zeiler, Matthew, Zhang, Sixin, Cun, Yann L, and Fergus, Rob. Regularization of neural networks using dropconnect. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pp. 1058–1066, 2013.
- Wang, Sida and Manning, Christopher. Fast dropout training. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pp. 118–126, 2013.
- Warde-Farley, David, Goodfellow, Ian J, Courville, Aaron, and Bengio, Yoshua. An empirical analysis of dropout in piecewise linear networks. *arXiv preprint arXiv:1312.6197*, 2013.