

# Dropout in Recurrent Neural Networks

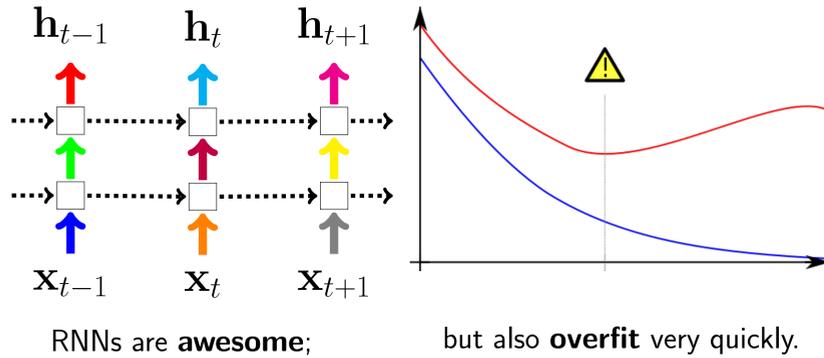
## A Theoretically Grounded Dropout Variant in RNNs using Variational Inference

Yarin Gal yg279@cam.ac.uk



UNIVERSITY OF CAMBRIDGE

### RNNs overfit quickly



This means...

- We can't use **large** models
- We have to use **early stopping**
- We can't use **small data**
- We have to **waste data** for validation sets

### Existing dropout in RNNs

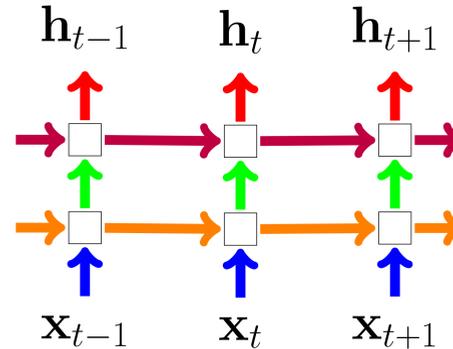
Let's use **dropout** then. But lots of research has claimed that that's a **bad idea**:

- **Pachitariu & Sahani**, 2013
  - noise added in the recurrent connections of an RNN leads to model **instabilities**
- **Bayer et al.**, 2013
  - with dropout, the RNN's **dynamics change** dramatically
- **Pham et al.**, 2014
  - dropout in recurrent layers **disrupts** the RNN's ability to model sequences
- **Zaremba et al.**, 2014
  - applying dropout to the **non-recurrent connections alone** results in improved performance
- **Bluche et al.**, 2015
  - **exploratory analysis** of the performance of dropout before, inside, and after the RNN's
- **Moon et al.**, 2015
  - Drop elements in the **LSTM's cell** using the same mask at every time step.

Many settled on using dropout for **inputs and outputs alone**.

### VI based dropout in RNNs

Uses the **same dropout mask** at each time step, including recurrent layers, and **drops word types** at random throughout the sentence:



### Why does it make sense?

- **Input**: sequence of vectors  $\mathbf{x} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$  with  $T$  time steps
- Let  $\omega = \{\text{all model weight matrices}\}$  and put prior  $p(\omega)$  (e.g. standard Gaussian)
- Define  $\mathbf{h}_t = \mathbf{f}_h^\omega(\mathbf{x}_t, \mathbf{h}_{t-1})$ 
  - single **recurrent unit** transition. E.g.  $\tanh(W\mathbf{x}_t + U\mathbf{h}_{t-1} + b)$  (similarly for **LSTM, GRU**)
- Set  $\mathbf{f}^\omega(\mathbf{x}) = \mathbf{f}_y^\omega(\mathbf{h}_T)$ 
  - model **output** (e.g. affine transformation of last state, or function of all states)
- Lastly, define  $p(\mathbf{y}|\mathbf{f}^\omega(\mathbf{x}))$ 
  - model **likelihood** on random function output  $\mathbf{f}^\omega(\mathbf{x})$ . E.g.  $\mathcal{N}(\mathbf{y}; \mathbf{f}^\omega(\mathbf{x}), \sigma^2)$

- **Variational interpretation of dropout** [Gal and Ghahramani, 2015]: Dropout objective minimises

$$\begin{aligned} \text{KL}(q(\omega)||p(\omega|\mathbf{X}, \mathbf{Y})) &\propto - \int q(\omega) \log p(\mathbf{Y}|\mathbf{X}, \omega) d\omega + \text{KL}(q(\omega)||p(\omega)) \\ &= - \sum_{i=1}^N \int q(\omega) \log p(\mathbf{y}_i|\mathbf{f}^\omega(\mathbf{x}_i)) d\omega + \text{KL}(q(\omega)||p(\omega)). \end{aligned}$$

with  $q(\omega)$  factorising over weight columns  $\mathbf{w}_{ik}$ , e.g.  $q(\mathbf{w}_{ik}) = p\delta_0 + (1-p)\delta_{m_{ik}}$ .

- But
 
$$\int q(\omega) \log p(\mathbf{y}|\mathbf{f}^\omega(\mathbf{x})) d\omega = \int q(\omega) \log p\left(\mathbf{y} \left| \mathbf{f}_y^\omega(\mathbf{f}_h^\omega(\mathbf{x}_T, \dots, \mathbf{f}_h^\omega(\mathbf{x}_1, \mathbf{h}_0) \dots)) \right.\right) d\omega,$$
- So using MC integration with  $\hat{\omega}_i \sim q(\omega)$ ,
 
$$\mathcal{L}_{VI} \approx - \sum_{i=1}^N \log p\left(\mathbf{y}_i \left| \mathbf{f}_y^{\hat{\omega}_i}(\mathbf{f}_h^{\hat{\omega}_i}(\mathbf{x}_{iT}, \dots, \mathbf{f}_h^{\hat{\omega}_i}(\mathbf{x}_{i1}, \mathbf{h}_0) \dots)) \right.\right) + \text{KL}(q_\theta(\omega) || p(\omega)).$$

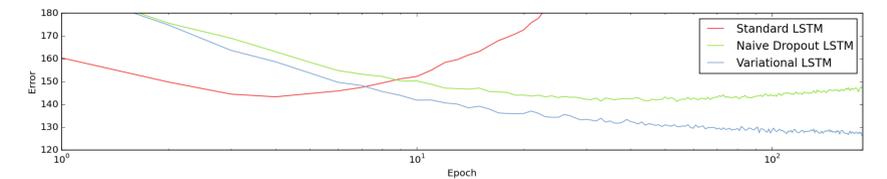
using random mask to set weight columns to zero (dropping units), repeating **the same mask at each time step** for **all weight matrices** (including embedding layer)

### Results

- **Penn Treebank** language modelling

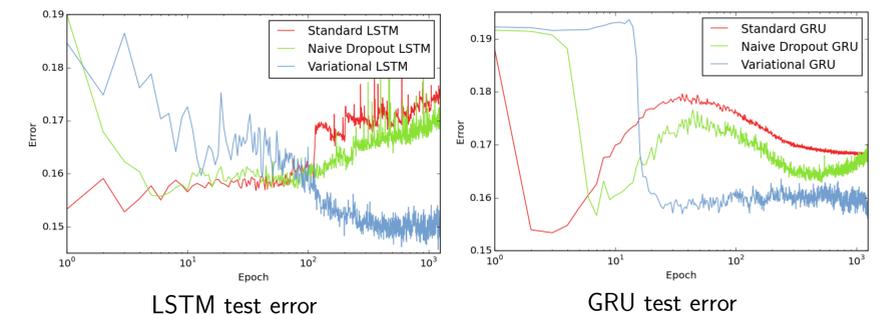
	Medium LSTM			Large LSTM		
	Validation	Test	WPS	Validation	Test	WPS
Non-regularized (early stopping)	121.1	121.7	5.5K	128.3	127.4	2.5K
Moon et al. [2015]	100.7	97.0	4.8K	122.9	118.7	3K
Moon et al. [2015] +emb dropout	88.9	86.5	4.8K	88.8	86.0	3K
Zaremba et al. [2014]	86.2	82.7	5.5K	82.2	78.4	2.5K
Variational (tied weights)	81.8 ± 0.2	79.7 ± 0.1	4.7K	77.3 ± 0.2	75.0 ± 0.1	2.4K
Variational (tied weights, MC)	–	79.0 ± 0.1	–	–	74.1 ± 0.0	–
Variational (untied weights)	81.9 ± 0.2	79.7 ± 0.1	2.7K	77.9 ± 0.3	75.2 ± 0.2	1.6K
Variational (untied weights, MC)	–	<b>78.6 ± 0.1</b>	–	–	<b>73.4 ± 0.0</b>	–

Single model perplexity (on test and validation sets). Two LSTM sizes are compared using Zaremba, Sutskever, and Vinyals [2014]'s setup.

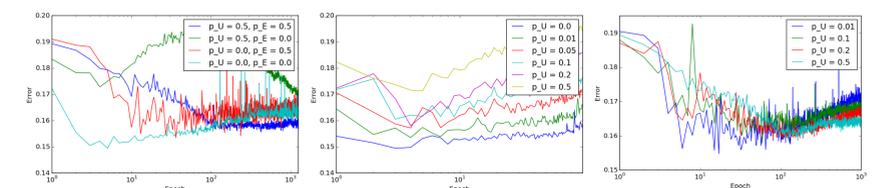


Validation perplexity (medium model) with dropout regularisation alone

- **Sentiment analysis** (raw Cornell film reviews corpus, Pang and Lee [2005])



Different dropout probabilities used with the recurrent layer ( $p_U$ ) and embedding layer ( $p_E$ ):



$p_E = 0, 0.5 \times p_U = 0, 0.5$      $p_U = 0, \dots, 0.5 \ \& \ p_E = 0$      $p_U = 0, \dots, 0.5 \ \& \ p_E = 0.5$